

A Regression Approach to Estimation in the Presence of Nonresponse

CARL ERIK SÄRNDAL¹

ABSTRACT

In the presence of unit nonresponse, two types of variables can sometimes be observed for units in the “intended” sample s , namely, (a) variables used to estimate the response mechanism (the response probabilities), (b) variables (here called co-variables) that explain the variable of interest, in the usual regression theory sense. This paper, based on Särndal and Swensson (1985 a, b), discusses nonresponse adjusted estimators with and without explicit involvement of co-variables. We conclude that the presence of strong co-variables in an estimator induces several favourable properties. Among other things, estimators making use of co-variables are considerably more resistant to nonresponse bias. We discuss the calculation of standard error and valid confidence intervals for estimators involving co-variables. The structure of the standard error is examined and discussed.

KEY WORDS: Response mechanism; Adjustment group method; Co-variate; Robustness.

1. INTRODUCTION

We consider a finite population $U = \{1, \dots, k, \dots, N\}$ from which a sample s of size n is drawn with a sampling design under which the k -th unit has the (strictly positive) probability π_k of being selected. The sampling weight associated with the k -th unit is thus π_k^{-1} . We may admit a complex sampling design, not necessarily self-weighting, for example, a three-stage design with stratified selection of primary units. The probability under the design of jointly including the units k and l is denoted π_{kl} ($\pi_{kl} > 0$ for all $k \neq l$, and π_{kk} is interpreted as equal to π_k).

Given s , a certain unit nonresponse is assumed to occur. The responding subset of s is denoted by r , its size by m . The variable of interest, y , is observed for $k \in r$ only. To counteract the biasing effects of the nonresponse, we assume for the purpose of this paper that the widely used adjustment group method is employed: the sample s is subdivided into H groups $s_1, \dots, s_h, \dots, s_H$ of respective sizes $n_1, \dots, n_h, \dots, n_H$. The response set r is correspondingly divided into the subsets $r_1, \dots, r_h, \dots, r_H$, of respective sizes $m_1, \dots, m_h, \dots, m_H$. The response rate in group h is denoted $f_h = m_h/n_h$. The method calls for attaching (in addition to the sampling weight) the “adjustment weight” f_h^{-1} to an observation coming from group h . (The sizes and the composition of the adjustment groups at the population level are here assumed unknown.) We have:

$$n = \sum_{h=1}^H n_h; \quad m = \sum_{h=1}^H m_h.$$

¹ Carl Erik Särndal, Department of Mathematics and Statistics, University of Montreal, Montreal, Quebec, Canada, H3C 3J7.

Let $t = \sum_U y_k$ be the unknown population total to be estimated. (If A is an arbitrary set of units, we shall systematically write $\sum_A y_k$ for $\sum_{k \in A} y_k$.) The usual adjustment class estimator of t then becomes

$$\hat{t} = \sum_{h=1}^H f_h^{-1} \sum_{r_h} \frac{y_k}{\pi_k}. \quad (1.1)$$

The adjustment group method is motivated theoretically by an assumption that units within the same group respond with the same (unknown) response probability. (More formally, this is expressed as Model A in Section 3 below.) The method clearly requires that group identity can be determined for each unit $k \in s$. The (categorical) variables that permit this grouping can thus be regarded as variables used for the estimation of an underlying response mechanism.

A different category of variable may be observable for each $k \in s$, namely, variables that explain y , in the ordinary regression theory sense. These variables will be termed co-variables. When incorporated in the estimator, such variables will not only reduce variance but also make the estimator more resistant to nonresponse bias. (They are not auxiliary variables in the usual sense of this term, since they are available not for the entire population U but only for the intended sample s .)

We shall thus keep a firm distinction in this paper between two types of variables observed for $k \in s$, those that are used to estimate the response mechanism, and those that explain the target variable y . Little (1983), in presenting a general framework for data with nonresponse, distinguishes several types of variables. One attempt to describe our situation in terms of Little's setup would be to say that the set of complete item variables in Little's terminology are, in our case, further subdivided into one subset of variables used to model the nonresponse mechanism, and another subset (the co-variables) serving as explanatory variables for the incomplete item variable y . Our approach to inference is that of "quasi-randomization" (Oh and Scheuren 1983), where "quasi" refers to the fact that the nonresponse selection phase must be modelled, whereas the sample selection phase is controlled by the sampler.

2. SOME SIMPLE NONRESPONSE ADJUSTED ESTIMATORS OF THE POPULATION TOTAL

A slight development of the often seen formula (1.1) leads to a (generally somewhat "better") alternative in which the sampling weights π_k^{-1} can be said to be more fully used:

$$\hat{t}_{\text{EXP}} = \left(\sum_s \frac{1}{\pi_k} \right) \frac{\sum_{h=1}^H f_h^{-1} \sum_{r_h} \frac{y_k}{\pi_k}}{\sum_{h=1}^H f_h^{-1} \sum_{r_h} \frac{1}{\pi_k}}.$$

The formula (which becomes identical to (1.1) for a self-weighting design) can be written as an expansion of the response set mean:

$$\hat{t}_{\text{EXP}} = \hat{N} \tilde{y}_r,$$

namely, if we let the expansion factor be $\hat{N} = \sum_s 1/\pi_k$, and

$$\tilde{y}_r = \frac{\sum_{h=1}^H f_h^{-1} \sum_{r_h} \frac{y_k}{\pi_k}}{\sum_{h=1}^H f_h^{-1} \sum_{r_h} \frac{1}{\pi_k}} \tag{2.1}$$

The symbol tilde will be used to indicate a properly weighted mean statistic. The ‘‘tilde mean’’ \tilde{y}_r , being a response set mean, is calculated by attaching to the k -th unit the multiplicative weight:

$$\text{sample weight} \times \text{non response adjustment weight} = \pi_k^{-1} f_h^{-1}$$

for each unit k in the h -th adjustment group.

The expansion estimator \hat{t}_{EXP} is appropriate for the nonresponse situation: it takes into account the sampling design and it makes an effort to adjust for nonresponse. However, \hat{t}_{EXP} can be improved upon if more information is at hand. Suppose that a single (and always positive) co-variate x is also observed for $k \in s$. In the image of the classical ratio estimator, we can then construct

$$\hat{t}_{\text{RA}} = \left(\sum_s \frac{x_k}{\pi_k} \right) \frac{\sum_{h=1}^H f_h^{-1} \sum_{r_h} \frac{y_k}{\pi_k}}{\sum_{h=1}^H f_h^{-1} \sum_{r_h} \frac{x_k}{\pi_k}} = \hat{N} \tilde{x}_s \frac{\tilde{y}_r}{\tilde{x}_r}$$

say, where the tilde mean \tilde{x}_r is formed according to (2.1) with x_k instead of y_k , and

$$\tilde{x}_s = \frac{\sum_s \frac{x_k}{\pi_k}}{\sum_s \frac{1}{\pi_k}}$$

The tilde mean \tilde{x}_s , being formed at the level of the intended sample s , employs sample weights only. (This type of mean can be calculated for the x -variable, which is observed for all $k \in s$, but obviously not for y -variable, which is observed for $k \in r$ only.)

The classical regression estimator formula corresponds, in our context, to

$$\hat{t}_{\text{REG}} = \hat{N} \{ \tilde{y}_r + b(\tilde{x}_s - \tilde{x}_r) \}$$

with

$$b = \frac{\sum_{h=1}^H f_h^{-1} \sum_{r_h} (y_k - \tilde{y}_r) (x_k - \tilde{x}_r) / \pi_k}{\sum_{h=1}^H f_h^{-1} \sum_{r_h} (x_k - \tilde{x}_r)^2 / \pi_k}$$

(Note: sample weighting as well as nonresponse weighting is used in b too.)

In summary, we have a series of three estimators

$$\hat{t}_{\text{EXP}} = \hat{N} \bar{y}_r, \quad (2.2a)$$

$$\hat{t}_{\text{RA}} = \hat{N} \bar{x}_s \frac{\bar{y}_r}{\bar{x}_r}, \quad (2.2b)$$

$$\hat{t}_{\text{REG}} = \hat{N} \{ \bar{y}_r + b (\bar{x}_s - \bar{x}_r) \}. \quad (2.2c)$$

All three are properly sample weighted and nonresponse weighted. The obvious differences have to do with the co-variate: \hat{t}_{EXP} uses no co-variate, whereas \hat{t}_{RA} and \hat{t}_{REG} do. It is also clear that \hat{t}_{RA} appeals to an underlying relationship between y and the co-variate x in the form of a line through the origin, the slope of which is estimated by \bar{y}_r/\bar{x}_r . In the case of \hat{t}_{REG} , the relationship is a regression with a non-zero intercept. We shall further explore the role of the co-variate.

If the population size N is known, it is in general better to replace \hat{N} by N in (2.2a) to (2.2c), yielding

$$\hat{t}_{\text{EXP}}^* = N \bar{y}_r, \quad (2.3a)$$

$$\hat{t}_{\text{RA}}^* = N \bar{x}_s \frac{\bar{y}_r}{\bar{x}_r}, \quad (2.3b)$$

$$\hat{t}_{\text{REG}}^* = N \{ \bar{y}_r + b (\bar{x}_s - \bar{x}_r) \}. \quad (2.3c)$$

For estimating the population total, N must be known in these three estimators, which may not be the case. However, for estimating the population mean \bar{Y} , they lead, by dividing by N , to the convenient expressions

$$\hat{Y}_{\text{EXP}} = \bar{y}_r, \quad (2.4a)$$

$$\hat{Y}_{\text{RA}} = \bar{x}_s \frac{\bar{y}_r}{\bar{x}_r}, \quad (2.4b)$$

$$\hat{Y}_{\text{REG}} = \bar{y}_r + b (\bar{x}_s - \bar{x}_r). \quad (2.4c)$$

The three series of estimators (2.2), (2.3), and (2.4) are easy to accept on intuitive grounds since all that is involved are elementary weighting principles, plus standard ratio feature or regression feature. Somewhat less elementary is to draw the proper consequences for variance estimation and the construction of valid confidence intervals. These questions are discussed in Section 4. (Contrary to what the rather informal presentation of the estimators (2.2) to (2.4) may suggest, the formulas are not “ad hoc” but the result of a formalized general estimation procedure (with a multivariate regression) for two phases of selection; see Särndal and Swensson (1985a). Most importantly, the variance estimators and confidence intervals follow directly from this theory.)

3. RESPONSE MODELS

The nonresponse weights in the estimators seen in Section 2 can be justified through a response mechanism model involving individual response probabilities that are constant for each unit in a given group. More formally, consider the response mechanism:

MODEL A:

- (1) The probability of response is constant (and equal to an unknown constant Θ_h) for all units $k \in s_h$; $h = 1, \dots, H$.
- (2) The units respond independently of each other.

The theoretical response probabilities Θ_h may vary considerably between groups. (An indication that large differences in response propensity may exist between different subsets is, of course, an incentive to set up adjustment groups, and to weight accordingly.)

Consider a fixed sample realization, s . The group frequencies $n_1, \dots, n_h, \dots, n_H$ are then fixed. Let us also consider a fixed value of the vector of group response frequencies $\underline{m} = (m_1, \dots, m_h, \dots, m_H)$. With s and \underline{m} fixed, the "selection" under Model A of a response set r_h can be shown to conform to a simple random selection of m_h from n_h . The conditional response probability of a unit k in the h -th group is therefore

$$\pi_{k|s,m} = \frac{m_h}{n_h} = f_h, \text{ all } k \in s_h. \tag{3.1}$$

(This consideration underlies the weight f_h^{-1} used in the estimators.) Similarly one can show that given s and \underline{m} , the probability under Model A that units k and l respond is

$$\pi_{kl|s,m} = \begin{cases} f_h & \text{if } k = l \\ \frac{f_h(m_h - 1)}{n_h - 1} & \text{if } k \neq l \in s_h \\ f_h f_{h'} & \text{if } k \in s_h; l \in s_{h'} \text{ (} h \neq h' \text{)} \end{cases} \tag{3.2}$$

($\pi_{kk|s,m}$ is by definition equal to $\pi_{k|s,m}$.) These quantities (which remind us of stratified random sampling with m_h units chosen from n_h in the h -th stratum) are important for the calculation of variance estimates and standard errors; see below.

In practice, the analyst decides how to set up his groups s_h . The decision is crucial, for it will determine the adjustment weights f_h^{-1} , and thus the numerical value of the estimate of t , the variance estimate, and the confidence interval. Two different groupings may lead to widely different point estimates and confidence intervals.

The analyst is not so naive as to think that response probabilities exist that are exactly equal within the group that he has identified. He does, however, believe (and usually with good reason) that more valid point estimates and confidence intervals will result with these groups (and thereby the weights f_h^{-1}) than without them. The adjustment group approach is a sound and firmly established practice.

On closer scrutiny, several things may be wrong with a response model such as Model A: the response probability is perhaps not constant within groups. And, even if it were, the particular groups postulated by the model are perhaps wrongly defined; there should have been more groups than assumed, etc. Two cases must therefore be distinguished for the continued discussion:

- (a) The assumed response mechanism (ARM; here in the form of Model A) is true. In practice, this is unlikely to be exactly the case.
- (b) The ARM is more or less false. This is the unpleasant truth in the majority of all practical situations, and it leads to nonresponse bias. In the case of Model A, the groups may be formed more or less incorrectly.

As is usual in statistics, the statistical analyst will formulate the model corresponding to the best of his judgement; accordingly, he will draw certain inferences (confidence statements, for example). Then he will wonder about the robustness of these conclusions, that is, how well do they hold up if the model is false? In the same order of things, let us consider these questions in our particular situation.

4. VARIANCE ESTIMATORS BASED ON A CERTAIN ASSUMED RESPONSE MECHANISM

Model A, with a specified set of groups, is assumed to hold. The response rates, $f_h = m_h/n_h$, $h = 1, \dots, H$, have been established. With this as a starting point, let us examine the variance estimators needed to construct a confidence interval at a specified $100(1 - \alpha)\%$ level. If \hat{t} is one of the estimators in Section 2, and Model A really holds, we have:

- (a) \hat{t} is unbiased (except for a usually unimportant technical bias)
- (b) an approximately $100(1 - \alpha)\%$ confidence interval for t is:

$$\hat{t} \pm z_{1-\alpha/2} \sqrt{\hat{V}(\hat{t})},$$

where the constant $z_{1-\alpha/2}$ is exceeded with probability $\alpha/2$ by the unit normal variate.

Under repeated draws of samples s and, for each fixed s , repeated realizations (obeying the assumed Model A) of response sets r , the interval will contain the true population total $100(1 - \alpha)\%$ of the time.

The variance and the estimated variance will be determined by two sets of selection probabilities:

1. π_k and π_{kl} , the probabilities of inclusion (first and second order) that accompany the sampling phase;
2. $\pi_{k|s,m}$, $\pi_{kl|s,m}$ the conditional response probabilities (first and second order) associated with the response Model A ("the nonresponse phase").

In our case, as a consequence of Model A, $\pi_{k|s,m}$, and $\pi_{kl|s,m}$, are given, respectively, by (3.1) and (3.2). As for π_k and π_{kl} , full generality is assumed; any design may be used for the sampling phase.

A detailed analysis will show that the total variance of any one of the estimators \hat{t} seen in Section 2 can be broken down into two components:

$$V(\hat{t}) = V_1(\hat{t}) + V_2(\hat{t})$$

where $V_1(\hat{t})$ may be termed the sampling variance and $V_2(\hat{t})$ the nonresponse variance. The exact formulas given in Särndal and Swensson (1985a) are not reproduced here, but one notes that the components have some reasonable properties:

1. $V_1(\hat{t}) = 0$ if the whole population U is observed (a census rather than a sample survey);

2. $V_2(\hat{t}) = 0$ if the response is complete ($r = s$);
3. $V_2(\hat{t})$ is greatly reduced in the presence of a strong co-variate, but $V_1(\hat{t})$ is not affected by the co-variate (naturally enough, since it is observed for $k \in s$ only).

Let us examine somewhat more closely the variance estimators. If $\hat{V}_i(\hat{t})$ denotes the estimator of $V_i(\hat{t})$, $i = 1, 2$, the total variance $V(\hat{t})$ will be estimated by an expression of the form

$$\hat{V}(\hat{t}) = \hat{V}_1(\hat{t}) + \hat{V}_2(\hat{t}).$$

Here, the estimated sampling variance component is

$$\hat{V}_1(\hat{t}) = \sum_{k \in r} \sum_{l \in r} \left(\frac{1}{\pi_k \pi_l} - \frac{1}{\pi_{kl}} \right) \frac{1}{\pi_{kl|s,m}} u_k u_l,$$

where $\pi_{kl|s,m}$ is given by (3.2), and π_k, π_{kl} are the inclusion probabilities of the sampling design. The estimated nonresponse variance component is

$$\hat{V}_2(\hat{t}) = \sum_{h=1}^H n_h^2 \left(\frac{1}{m_h} - \frac{1}{n_h} \right) S_{wrh}^2$$

with

$$S_{wrh}^2 = \frac{1}{m_h - 1} \sum_{r_h} (w_k - \bar{w}_{r_h})^2$$

The quantities u_k and w_k differ from one estimator \hat{t} to another. Let us look first at the estimated nonresponse variance, $\hat{V}_2(\hat{t})$. This component is of ‘‘stratified form’’: the factor $n_h^2(1/m_h - 1/n_h)$ is characteristic of a stratified simple random selection with m_h units chosen from n_h in the h -th stratum. The reason for this structure lies in the conditional response probabilities $\pi_{kl|s,m}$ given by (3.2).

The quantities w_h have the following appearance:

$$\text{For } \hat{t}_{\text{EXP}} \text{ and } \hat{t}_{\text{EXP}}^*: w_k = \frac{y_k - \bar{y}_r}{\pi_k},$$

$$\text{For } \hat{t}_{\text{RA}} \text{ and } \hat{t}_{\text{RA}}^*: w_k = \frac{y_k - (\bar{y}_r/\bar{x}_r)x_k}{\pi_k},$$

$$\text{For } \hat{t}_{\text{REG}} \text{ and } \hat{t}_{\text{REG}}^*: w_k = \frac{y_k - \bar{y}_r - b(x_k - \bar{x}_r)}{\pi_k}.$$

The expressions for w_k are sample weighted regression residuals. Consequently, if x_k is a powerful explanatory variable for y_k , one will ordinarily have that the variance of the w_k (and thus $\hat{V}_2(\hat{t})$) is smaller for the RA and REG estimators than for the EXP estimator, where the quantity w_k is just a deviation of y_k from the response set mean \bar{y}_r . Consequently, in fortunate circumstances, the part of the standard error that is due to the nonresponse will be reduced to near-zero levels, namely, when x and y have near perfect correlation.

The estimated sampling variance component $\hat{V}_1(\hat{t})$ is of less interest in this discussion, since it is not directly influenced by the co-variate. It should be mentioned, however, that the

u_k are determined as follows: \hat{t}_{EXP} , \hat{t}_{RA} , and \hat{t}_{REG} , $u_k = y_k$, while for the “starred” series of estimators \hat{t}_{EXP}^* , \hat{t}_{RA}^* , and \hat{t}_{REG}^* , $u_k = y_k - \hat{y}_s$, where $\hat{y}_s = (\sum_s \hat{y}_k / \pi_k) / (\sum_s 1 / \pi_k)$ is the mean of the predicted values from the regression fit, so that for \hat{t}_{EXP}^* , $\hat{y}_k = \hat{y}_r$ for all k ; for \hat{t}_{RA}^* , $\hat{y}_k = (\hat{y}_r / \hat{x}_r) x_k$; and for \hat{t}_{REG}^* , $\hat{y}_k = \hat{y}_r - b(x_k - \hat{x}_r)$.

A special case arises when $m_h = n_h$ for all h (that is, no nonresponse). Then $\hat{V}_2(\hat{t}) = 0$ (as is reasonable), and $\pi_{kl|s,m} = 1$ for all k and l , leaving the non-zero component

$$\hat{V}_1(\hat{t}) = \sum_{k \in r} \sum_{l \in r} \left(\frac{1}{\pi_k \pi_l} - \frac{1}{\pi_{kl}} \right) u_k u_l$$

which is the well-known variance estimator for the case of full response.

5. ROBUSTNESS PROPERTIES WHEN THE ASSUMED RESPONSE MECHANISM IS FALSE

Unbiased estimates and valid confidence intervals can be obtained with the aforementioned estimators, provided the ARM (given by Model A) holds. The presence of a strong co-variate brings about a reduction of the nonresponse component of the variance.

More interesting in a real-life situation is the case where the ARM breaks down. This case must be considered, because even the most careful judgement in setting up adjustment groups is bound to be less than perfect. The extent of the departure of the true response behaviour from that of the ARM will now determine behaviour of the various estimators. The statistical properties (bias, coverage rate achieved by confidence intervals, etc.) are in other words functions of the extent of model breakdown.

In Särndal and Swensson (1985a), a small scale Monte Carlo experiment was carried out to study the impact of certain types of breakdown in Model A. For purposes of illustration, we cite a few results from this study.

The true ARM in the experiment had $H = 4$ adjustment groups, with different response probabilities between groups (but constant response probability for all units in the same group). 1,000 simple random samples were drawn, and each sample was exposed to simulated nonresponse according to the true ARM (which is taken as known, since this is a controlled experiment).

As expected from theory, when the ARM underlying \hat{t}_{EXP} and \hat{t}_{RA} was true, there is essentially no bias, and the empirical coverage rates of the confidence intervals agree essentially with the nominal 95% rate. The advantage of \hat{t}_{RA} lies in a smaller component of variance due to nonresponse. (See “ARM is true” in Table 1.)

False ARM’s were created by joining together groups of the true ARM. The estimator and the confidence interval (based on the false ARM) will then be calculated on the basis of fewer groups than ought to be the case. The case “ARM is false” in Table 1 represents the extreme situation where all four groups of the true ARM were joined into one, meaning that one acts in the estimation process as if all units throughout the population had the same (unknown, but estimated) response probability. The table shows that the co-variate estimator, \hat{t}_{RA} , when compared to the no-co-variate estimator, \hat{t}_{EXP} , has the following (not unexpected) advantages: (a) strong resistance to nonresponse bias (1.26 versus 4.85); (b) much better preservation of the nominal 95% confidence coefficient (92.6% versus 46.3% empirical coverage rate). In addition, \hat{t}_{RA} has a variance advantage, and therefore shorter confidence intervals on the average.

Table 1
Comparison of \hat{t}_{EXP} and \hat{t}_{RA}

Estimator	Absolute bias	Mean of the variance component \hat{V}_2	Empirical coverage rate (95% nominal)
ARM is true	\hat{t}_{EXP}	0.00	95.2%
	\hat{t}_{RA}	-0.01	95.5%
ARM is false	\hat{t}_{EXP}	4.85	46.3%
	\hat{t}_{RA}	1.26	92.6%

6. CONCLUSION

In summary, we have argued in this paper that two different categories of variables (observed for k in the intended sample s) are of importance:

- variables suitable for estimating the response mechanism (in the case of Model A, these variables allow the construction of the adjustment groups);
- variables (here called co-variates) that are powerful predictors of the y -variable; when used in the estimator formula, they reduce variance and improve the robustness properties.

Whenever possible, one should thus be on the outlook for suitable co-variates. One should also note that when several y -totals are to be estimated, the appropriate co-variates may differ from one y -variable to the other, whereas the weighting classes would probably be set up to apply uniformly for all variables of interest.

REFERENCES

- LITTLE, R.J.A. (1983). Models for nonresponse in sample surveys. *Journal of the American Statistical Association*, 77, 237-250.
- SÄRNDAL, C.E., and SWENSSON, B. (1985a). A general view of estimation for two phases of selection. Part I: Randomized subsample selection (Two-phase sampling). Part II: Nonrandomized subsample selection (Nonresponse). Promemorior fran P/STM no. 20, Statistics Sweden.
- SÄRNDAL, C.E., and SWENSSON, B. (1985b). Incorporating nonresponse modelling in a general randomization theory approach. *Bulletin of the International Statistical Institute* (45th session), 51:3, 15.2.1-16.
- OH, H.L., and SCHEUREN, F.J. (1983). Weighting adjustment for unit non-response. In *Incomplete Data in Sample Surveys*, Vol. 2, (Eds. W.G. Madow, I. Olkin, and D.B. Rubin), New York: Academic Press, 143-183.