

**TESTING BIAS OF MEAN ESTIMATES
DUE TO NOT MISSING AT RANDOM NONRESPONSE
WITH APPLICATION TO THE CONSUMER EXPENDITURE SURVEY**

Michael Sverchkov
Bureau of Labor Statistics, Washington DC, USA

1. INTRODUCTION

There exists almost no survey without nonresponse. In practice most survey methods that adjust for nonresponse assume either explicitly or implicitly that the missing data are ‘missing at random’ (MAR). That is, they assume response does not depend on the variable of interest (the outcome variable) given some auxiliary information known for the whole population. However, in many practical situations, this assumption is not valid, since the probability of responding often depends on the outcome value, even after conditioning on available covariate information. In such cases, the use of methods that assume nonresponse is MAR can lead to large bias of parameter estimators and distort subsequent inference.

The case where the missing data are not MAR (NMAR) can be treated by postulating a parametric model for the distribution of the outcomes before nonresponse and a model for the response mechanism. These two models define a parametric model for the observed outcomes, so that the parameters of these models can be estimated from the observed data. Once the parameters are estimated, the first model can be used for inference. See, for example, Pfeiffermann and Sverchkov (2009) for details, with overview of related literature.

Modeling the distribution of the outcomes before nonresponse is difficult since only the observed data are available. Sverchkov (2008) proposes an alternative approach, which allows the parameters of the response model to be estimated without postulating a parametric model for the distribution of the outcomes before nonresponse. To account for the nonresponse, Sverchkov (2008) assumes a response model and estimates the response probabilities by applying the missing information principle (MIP), which consists of defining the likelihood as if there was complete response, and then integrating out the

unobserved outcomes from the likelihood, employing the relationship between the distributions of the observed and unobserved data. We describe the main steps of this approach in Sections 2, 3 and 4.

In this paper, Section 5, we show how this approach can be applied to testing whether probability weighted estimators, corrected for MAR nonresponse, are still biased due to ignoring that nonresponse is NMAR actually. In Section 6 We illustrate this approach on real data (Consumer Expenditure Survey) example.

2. NOTATION AND MODELS

Let $\{y_i, \mathbf{x}_i; i = 1, \dots, N\}$ represent the data in a finite population of N units, where y_i is the value of the outcome variable for unit i and $\mathbf{x}_i' = (x_{i,1}, \dots, x_{i,K})$ is a vector of corresponding K covariates. Suppose that the population outcome values follow model (2.1):

$$y_i | \mathbf{x}_i \sim f(y_i | \mathbf{x}_i), \quad i = 1, \dots, N \quad (2.1)$$

The target is to estimate the population mean $\bar{Y} = N^{-1} \sum_{i=1}^N y_i$, based on a sample s of n units with inclusion probabilities $\pi_i = \Pr(i \in s)$. Denote by I_i the sample indicator; $I_i = 1$ if unit i is selected in the sample and 0 otherwise. Let $w_i = 1/\pi_i$ denote the sampling weights.

In practice, not every unit in the sample responds. Define the response indicator; $R_i = 1$ if unit $i \in s$ responds and $R_i = 0$ otherwise. The sample of respondents is thus $R = \{i : I_i = 1, R_i = 1\}$ and the sample of nonrespondents among the sampled units is $R^c = \{i : I_i = 1, R_i = 0\}$.

Assumption 1. The response process is assumed to occur stochastically, independently between units, and $\sum_{i=1}^n R_i > 0$. (This assumption is used in the following Eq. 3.1 in Section 3 for simplification.)

The sample of respondents defines therefore a second, self-selected stage of the sampling process with unknown response probabilities. (Särndal and Swensson, 1987).

Under the population model (2.1), the observed data follow ‘respondents’ model:

$$f_R(y_i | \mathbf{x}_i) = f(y_i | \mathbf{x}_i, i \in R) \quad (2.2)$$

The model (2.2) is again general and all that we state at this stage is that under informative sampling and/or NMAR nonresponse, the models for the respondents and for the population differ; $f_R(y_i | \mathbf{x}_i) \neq f(y_i | \mathbf{x}_i)$.

Remark 1. The respondents’ model refers to the observed data and hence can be estimated and tested by standard methods.

Let $p_r(y_i, \mathbf{x}_i) = \Pr(R_i = 1 | y_i, \mathbf{x}_i, i \in s)$. If the probabilities $p_r(y_i, \mathbf{x}_i)$ were known, the sample of respondents could be considered as a sample from the finite population with known sampling probabilities $\tilde{\pi}_i = \pi_i p_r(y_i, \mathbf{x}_i)$. In this case, the population mean \bar{Y} can

be estimated, for example, by probability weighted estimator, $\hat{Y} = \frac{\sum_{i=1}^n y_i / \tilde{\pi}_i}{\sum_{i=1}^n 1 / \tilde{\pi}_i}$. Also, if

known, the response probabilities could be used for imputation of the missing data within the selected areas, by applying the relationship between the sample and sample-complement distributions, (Sverchkov and Pfeffermann, 2004);

$$f(y_i | \mathbf{x}_i, i \in R^c) = \frac{[p_r^{-1}(y_i, \mathbf{x}_i) - 1]f(y_i | \mathbf{x}_i, i \in R)}{E\{[p_r^{-1}(y_i, \mathbf{x}_i) - 1] | \mathbf{x}_i, i \in R\}} \quad (2.3)$$

(here and in what follows $a^{-1} = 1/a$).

Remark 2. Equation (2.3) is a key relationship for derivation of the main result of the next section, Equation (3.2). It allows the distribution of the unobserved outcomes to be written as a function of the distribution of the observed outcomes and response and/or selection model.

3. ESTIMATION OF RESPONSE PROBABILITIES

Unlike the sampling probabilities, the response probabilities are generally unknown. We assume therefore a parametric model, which is allowed to depend on the outcome and the covariate values; $\Pr[R_i = 1 | y_i, \mathbf{x}_i, i \in s; \boldsymbol{\gamma}] = p_r(y_i, \mathbf{x}_i; \boldsymbol{\gamma})$, where $\boldsymbol{\gamma}$ is a vector of unknown coefficients. We assume that $p_r(y_i, \mathbf{x}_i; \boldsymbol{\gamma})$ is differentiable with respect to $\boldsymbol{\gamma}$. In Section 4 we suggest a method for selection the response model.

Under these assumptions and Assumption 1, if the missing outcome values were observed, $\boldsymbol{\gamma}$ could be estimated by solving the likelihood equations:

$$\sum_{i \in R} \frac{\partial \log p_r(y_i, \mathbf{x}_i; \boldsymbol{\gamma})}{\partial \boldsymbol{\gamma}} + \sum_{k \in R^c} \frac{\partial \log [1 - p_r(y_k, \mathbf{x}_k; \boldsymbol{\gamma})]}{\partial \boldsymbol{\gamma}} = 0. \quad (3.1)$$

In practice, the missing data are unobserved and hence the likelihood equations (3.1) are not operational. However, one may apply in this case the missing information principle:

Missing Information Principle (MIP, Cepillini et al. 1955, Orchard and Woodbury, 1972): Let $O = \{y_i, n, i \in R; \mathbf{x}_i, t = 1, \dots, n\}$ represent the known observed data used below. Since no observations are available for $i \in R^c$, solve instead,

$$\begin{aligned} E \left\{ \left[\sum_{i \in R} \frac{\partial \log p_r(y_i, \mathbf{x}_i; \boldsymbol{\gamma})}{\partial \boldsymbol{\gamma}} + \sum_{k \in R^c} \frac{\partial \log [1 - p_r(y_k, \mathbf{x}_k; \boldsymbol{\gamma})]}{\partial \boldsymbol{\gamma}} \right] \middle| O \right\} & \stackrel{\text{by (2.3)}}{=} \sum_{i \in R} \frac{\partial \log p_r(y_i, \mathbf{x}_i; \boldsymbol{\gamma})}{\partial \boldsymbol{\gamma}} \\ + \sum_{k \in R^c} \frac{E \left\{ [p_r^{-1}(y_k, \mathbf{x}_k; \boldsymbol{\gamma}) - 1] \frac{\partial \log [1 - p_r(y_k, \mathbf{x}_k; \boldsymbol{\gamma})]}{\partial \boldsymbol{\gamma}} \middle| \mathbf{x}_k, k \in R \right\}}{E \{ [p_r^{-1}(y_k, \mathbf{x}_k; \boldsymbol{\gamma}) - 1] \mid \mathbf{x}_k, k \in R \}} & = 0. \end{aligned} \quad (3.2)$$

Notice that the expectations in the last expression are with respect to the model holding for the observed data for the respondents. This result is possible because of key relation (2.3), see Remark 2 See Sverchkov (2008) for detailed derivation of (3.2).

Remark 3. When the response probabilities $p_r(y_i, \mathbf{x}_i; \boldsymbol{\gamma})$ depend on only \mathbf{x}_i , they are referred to as *propensity scores*, and the missing data are missing at random. The estimating equations in (3.2) reduce in this case to the common log-likelihood equations,

$$\sum_{i \in R} \frac{\partial \log p_r(\mathbf{x}_i; \boldsymbol{\gamma})}{\partial \boldsymbol{\gamma}} + \sum_{k \in R^c} \frac{\partial \log [1 - p_r(\mathbf{x}_k; \boldsymbol{\gamma})]}{\partial \boldsymbol{\gamma}} = 0, \quad (3.3)$$

where $p_r(\mathbf{x}_i; \boldsymbol{\gamma}) = \Pr(R_i = 1 | \mathbf{x}_i; \boldsymbol{\gamma})$.

Imagine an unrealistic situation when we know response probabilities, $\Pr[R_i = 1 | y_i, \mathbf{x}_i, i \in S; \boldsymbol{\gamma}^*] = p_r(y_i, \mathbf{x}_i; \boldsymbol{\gamma}^*)$, where $\boldsymbol{\gamma}^*$ is a true value of $\boldsymbol{\gamma}$, but we still want to estimate $\boldsymbol{\gamma}$ by solving (3.2). Then the equations (3.2) can be solved by maximizing the log-likelihood,

$$l(\boldsymbol{\gamma}) = E \left(\sum_{i \in R} \log p_r(y_i, \mathbf{x}_i; \boldsymbol{\gamma}) + \sum_{i \in R^c} [1 - \log p_r(y_i, \mathbf{x}_i; \boldsymbol{\gamma})] \middle| \mathcal{O} \right) = \sum_{i \in R} \log p_r(y_i, \mathbf{x}_i; \boldsymbol{\gamma}) + \sum_{k \in R^c} \frac{E\{[p_r^{-1}(y_k, \mathbf{x}_k; \boldsymbol{\gamma}^*) - 1] \log[1 - p_r(y_k, \mathbf{x}_k; \boldsymbol{\gamma})] | \mathbf{x}_k, k \in R\}}{E_{re}\{[p_r^{-1}(y_k, \mathbf{x}_k; \boldsymbol{\gamma}^*) - 1] | \mathbf{x}_k, k \in R\}}. \quad (3.4)$$

Although such scenario is unrealistic, it suggests the following iteration algorithm of solving (3.2): starting with some initial estimate of parameter $\boldsymbol{\gamma}$, say $\boldsymbol{\gamma}^0$, for example $\boldsymbol{\gamma}^0$ is a solution of (3.3), maximize in the (q+1) iteration the expression,

$$\sum_{i \in R} \log p_r(y_i, \mathbf{x}_i; \boldsymbol{\gamma}^{(q+1)}) + \sum_{k \in R^c} \frac{E\{[p_r^{-1}(y_k, \mathbf{x}_k; \boldsymbol{\gamma}^{(q)}) - 1] \log[1 - p_r(y_k, \mathbf{x}_k; \boldsymbol{\gamma}^{(q+1)})] | \mathbf{x}_k, k \in R\}}{E\{[p_r^{-1}(y_k, \mathbf{x}_k; \boldsymbol{\gamma}^{(q)}) - 1] | \mathbf{x}_k, k \in R\}} \quad (3.5)$$

with respect to $\boldsymbol{\gamma}^{(q+1)}$. The maximization can be carried out, for example, by SAS Proc NLIN.

Remark 4. A fundamental question regarding the solution of the MIP equations is the existence of a unique global solution or more generally, the identifiability of the response model. Riddles et al. (2016) propose a similar approach to deal with NMAR nonresponse in the general context of survey sampling inference and establish the following fundamental condition for the response model identifiability: the covariates \mathbf{x} can be decomposed as $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)$ with $\dim(\mathbf{x}_2) \geq 1$, such that $\Pr(R_i = 1 | y_i, \mathbf{x}_i) = \Pr(R_i = 1 | y_i, \mathbf{x}_{1i})$. In other words, the covariates in \mathbf{x}_2 that appear in the outcome model do not affect the response probabilities, given the outcome and the other covariates.

Riddles et al. (2016) prove asymptotic normality of the estimate $\hat{\boldsymbol{\gamma}}$ under general regularity conditions.

4. SELECTION OF A RESPONSE MODEL

When following the approach proposed in Section 3, the likelihood (3.4) suggests at least two procedures for the selection of the response model under NMAR nonresponse. One procedure involves comparing different models based on information criteria such as the Akaike information criterion, $AIC = -2l(\boldsymbol{\gamma}) + 2\dim(\boldsymbol{\gamma})$, or Schwarz information criterion, $BIC = -2l(\boldsymbol{\gamma}) + \dim(\boldsymbol{\gamma})\log(n)$, $n = \sum_{i \in s} n_i$; a second procedure involves testing a saturated versus a nested model based on the likelihood ratio test. In Section 5 we illustrate via a simulation study how the likelihood (3.4) can be used for the application of these selection procedures.

5. TESTING BIAS OF MEAN ESTIMATES DUE TO NOT MISSING AT RANDOM NONRESPONSE

Let \hat{w}_i^{MAR} be a good estimate of $w_i \frac{1}{p_r(\mathbf{x}_i)} = w_i \frac{1}{\Pr(R_i = 1 | \mathbf{x}_i, i \in s)}$, in other words

\hat{w}_i^{MAR} is a good estimate of final inclusion weight if (or “assuming”) the nonresponse is

MAR. Then the obvious estimate for the mean will be $\hat{Y}^{MAR} = \frac{\sum_{i=1}^n \hat{w}_i^{MAR} y_i}{\sum_{i=1}^n \hat{w}_i^{MAR}}$ and if the

response is really MAR then this estimate is consistent.

How one can check in real situation that the latter estimate is not biased? Formally one can estimate response probabilities $\Pr[R_i = 1 | y_i, \mathbf{x}_i, i \in s; \boldsymbol{\gamma}] = p_r(y_i, \mathbf{x}_i; \boldsymbol{\gamma})$ following the same approach as in Section 3 and then compare model $p_r(y_i, \mathbf{x}_i; \boldsymbol{\gamma})$ with $p_r(\mathbf{x}_i; \boldsymbol{\gamma})$ by some information criterion as in Section 4 or check whether MAR and NMAR estimates,

$$\hat{Y}^{MAR} = \frac{\sum_{i=1}^n \hat{W}_i^{MAR} y_i}{\sum_{i=1}^n \hat{W}_i^{MAR}} \quad \text{and} \quad \hat{Y}^{NMAR} = \frac{\sum_{i=1}^n \hat{W}_i y_i}{\sum_{i=1}^n \hat{W}_i}, \quad \hat{W}_i = w_i \frac{1}{p_r(y_i, \mathbf{x}_i, \hat{\gamma})},$$

are significantly different (use $\widehat{\text{var}}(\hat{Y}^{MAR})$ for the test, the latter usually available).

As one can see, estimation procedure (iterative maximization of the likelihood (3.5)) is complicated optimization problem if dimension of vector of covariates \mathbf{x}_i is big, which is usually the case in real applications. Moreover, often correction to the response mechanism is made through different calibration procedures and so that response model, $\Pr[R_i = 1 | y_i, \mathbf{x}_i, i \in s; \gamma] = p_r(y_i, \mathbf{x}_i; \gamma)$, can be hard to identify. Assuming that NMAR

estimate is consistent,
$$\text{Bias}(\hat{Y}^{MAR}) = \frac{\sum_{j=1}^N y_j}{N} - \hat{Y}^{MAR} = \frac{\sum_{j=1}^N (y_j - \hat{Y}^{MAR})}{N} \cong - \frac{\sum_{i=1}^n \hat{W}_i \varepsilon_i^{MAR}}{\sum_{i=1}^n \hat{W}_i},$$

where $\varepsilon_i^{MAR} = y_i - \frac{\sum_{i=1}^n \hat{W}_i^{MAR} y_i}{\sum_{i=1}^n \hat{W}_i^{MAR}}$, and therefore one can test whether \hat{Y}^{MAR} is biased or not

by testing if $\Pr(R_i = 1 | \varepsilon_i^{MAR})$ depends on ε_i^{MAR} or not which can be done by the approach suggested in Sections 3 and 4. Therefore dimension of optimization procedure can be reduced dramatically. Note that although $\Pr(R_i = 1 | \varepsilon_i^{MAR})$ does not involve any covariates, we still need a covariate \mathbf{x}_{2i} such that ε_i^{MAR} correlates with \mathbf{x}_{2i} but $\Pr(R_i = 1 | \varepsilon_i^{MAR}, \mathbf{x}_{2i}) = \Pr(R_i = 1 | \varepsilon_i^{MAR})$, otherwise the solution of (3.5) can be not unique, see Remark 4.

6. APPLICATION TO CONSUMER EXPENDITURE SURVEY

We consider estimation of the following expenditures, Total, Food, Housing, Health, based on 2019, 2020 and 2021 samples. We assume truncated logit model for response probability

$$\Pr[R_i = 1 | \varepsilon_i^{MAR}, i \in s; \gamma] = p_r(\varepsilon_i^{MAR}; \gamma) = \frac{\left[0.0001 + \frac{\exp(\gamma_0 + \gamma_1 \varepsilon_i^{MAR})}{1 + \exp(\gamma_0 + \gamma_1 \varepsilon_i^{MAR})} \right]}{1.0001}$$

where ε_i^{MAR} is defined as in Section 5 by a particular expenditure and MAR final sample weight, FINLWT21, used currently by Consumer Expenditure Survey. We use “median income” variable as instrumental variable \mathbf{x}_{2i} (see Remark 4), correlation coefficients between \mathbf{x}_{2i} and Total, Food, Housing, Health are equal to 0.31, 0.26, 0.35, 0.08 respectively for year 2021, and simple linear model explain relation between the expenditures and \mathbf{x}_{2i} reasonably although with low R-square statistics, 0.15, 0.07, 0.16, 0.20. Parameter estimates for response model, $\Pr[R_i = 1 | \varepsilon_i^{MAR}, i \in s; \gamma] = p_r(\varepsilon_i^{MAR}; \gamma)$, are as follows (all significant with p-value < 0.0001, p-values are produced by SAS Proc NLIN used in maximization of the likelihood (3.4)-(3.5), this suggests that the response model is NMAR not MAR):

Total: $\gamma_0 = -0.513$, $\gamma_1 = 0.0000085$.

Food: $\gamma_0 = -0.510$, $\gamma_1 = 0.0000086$.

Housing: $\gamma_0 = -0.496$, $\gamma_1 = 0.0000303$.

Health: $\gamma_0 = -0.619$, $\gamma_1 = 0.0007000$.

(Similar results for these relationships for samples of year 2019 and 2020).

AIC and BIC criteria for all 4 expenditures and all years considered in this research prefer response model that include ε_i^{MAR} as a covariate (NMAR nonresponse) versus response model that assume MAR response, results summarized in the table below.

		AIC		BIC		$\frac{Bi\hat{a}s(\hat{Y}^{MAR})}{\hat{std}(\hat{Y}^{MAR})}$
		MAR	NMAR	MAR	NMAR	
Total	2019	68372.9	66794.6	68380.9	66810.6	1.00
Total	2020	70538.0	69337.4	70545.9	69353.3	0.49
Total	2021	71907.6	71082.4	71915.6	71098.2	1.67
Food	2019	68372.9	66872.5	68380.9	66888.4	0.71
Food	2020	70538.0	69337.4	70545.9	69353.3	2.50
Food	2021	71907.6	71159.3	71915.6	71175.1	0.65
Housing	2019	68372.9	66858.2	68380.9	66874.2	0.01
Housing	2020	70538.0	69366.4	70545.9	69382.0	2.86
Housing	2021	71907.6	71049.2	71915.6	71065.1	3.15
Health	2019	68372.9	66872.1	68380.9	66888.2	-1.11
Health	2020	70538.0	69337.4	70545.9	69353.3	-0.66
Health	2021	71907.6	71142.1	71915.6	71157.9	0.28

The last column of the table presents an estimate of normalized bias of MAR estimator,

$$\frac{Bi\hat{a}s(\hat{Y}^{MAR})}{\hat{std}(\hat{Y}^{MAR})}, \text{ where } Bi\hat{a}s(\hat{Y}^{MAR}) = -(\hat{Y}^{NMAR} - \hat{Y}^{MAR}) = -\frac{\sum_{i=1}^n \hat{w}_i \varepsilon_i^{MAR}}{\sum_{i=1}^n \hat{w}_i}, \text{ and } \hat{std}(\hat{Y}^{MAR}) \text{ is}$$

a BRR standard error estimate of \hat{Y}^{MAR} . Note that the later statistics can be used only as

crude approximation of normalize bias $\frac{Bias(\hat{Y}^{MAR})}{std[Bias(\hat{Y}^{MAR})]}$, since it does not include errors

due to \hat{Y}^{NMAR} estimation.

The above three statistics (naïve approximations) suggest that the biases of the estimate due to ignoring NMAR nonresponse are not (or slightly, see Hosing) significant based on the standard errors of the final estimates. Note that standard error of $Bi\hat{a}s(\hat{Y}^{MAR})$ is a standard error of the difference of MAR and NMAR estimators, therefore the standard error of \hat{Y}^{MAR} can be used only for crude comparison of two estimators.

Conclusion: Although current Consumer Expenditure Survey nonresponse adjustment of sample weights does not remove nonresponse bias (MAR estimates for the mean total expenditures are still biased after the adjustment), the bias is not or only slightly significant based on the standard error of the final estimates.

REFERENCES

- Cepillini, R., Siniscalco, M., and Smith, C.A.B. (1955). The estimation of gene frequencies in a random mating population. *Annals of Human Genetics*, **20**, 97-115.
- Orchard, T., and Woodbury, M.A. (1972). A missing information principle: theory and application. *Proceedings of the 6th Berkeley Symposium on Mathematical Statistics and Probability*, **1**, 697-715.
- Pfeffermann, D., and Sverchkov, M. (2009). Inference under Informative Sampling. In: *Handbook of Statistics 29B; Sample Surveys: Inference and Analysis*. Eds. D. Pfeffermann and C.R. Rao. North Holland. pp. 455-487.
- Riddles, K.M., Kim, J.K. and Im, J. (2016). A propensity-score adjustment method for nonignorable nonresponse. *Journal of Survey Statistics and Methodology*, **4**, 215-245.
- Särndal, C.E. and Swensson B. (1987). A general view of estimation for two phases of selection with applications to two-phase sampling and nonresponse. *International Statistical Review*, **55**, 279-294.

Sverchkov, M. (2008). A new approach to estimation of response probabilities when missing data are not missing at random. *Joint Statistical Meetings, Proceedings of the Section on Survey Research Methods*, 867-874.

Sverchkov, M., and Pfeffermann, D. (2004). Prediction of Finite Population Totals Based on the Sample Distribution. *Survey Methodology*, **30**, 79-92.