

## OPTIMAL ESTIMATORS UNDER REGRESSION MODEL

By V. R. PADMAWAR

Indian Statistical Institute

**SUMMARY.** In view of the nonexistence of an optimal design unbiased strategy under regression model, certain optimal strategies are obtained under different restricted set-ups.

### 1. INTRODUCTION

Consider a finite population  $U = \{1, 2, \dots, N\}$  of size  $N$ . Let  $y$  be a study variate taking values  $y_i$  on units  $i$ ,  $1 \leq i \leq N$ . Let  $x$  be an auxiliary variate, closely related to  $y$ , taking values  $x_i$  on units  $i$ ,  $1 \leq i \leq N$ . We assume that  $y_1, y_2, \dots, y_N$  are a realisation of the variables  $Y_1, Y_2, \dots, Y_N$ ; the joint distribution of which is unknown but specified by the first two moments as follows. If  $E_{\xi}$  and  $V_{\xi}$  denote expectation and variance, respectively, w.r.t. the model  $\xi$  that defines a class of distributions for  $Y_1, Y_2, \dots, Y_N$ ; then

$$\begin{aligned} E_{\xi}(Y_i) &= \beta x_i & i &= 1, 2, \dots, N \\ V_{\xi}(Y_i) &= \sigma^2 x_i^2 & i &= 1, 2, \dots, N \quad \dots (1.1) \\ E_{\xi}(Y_i Y_j) &= \beta^2 x_i x_j & i \neq j &= 1, 2, \dots, N \end{aligned}$$

where  $\beta$  and  $\sigma^2 > 0$  are the model parameters and  $g \in [0, 2]$ . Model (1.1) is called regression model.

Our aim is to estimate the population mean  $\bar{Y} = \frac{1}{N} \sum_{i=1}^N y_i$ . Given a strategy  $(p, t)$  that consists of a design  $p$  and an estimator  $t(s, y)$ ,  $y = (y_1, y_2, \dots, y_N)$ , we use the following measure of uncertainty to measure its performance.

$$M(p, t) = E_p E_y (t - \bar{Y})^2 \quad \dots (1.2)$$

where  $E_p$  denotes expectation w.r.t. the design  $p$ .

A strategy  $(p, t)$  is said to be  $p$ -unbiased or design unbiased for estimating the population mean  $\bar{Y}$  if

$$E_p t(s, y) = E_p t = \sum_{s \in S} p(s) t(s, y) = \bar{Y} \quad \forall y$$

---

*AMS (1980) subject classification:* 62D05.

*Key words:* Regression model, Population mean, Measure of uncertainty, Optimal estimators.

where  $S$  is the collection of all samples.

A strategy  $(p, t)$  is said to be  $\xi$ -unbiased or modol unbiased for the population mean  $\bar{Y}$  if

$$E_t(t(s, y)) = E_t(\bar{Y}) \quad \forall s \text{ with } p(s) > 0$$

Now for a  $p$ -unbiased strategy  $(p, t)$

$$M(p, t) = E_p V_t(t) + E_p (E_t(t - \bar{Y}))^2 - V_t(\bar{Y}) \quad \dots (1.3)$$

and for a  $p$  as well as  $\xi$ -unbiased strategy  $(p, t)$

$$M(p, t) = E_p V_t(t) - V_t(\bar{Y}) \quad \dots (1.4)$$

Ramachandran (1982) obtained a strategy for estimating the population mean  $\bar{Y}$  that, if exists, is optimal in the class of all  $p$ -unbiased strategies of given average size in the sense of minimum  $M(p, t)$ . Her strategy consisted of Horvitz-Thompson estimator and a sampling design that satisfies the following conditions

$$(i) \sum_{i=1}^N p_i(s) = \frac{n x_i^{g/2}}{\sum_{j=1}^N x_j^{g/2}}, \quad 1 \leq i \leq N$$

and

$$(ii) \sum_{i=1}^N x_i^{1-g/2} = n \sum_{i=1}^N x_i / \sum_{i=1}^N x_i^{g/2} \quad \forall s \text{ with } p(s) > 0.$$

Condition (i) requires the design to be a  $\pi p x^{g/2}$  design (inclusion probabilities proportional to  $x^{g/2}$ ) of average size  $n$ .

As such it is difficult to construct a design satisfying the first condition. Further condition (ii) puts far too many constraints on the fixed set of constants  $x_1, x_2, \dots, x_N$  and  $g$  to be able to satisfy them. Therefore the existence of the optimal strategy due to Ramachandran is ruled out in most practical situations.

In view of this it is in order to either compare different known strategies or consider the problem of obtaining optimal strategies in different restricted set-ups. The first aspect is studied in Padmawar (1981). Here we deal with the second aspect. Padmawar (1984) proved the existence of certain optimal estimators under the regression model in the continuous survey sampling set up. We take this opportunity to see their implications in the finite set up.

## 2. OPTIMAL ESTIMATORS

Given a design  $p$ , in this section, we first establish the existence and uniqueness of best linear  $p$  as well as  $\xi$ -unbiased estimator in the sense of minimum  $M(p, t)$ . Next, assuming the ratio  $\sigma^2/\beta^2 = r$  (say), in the model (1.1), to be known we prove the existence of a best linear  $p$ -unbiased estimator again in the sense of minimum  $M(p, t)$ .

A linear estimator is of the type  $t(s, y) = \sum_{i \in S} b(s, i)y_i$ . The condition of  $p$ -unbiasedness is equivalent to

$$\sum_{i \in I} b(s, i)p(s) = \frac{1}{N} \quad \forall i = 1, 2, \dots, N \quad \dots (2.1)$$

and the condition of  $\xi$ -unbiasedness is equivalent to

$$\sum_{i \in S} b(s, i)x_i = \bar{X} \quad \forall s \text{ with } p(s) > 0 \quad \dots (2.2)$$

where 
$$X = N\bar{X} = \sum_{i=1}^N x_i.$$

For a given design  $p$ , to tackle the first problem, we minimise  $M(p, t)$  subject to the conditions (2.1) and (2.2). In view of (1.4) the problem reduces to the following minimisation problem.

$$\begin{aligned} \text{Minimise} \quad & \sum_{i=1}^N q_i \sum_{s \in I} a^2(s, i)p(s) \\ \text{subject to} \quad & \sum_{s \in I} a(s, i)p(s) = 1 \quad \forall i = 1, 2, \dots, N \quad \dots (2.3) \\ \text{and} \quad & \sum_{s \in S} a(s, i)p_i = 1 \quad \forall s \text{ with } p(s) > 0 \end{aligned}$$

where  $a(s, i) = Nb(s, i)$ ;  $x_i = Np_i$  and  $N^2q_i = x_i^2$ .

We employ Lagrangian multipliers' technique to solve (2.3). Let

$$\begin{aligned} Q = \sum_{i=1}^N q_i \sum_{s \in I} a^2(s, i)p(s) - 2 \sum_{i=1}^N \lambda_i \left[ \sum_{s \in I} a(s, i)p(s) - 1 \right] \\ - 2 \sum_{s \in S} p(s)x_s \left[ \sum_{i \in I} a(s, i)p_i - 1 \right] \end{aligned}$$

where  $\lambda_i$ ,  $1 \leq i \leq N$  and  $\alpha_s$ ;  $s \in S$  and  $p(s) > 0$  are the Lagrangian multipliers. Solving (2.3) is equivalent to minimising  $Q$  unconditionally. Differ-

entiating and simplifying the problem boils down to solving the following system of equations

$$\begin{aligned} g_{i1} a(s, i) &= \lambda_i + \alpha_s & i \in S, p(s) > 0 \\ \sum_{s \in S} a(s, i) p(s) &= 1 & 1 \leq i \leq N & \dots (2.4) \\ \text{and} \quad \sum_{i \in S} a(s, i) p_i &= 1 & s \in S; p(s) > 0. \end{aligned}$$

After a little algebra one gets

$$a(s, i) = z_i + \eta_i \left( 1 - \sum_{j \neq i} z_j p_j \right) m(s) \quad i \in S; s \in S, p(s) > 0 \quad \dots (2.5)$$

where  $\eta_i = p_i / m(s) = \frac{1}{\sum_{j \in S} \eta_j p_j}$  and  $\mathbf{z} = (z_1, z_2, \dots, z_N)$  is a solution to the system of equations

$$A \mathbf{z} = \mathbf{d} \quad \dots (2.6)$$

with

$$\begin{aligned} a_{ii} &= \pi_i - \eta_i p_i \sum_{s \in S} p(s) m(s) & i = 1, 2, \dots, N \\ a_{ij} &= -\eta_i p_j \sum_{s \in S, j} p(s) m(s) & i \neq j = 1, 2, \dots, N & \dots (2.7) \end{aligned}$$

and

$$d_i = 1 - \eta_i \sum_{s \in S} p(s) m(s) \quad i = 1, 2, \dots, N$$

where

$$\pi_i = \sum_{s \in S} p(s), \quad 1 \leq i \leq N.$$

Following Patel and Dharmadhikari (1977, 1978) it can be shown that the system (2.6) is consistent and the solutions (2.5) do not depend on a choice of  $\mathbf{z}$  satisfying (2.6). We thus have the following theorem

**Theorem 2.1:** *Given a design  $p$  there exists a unique best linear  $p$  as well as  $\xi$ -unbiased estimator in the sense of minimum  $M(p, t)$ .*

**Remark 2.1:** Theorem 2.1 is true even when we have a more general variance function  $\sigma^2 v(x)$ , instead of  $\sigma^2 x^2$ , in the model (1.1). In that case we simply set  $N^2 g_i = v(x_i)$  in (2.3).

**Remark 2.2:** In the continuous set up (vide Padmawar (1984)) the existence of an analogous optimal estimator depends on the solvability of an integral equation whereas in the finite set up existence as well as uniqueness are always guaranteed. However, the construction of such an optimal estimator depends on the tractability of solutions to the system of equations (2.6). In what follows we consider two examples so as to get some idea of the above optimality result.

*Example 2.1:* Consider a fixed size design of size  $n$  for which  $\pi_i = nx_i/X$ ,  $1 \leq i \leq N$ .  $z_i = X/nx_i$ ,  $1 \leq i \leq N$ , is a solution to (2.6) for the  $i$ th entry in  $Az$  is then given by

$$\begin{aligned} & \frac{X}{nx_i} \left[ \pi_i - \eta_i p_i \sum_{s \neq i} p(s) m(s) \right] - \sum_{j \neq i} \frac{X}{nx_j} \eta_i p_j \sum_{s \neq i, j} p(s) m(s) \\ &= 1 - \frac{1}{n} \eta_i \sum_{s \neq i} p(s) m(s) - \frac{n-1}{n} \eta_i \sum_{s \neq i} p(s) m(s) \\ &= 1 - \eta_i \sum_{s \neq i} p(s) m(s) \\ &= d_i. \end{aligned}$$

Therefore the optimal estimator, as expected, is given by  $\frac{X}{nN} \sum_{i=1}^N \frac{y_i}{z_i}$ .

*Example 2.2:* Consider the Midzuno-Sen sampling design of fixed size  $n$  for which  $p(s)$  is given by  $w_s/M_1 W$  where  $w_s = \sum_{i \in s} w_i$ ,  $W = \sum_{i=1}^N w_i$  and  $M_1 = \binom{N-1}{n-1}$  with  $w_i = x_i^{1-\theta}$ ,  $1 \leq i \leq N$ . For this design the coefficients (2.7) simplify to

$$\begin{aligned} a_{ii} &= \pi_i - w_i/W & i &= 1, 2, \dots, N \\ a_{ij} &= -\frac{n-1}{N-1} x_i^{1-\theta} x_j/W & i \neq j &= 1, 2, \dots, N \\ d_i &= 1 - x_i^{1-\theta} X/W & i &= 1, 2, \dots, N \end{aligned}$$

where  $\pi_i = \frac{n-1}{N-1} + \frac{N-n}{N-1} w_i/W$ ,  $1 \leq i \leq N$ , and  $X = \sum_{i=1}^N x_i$

$$\begin{aligned} \text{Now } \sum_{j=1}^N a_{ij} &= \pi_i - w_i/W - \frac{n-1}{N-1} x_i^{1-\theta} (X - x_i)/W \\ &= \frac{n-1}{N-1} \{1 - x_i^{1-\theta} X/W\} \\ &= \frac{n-1}{N-1} d_i. \end{aligned}$$

Thus  $z_i = \frac{N-1}{n-1}$ ;  $1 \leq i \leq N$ ; is a solution to (2.6) and the optimal estimator is given by

$$\frac{1}{N} \frac{N-1}{n-1} \sum_{i=1}^N y_i + \frac{1}{N} \left\{ 1 - \frac{N-1}{n-1} \frac{1}{X} \sum_{i=1}^N x_i \right\} \sum_{i=1}^N x_i^{1-\theta} y_i \frac{X}{\sum_{i=1}^N x_i^{1-\theta}}$$

We have the following interesting cases.

*Case 1:* For  $g = 2$  in Example 2.2 the sampling design reduces to simple random sampling and the optimal estimator simplifies to

$$\frac{1}{N} \frac{N-1}{n-1} \sum_{i \in s} y_i + \frac{1}{N} \left\{ 1 - \frac{N-1}{n-1} \frac{1}{\bar{X}} \sum_{i \in s} x_i \right\} \frac{\bar{X}}{n} \sum_{i \in s} y_i / x_i$$

which can be rewritten as

$$\bar{r}\bar{X} + \frac{n(N-1)}{N(n-1)} (y - \bar{r}\bar{X})$$

where  $n\bar{x} = \sum_{i \in s} x_i$ ,  $n\bar{y} = \sum_{i \in s} y_i$ ,  $n\bar{r} = \sum_{i \in s} y_i/x_i$ ,  $N\bar{X} = X$ . Note that this optimal estimator is same as the Hartley-Ross (1954) ratio-type estimator for the population mean  $\bar{Y}$ .

*Case 2:* For  $g = 1$  in Example 2.2 the sampling design reduces to the usual Midzuno-Sen sampling scheme with  $p(s) = \frac{x_s}{M_1 \bar{X}}$ ,  $x_s = \sum_{i \in s} x_i$  and the optimal estimator reduces to

$$\frac{1}{N} \frac{N-1}{n-1} \sum_{i \in s} y_i + \frac{1}{N} \left\{ 1 - \frac{N-1}{n-1} \frac{1}{\bar{X}} \sum_{i \in s} x_i \right\} \frac{\bar{X}}{\sum_{i \in s} x_i} \sum_{i \in s} y_i$$

or

$$\frac{\sum_{i \in s} y_i}{\sum_{i \in s} x_i} \bar{X}$$

which is the standard ratio estimator for estimating the population mean.

We now assume that the ratio  $\sigma^2/\beta^2 = r$  in (1.1) is known. Given a design  $p$  we prove the existence of an optimal  $p$ -unbiased linear estimator in the sense of minimum  $J(p, l)$ .

In view of (1.3) our problem is to

$$\text{Minimise} \quad r \sum_{i=1}^N \alpha_i^2 \sum_{s \ni i} b^2(s, i) p(s) + \sum_{s \in S} p(s) \left[ \sum_{i \in s} U(s, i) x_i \right]^2 \quad \dots (2.8)$$

$$\text{subject to} \quad \sum_{s \ni i} b(s, i) p(s) = 1/N \quad \forall i = 1, 2, \dots, N$$

We again employ the Lagrangian multipliers' technique to solve the problem (2.8). Let

$$Q^* = r \sum_{i=1}^N \alpha_i^2 \sum_{s \ni i} b^2(s, i) p(s) + \sum_{s \in S} p(s) \left[ \sum_{i \in s} b(s, i) x_i \right]^2 \\ - 2 \sum_{i=1}^N \alpha_i \left[ \sum_{s \ni i} b(s, i) p(s) - \frac{1}{N} \right]$$

where  $\alpha_i$ ,  $1 < i < N$ , are the Lagrangian multipliers.

Solving (2.8) is equivalent to minimizing  $Q^*$  unconditionally. Differentiating and simplifying we get

$$b(\theta, i) = \frac{1}{r x_i^2} \left[ \alpha_i - x_i c(\theta) \sum_{j \neq i} \alpha_j x_j^{-\theta} \right] \quad \dots (2.9)$$

where  $r + \sum_{i \neq i} x_i^{2-\theta} = \frac{1}{c(\theta)}$  and  $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_N)$  is a solution to the following system of equations

$$B\alpha = h \quad \dots (2.10)$$

where

$$b_{ii} = \pi_i - x_i^{2-\theta} \sum_{j \neq i} p(\theta) c(\theta) \quad i = 1, 2, \dots, N$$

$$b_{ij} = -x_i x_j^{-1\theta} \sum_{s \neq i, j} p(\theta) c(\theta) \quad i \neq j = 1, 2, \dots, N$$

and  $\dots (2.11)$

$$N h_i = r x_i^2 \quad i = 1, 2, \dots, N$$

$$\pi_i = \sum_{j \neq i} p(\theta) \quad i = 1, 2, \dots, N$$

Note that the left hand side of the equation  $\det B = 0$  can be thought of as a polynomial in  $r$  of degree  $k$  (say). Therefore  $\det B = 0$  has at most  $k$  solutions  $r_1, r_2, \dots, r_k$  (say). For all other values of  $r$   $\det B$  is nonzero and hence  $B\alpha = h$  will have a unique solution. We therefore have the following theorem.

**Theorem 2.2 :** *If the ratio  $r$  in (1.1) is known then for a given design  $p$  there exists a unique best linear  $p$ -unbiased estimator for estimating the population mean  $\bar{Y}$  for all but a finite number of values of  $r$  in the sense of minimum  $M(p, t)$ .*

**Remark 2.3 :** Theorem 2.2 is true for a more general variance function  $\sigma^2(x)$ , instead of  $\sigma^2 x^2$ , in the model (1.1).

**Remark 2.4 :** In the continuous set up (vide Padmawar (1984)) the existence of an analogous optimal estimator depends on the solvability of an integral equation whereas in the finite set up existence as well as uniqueness are guaranteed for all but a few values of  $r$ . However, the construction of such an optimal estimator depends on the tractability of a solution to the system of equations (2.10). We consider an example so as to get an idea of the above optimality result.

*Example 2.3:* Consider a Midzuno-Son sampling design of fixed size  $n$  for which  $p(s)$  is given by  $\frac{w_s}{M_1 W}$  where  $w_s = \sum_{i \in s} w_i$ ,  $W = \sum_{i=1}^N w_i$  and  $M_1 = \binom{N-1}{n-1}$  with  $w_i = x_i^{2-\theta} + r/n$ ,  $1 \leq i \leq N$ . For this design the quantities in (2.11) are given by

$$\begin{aligned} b_{ii} &= \pi_i - x_i^{2-\theta} / W & i &= 1, 2, \dots, N \\ b_{ij} &= -(n-1)x_i x_j^{1-\theta} / (N-1)W & i &\neq j = 1, 2, \dots, N \end{aligned}$$

where  $\pi_i = (n-1)/(N-1) + (N-n)w_i/(N-1)W$ ,  $1 \leq i \leq N$ .

It is easy to see that  $\alpha_i = ax_i + bx_i^2$ ,  $1 \leq i \leq N$ , is a solution to (2.10).

where  $a = \frac{(n-1)b\bar{X}}{r(N-1)}$  and  $b = \frac{N-1}{N} \left/ \left( \frac{n-1}{r} + \frac{N-n}{nW} \right) \right.$  ... (2.12)

$$\begin{aligned} \text{Note that } \sum_{j=1}^N b_{ij} \alpha_j &= \pi_i \alpha_i - \alpha_i x_i^{2-\theta} / W - \frac{(n-1)x_i}{(N-1)W} \sum_{j \neq i}^N \alpha_j x_j^{1-\theta} \\ &= \left[ \frac{n-1}{N-1} + \frac{(N-n)r}{n(N-1)W} \right] (ax_i + bx_i^2) - \frac{(n-1)x_i}{(N-1)W} \left[ a \sum_{j=1}^N x_j^{2-\theta} + b \sum_{j=1}^N x_j \right] \\ &= \left[ \frac{n-1}{N-1} + \frac{(N-n)r}{n(N-1)W} \right] bx_i^2 + \frac{x_i}{(N-1)W} \left[ r(N-1)a - (n-1)b \sum_{j=1}^N x_j \right] \\ &= \frac{r}{N} x_i^2. \end{aligned}$$

Therefore (2.10) is consistent and the optimal estimator is of the form

$$\frac{b}{r} \sum_{i \in s} y_i + \left( a - \frac{b}{r} \sum_{i \in s} x_i \right) \frac{\sum_{i \in s} x_i^{1-\theta} y_i}{r + \sum_{j \in s} x_j^{2-\theta}}$$

where  $a$  and  $b$  are given by (2.12).

*Remark 2.5:* Note that for the sampling design in the above example the system (2.10) is consistent for all values of  $r$ . Therefore the optimal estimator exists for all values of  $r$ .

*Remark 2.6:* Instead of assuming the ratio  $r$  to be known if we assume a prior  $\rho(\sigma^2, \beta)$  for  $\sigma^2$  and  $\beta$  jointly with known  $E_\rho(\sigma^2)$  and  $E_\rho(\beta^2)$ , where  $E_\rho$  denotes expectation under the prior  $\rho$ , then  $E_\rho M(p, t)$  is a natural choice for an optimality criterion. In this set up we can compare any two strategies for estimating the population mean w.r.t. the optimality criterion  $E_\rho M(p, t)$ .



The strategy that consists of a design for which  $\pi_i x_i^{t/2}$ ,  $1 < i < N$ , and the corresponding Horvitz-Thompson estimator is relatively more and more efficient as the ratio  $r_p = E_p(\sigma^2)/E_p(\beta^2)$  increases and for large values of  $r_p$  it is likely to be the best strategy in the class of all  $p$ -unbiased strategies of given average size in the sense of minimum  $E_p M(p, t)$ . For a linear  $p$ -unbiased strategy  $E_p M(p, t)$  is given by

$$E_p(\sigma^2) \sum_{i=1}^N x_i^2 \left\{ \sum_{s, t} a^2(s, i) p(s) - \frac{1}{N^2} \right\} + E_p(\beta^2) V(p, t) |_{\mathbf{x}}$$

where  $V(p, t) |_{\mathbf{x}}$  denotes the variance of the strategy  $(p, t)$  at  $\mathbf{x} = (x_1, \dots, x_N)$ .

The aforesaid strategy minimises the first term in the class of all  $p$ -unbiased linear strategies of given average size  $n$ .

*Acknowledgements.* The author is thankful to the referee for his useful suggestions.

#### REFERENCES

- HARTLEY, H. C. and ROSS, A. (1954): Unbiased ratio estimates. *Nature*, 174, 270-271.
- PADMAWAR, V. R. (1981): A note on the comparison of certain sampling strategies. *J. R. Statist. Soc. B*, 43, 321-326.
- (1984): Two existence theorems in survey sampling of continuous populations. *Sankhyā*, B, 46, 217-227.
- PATEL, N. C. and DHARMADHIKARI, S. W. (1977): On linear invariant unbiased estimators in survey sampling. *Sankhyā*, C, 39, 21-27.
- (1978): Admissibility of Murthy's and Midzuno's estimators within the class of linear unbiased estimators of finite population totals. *Sankhyā*, C, 40, 21-28.
- RAMACHANDRAN, G. (1982): Horvitz-Thompson estimator and generalised  $\Pi$ s designs. *J.S.P.I.*, 7, 151-153.

*Paper received: May, 1985.*

*Revised: July, 1987.*