

RESPONSE MODEL SELECTION IN SMALL AREA ESTIMATION UNDER NOT MISSING AT RANDOM NONRESPONSE

By Michael Sverchkov

Calcutta Statistical Association Bulletin

This manuscript has been submitted to Calcutta Statistical Association Bulletin

Journal Name: Calcutta Statistical Association Bulletin	Manuscript ID: CSA-2023-0009.RV1
Manuscript Type: Original articles in statistics & probability	Manuscript Title: RESPONSE MODEL SELECTION IN SMALL AREA ESTIMATION UNDER NOT MISSING AT RANDOM NONRESPONSE
Author Names: Michael Sverchkov, Danny Pfeffermann	
Keywords: efficiency, probability, statistical application, Statistics	
MeSH terms:	
<p>Abstract: Sverchkov and Pfeffermann (2018) consider Small Area Estimation (SAE) under informative probability sampling of areas and within the sampled areas, and not missing at random (NMAR) nonresponse. To account for the nonresponse, the authors assume a given response model, which contains the outcome values as one of the covariates and estimate the corresponding response probabilities by application of the Missing Information Principle, which consists of defining the likelihood as if there was complete response and then integrating out the unobserved outcomes from the likelihood by employing the relationship between the distributions of the observed and unobserved data.</p> <p>A key condition for the success of this approach is the “correct” specification of the response model. In this presentation we consider the likelihood ratio test and information criteria based on the appropriate likelihood and show how they can be used for the selection of the response model. We illustrate the approach by a small simulation study.</p>	

Peer Review

3 1. INTRODUCTION

29

30 There exists almost no survey without nonresponse, but in practice most methods
31 that deal with this problem assume either explicitly or implicitly that the missing
32 data are 'missing at random' (MAR). However, in many practical situations, this
33 assumption is not valid, since the probability to respond often depends on the
34 outcome value, even after conditioning on available covariate information. In such
35 cases, the use of methods that assume that the nonresponse is MAR can lead to
36 large bias of parameter estimators and distort subsequent inference.

37

38 The case where the missing data are not MAR (NMAR) can be treated by
39 postulating a parametric model for the distribution of the outcomes before non-
40 response and a model for the response mechanism. These two models define a
41 parametric model for the observed outcomes, so that the parameters of these
42 models can be estimated from the observed data. See, for example, Pfeiffermann
43 and Sverchkov (2009) for details, with overview of related literature.

44

45 Modeling the distribution of the outcomes before non-response can be problematic
46 since only the observed data are available. Sverchkov (2008) proposes an
47 alternative approach, which allows to estimate the parameters of the response
48 model without postulating a parametric model for the distribution of the outcomes
49 before nonresponse. To account for the nonresponse, Sverchkov (2008) assumes
50 a given response model and estimates the corresponding response probabilities
51 by application of the missing information principle (MIP), which consists of defining
52 the likelihood as if there was complete response, and then integrating out the
53 unobserved outcomes from the likelihood, employing the relationship between the
54 distributions of the observed and unobserved data. Sverchkov and Pfeiffermann
55 (2018) apply this approach for small area estimation (SAE) under informative
56 probability sampling of areas and within the sampled areas, and NMAR
57 nonresponse. We describe the main steps of this approach in Sections 2 and 3.

58

2
 59 A key condition for the success of this approach is the “correct” specification of the
 60 response model. In section 4 we consider the likelihood ratio test and information
 61 criteria based on the appropriate likelihood and show how they can be used for the
 62 selection of the response model. Section 5 illustrates the application of the
 63 approach by a small simulation study.

2. NOTATION AND MODELS

5
 65 Let $\{y_{ij}, \mathbf{x}_{ij}; i = 1, \dots, M, j = 1, \dots, N_i\}$ represent the data in a finite population of N
 66 units, comprised of M areas with N_i units in area i , $\sum_{i=1}^M N_i = N$, where y_{ij} is
 67 the value of the outcome variable for unit j in area i and $\mathbf{x}'_{ij} = (x_{ij,1}, \dots, x_{ij,K})$ is a
 68 vector of corresponding K covariates. We assume that the covariates are known
 69 for every unit in the population. Suppose that the population outcome values follow
 70 the generic two-level model:

$$1 \quad y_{ij} | \mathbf{x}_{ij}, u_i^U \sim f(y_{ij} | \mathbf{x}_{ij}, u_i^U), \quad i = 1, \dots, M, j = 1, \dots, N_i \quad (2.1)$$

$$u_i^U \sim f(u_i^U); \quad E(u_i^U) = 0, \quad V(u_i^U) = \sigma_{u^U}^2,$$

1
 72 where u_i^U is the i^{th} area level random effect. The target is to estimate the area
 73 means $\bar{Y}_i = N_i^{-1} \sum_{j=1}^{N_i} y_{ij}, i = 1, \dots, M$, based on a sample obtained by the following
 74 two-stage sampling scheme: **i)** select a sample s of m out of the M population
 75 areas with inclusion probabilities $\pi_i = \Pr(i \in s)$; **ii)** select a sample s_i of $n_i > 0$
 76 units from selected area i with probabilities $\pi_{ji} = \Pr(j \in s_i | i \in s)$. Denote by I_i ,
 77 I_{ij} the sample indicators; $I_i = 1$ if area i is selected in the first stage and 0
 78 otherwise, $I_{ij} = 1$ if unit j of selected area i is sampled in the second stage and
 79 $I_{ij} = 0$ otherwise. Let $w_i = 1 / \pi_i$, $w_{ji} = 1 / \pi_{ji}$ denote the first- and second-stage
 80 sampling weights.

81
 82 In practice, not every unit in the sample responds. Define the response indicator;
 83 $R_{ij} = 1$ if unit $j \in s_i$ responds and $R_{ij} = 0$ otherwise. The sample of respondents

84 is thus $R = \{(i, j) : I_i = 1, I_{ij} = 1, R_{ij} = 1\}$ and the sample of nonrespondents among
 85 the sampled units is $R^c = \{(i, k) : I_i = 1, I_{ik} = 1, R_{ik} = 0\}$. The response process is
 86 assumed to occur stochastically, independently between units. We assume
 87 $\sum_{j=1}^{n_i} R_{ij} > 0$ in all the sampled areas. The sample of respondents defines therefore
 88 a third, self-selected stage of the sampling process with unknown response
 89 probabilities. (Särndal and Swensson, 1987).

90 Define, $u_i = u_i^U - E(u_i^U | i \in s)$. Then, under the population model (2.1), the
 91 observed data follow the two-level 'respondents' model:

$$\begin{aligned}
 f_R(y_{ij} | \mathbf{x}_{ij}, u_i) &= f(y_{ij} | \mathbf{x}_{ij}, u_i, (i, j) \in R); \\
 u_i &\sim f(u_i | i \in s), E(u_i | i \in s) = 0.
 \end{aligned}
 \tag{2.2}$$

93 The model (2.2) is again general and all that we state at this stage is that under
 94 informative sampling and/or NMAR nonresponse, the population and the
 95 respondents' models differ; $f_R(y_{ij} | \mathbf{x}_{ij}, u_i) \neq f(y_{ij} | \mathbf{x}_{ij}, u_i^U)$.

96 *Remark 1.* The respondents' model refers to the observed data and hence can be
 97 estimated and tested by standard SAE methods. See Pfeffermann (2013) and Rao
 98 and Molina (2015) for estimation and testing procedures in SAE, with references.

99 Let $p_r(y_{ij}, \mathbf{x}_{ij}) = \Pr[R_{ij} = 1 | y_{ij}, \mathbf{x}_{ij}, i \in s, j \in s_i]$. If the probabilities $p_r(y_{ij}, \mathbf{x}_{ij})$ were
 100 known, the sample of respondents could be considered as a two-stage sample
 101 from the finite population with known sampling probabilities π_i and
 102 $\tilde{\pi}_{j|i} = \pi_{j|i} p_r(y_{ij}, \mathbf{x}_{ij})$. In this case, the area means \bar{Y}_i can be estimated as in
 103 Pfeffermann and Sverchkov (2007). Also, if known, the response probabilities
 104 could be used for imputation of the missing data within the selected areas, by
 105 application of the relationship between the sample and sample-complement
 106 distributions, (Sverchkov and Pfeffermann, 2004);

$$f(y_{ij} | \mathbf{x}_{ij}, u_i, (i, j) \in R^c) = \frac{[p_r^{-1}(y_{ij}, \mathbf{x}_{ij}) - 1] f(y_{ij} | \mathbf{x}_{ij}, u_i, (i, j) \in R)}{E\{[p_r^{-1}(y_{ij}, \mathbf{x}_{ij}) - 1] | \mathbf{x}_{ij}, u_i, (i, j) \in R\}}.
 \tag{2.3}$$

21

108 See Sverchkov and Pfeffermann (2018), and Pfeffermann and Sverchkov (2019)
 109 for details.

110 3. ESTIMATION OF RESPONSE PROBABILITIES

111 Unlike the sampling probabilities, the response probabilities are generally
 112 unknown. We assume therefore a parametric model, which is allowed to depend
 113 on the outcome and the covariate values; $\Pr[R_{ij} = 1 | y_{ij}, \mathbf{x}_{ij}, i \in S, j \in S_i; \boldsymbol{\gamma}]$
 114 $= p_r(y_{ij}, \mathbf{x}_{ij}; \boldsymbol{\gamma})$, where $\boldsymbol{\gamma}$ is a vector of unknown coefficients. We assume that
 115 $p_r(y_{ij}, \mathbf{x}_{ij}; \boldsymbol{\gamma})$ is differentiable with respect to $\boldsymbol{\gamma}$ and satisfies the same mild
 116 regularity conditions as in Sverchkov and Pfeffermann (2018).

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

$$119 \sum_{(i,j) \in R} \frac{\partial \log p_r(y_{ij}, \mathbf{x}_{ij}; \boldsymbol{\gamma})}{\partial \boldsymbol{\gamma}} + \sum_{(i,k) \in R^c} \frac{\partial \log [1 - p_r(y_{ik}, \mathbf{x}_{ik}; \boldsymbol{\gamma})]}{\partial \boldsymbol{\gamma}} = 0. \quad (3.1)$$

120 In practice, the missing data are unobserved and hence the likelihood equations
 121 (3.1) are not operational. However, one may apply in this case the missing
 122 information principle:

123 **Missing Information Principle (MIP, Cepillini et al. 1955, Orchard and Woodbury,**
 124 **1972):** Let $O = \{y_{ij}, n_i, (i, j) \in R; \mathbf{x}_{ht}, h = 1, \dots, M, t = 1, \dots, N_i\}$ represent the known
 125 observed data used below. Since no observations are available for $(i, k) \in R^c$,
 126 solve instead,

$$\begin{aligned}
127 \quad & E_U \left\{ \left[\sum_{(i,j) \in R} \frac{\partial \log p_r(y_{ij}, \mathbf{x}_{ij}; \gamma)}{\partial \gamma} + \sum_{(i,k) \in R^c} \frac{\partial \log [1 - p_r(y_{ik}, \mathbf{x}_{ik}; \gamma)]}{\partial \gamma} \right] \middle| O \right\} \\
128 \quad & \stackrel{\text{by (2.3)}}{=} \sum_{(i,j) \in R} \frac{\partial \log p_r(y_{ij}, \mathbf{x}_{ij}; \gamma)}{\partial \gamma} \\
129 \quad & + \sum_{(i,k) \in R^c} E_s \left(\frac{E_{re} \left\{ [p_r^{-1}(y_{ik}, \mathbf{x}_{ik}; \gamma) - 1] \frac{\partial \log [1 - p_r(y_{ik}, \mathbf{x}_{ik}; \gamma)]}{\partial \gamma} \middle| \mathbf{x}_{ik}, u_i, (i,k) \in R \right\}}{E_{re} \{ [p_r^{-1}(y_{ik}, \mathbf{x}_{ik}; \gamma) - 1] \mid \mathbf{x}_{ik}, u_i, (i,k) \in R \}} \right) \middle| O = 0.
\end{aligned}$$

(3.2)

131 See Sverchkov (2008) and Sverchkov and Pfeffermann (2018) for derivation of
132 (3.2). In these equations, E_U, E_s, E_{re} define respectively expectations with respect
133 to the population distribution, the sample distribution and the respondents'
134 distribution. Notice that the internal expectations in the last expression are with
135 respect to the model holding for the observed data for the respondents.

136
137 *Remark 2.* When the response probabilities $p_r(y_{ij}, \mathbf{x}_{ij}; \gamma)$ depend on only \mathbf{x}_{ij} , they
138 are referred to as *propensity scores*, and the missing data are missing at random.
139 This kind of response mechanism may hold in establishment surveys, for example,
140 when the response probability is related to the known size of the establishment.
141 The estimating equations in (3.2) reduce in this case to the common log-likelihood
142 equations,

$$143 \quad \sum_{(i,j) \in R} \frac{\partial \log p_r(\mathbf{x}_{ij}; \gamma)}{\partial \gamma} + \sum_{(i,k) \in R^c} \frac{\partial \log [1 - p_r(\mathbf{x}_{ik}; \gamma)]}{\partial \gamma} = 0, \quad (3.3)$$

144 where $p_r(\mathbf{x}_{ij}; \gamma) = \Pr(R_{ij} = 1 \mid \mathbf{x}_{ij}; \gamma)$.

145 Sverchkov and Pfeffermann (2018) propose to solve the equations (3.2) by
146 maximizing the log-likelihood leading to them, i.e., maximizing,

$$\begin{aligned}
147 \quad l(\gamma) &= \sum_{(i,j) \in R} \log p_r(y_{ij}, \mathbf{x}_{ij}; \gamma) \\
148 \quad &+ \sum_{(i,k) \in R^c} E_s \left(\frac{E_{re} \{ [p_r^{-1}(y_{ik}, \mathbf{x}_{ik}; \gamma^*) - 1] \log [1 - p_r(y_{ik}, \mathbf{x}_{ik}; \gamma)] \mid \mathbf{x}_{ik}, u_i, (i,k) \in R \}}{E_{re} \{ [p_r^{-1}(y_{ik}, \mathbf{x}_{ik}; \gamma^*) - 1] \mid \mathbf{x}_{ik}, u_i, (i,k) \in R \}} \right) \Big| O \Big). \quad (3.4)
\end{aligned}$$

149 We distinguish between γ^* and γ because by (3.2), the derivatives should only
150 be taken with respect to γ .

151 We maximize the likelihood (3.4) by replacing u_i by \hat{u}_i , obtained by fitting a model
152 of the form (2.2), and dropping the external expectation- E_s . The maximization is
153 carried out iteratively, by maximizing in the $(q+1)$ iteration the expression,

$$\begin{aligned}
154 \quad &\sum_{(i,j) \in R} \log p_r(y_{ij}, \mathbf{x}_{ij}; \gamma^{(q+1)}) \\
155 \quad &+ \sum_{(i,k) \in R^c} \frac{E_{re} \{ [p_r^{-1}(y_{ik}, \mathbf{x}_{ik}; \gamma^{(q)}) - 1] \log [1 - p_r(y_{ik}, \mathbf{x}_{ik}; \gamma^{(q+1)})] \mid \mathbf{x}_{ik}, \hat{u}_i, (i,k) \in R \}}{E_{re} \{ [p_r^{-1}(y_{ik}, \mathbf{x}_{ik}; \gamma^{(q)}) - 1] \mid \mathbf{x}_{ik}, \hat{u}_i, (i,k) \in R \}}
\end{aligned} \quad (3.5)$$

156 with respect to $\gamma^{(q+1)}$. The maximization can be carried out, for example, by SAS
157 Proc NLIN. See Sverchkov (2022) and the examples following Remark 3 for
158 details.
159 details.

160 **Remark 3.** A fundamental question regarding the solution of the MIP equations is
161 the existence of a unique solution or more generally, the identifiability of the
162 response model. Riddles et al. (2016) propose a similar approach to deal with
163 NMAR nonresponse in the general context of survey sampling inference and
164 establish the following fundamental condition for the response model identifiability:
165 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
166 $\Pr(R_{ij} = 1 \mid y_{ij}, \mathbf{x}_{ij}) = \Pr(R_{ij} = 1 \mid y_{ij}, \mathbf{x}_{1ij})$. In other words, the covariates in \mathbf{x}_2 that
167 appear in the outcome model do not affect the response probabilities, given the
168 outcome and the other covariates. Variable(s) of this property may or may not exist
169 in a general set up, but interesting enough, SAE models actually contain such a
170 variable, namely, the random effects. The random effects play a fundamental role
171 in SAE models, so the outcome clearly depends on them, but it is reasonable to

172 assume that the response probabilities do not depend on the random effect, given
 173 the outcome value. In practice, the random effects are unobservable, but we
 174 estimate them and then solve the equations (3.5) by conditioning on the estimated
 175 effects. So, it is actually the estimated random effects that play the role of the
 176 covariates \mathbf{x}_2 . (Other covariates that are predictive of the outcome but not of the
 177 response might exist as well).

178 Clearly, the larger is the absolute values of the random effects, the more they affect
 179 the values of the outcome values and hence also the values of the response
 180 probabilities. In the simulation study of Sverchkov and Pfeffermann (2018), the
 181 authors study the effect of the magnitude of the variance of the random effects on
 182 the prediction of the area means. The conclusions from that study is that although
 183 the estimators of response model parameters become biased as the variance of
 184 the random effects increases, the biases are relatively very small and so are the
 185 standard deviations of the estimators. Increasing the variance of the random
 186 effects has negligible effect on the estimation of the true response probabilities
 187 and the predictors of the true small area means remain virtually unbiased in each of
 188 the areas.

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

192 **Example 1.** (Sverchkov and Pfeffermann 2018): *Mixed logistic model for the*
 193 *outcome variable.*

194 Suppose that the model fitted to the observed data of the respondents is the mixed
 195 generalized logistic model,

$$196 \quad p_y(x_{ij}, u_i) = \Pr(y_{ij} = 1 | x_{ij}, u_i, (i, j) \in R; \boldsymbol{\beta}) = \frac{\exp(\beta_0 + \beta_1 x_{ij} + u_i)}{1 + \exp(\beta_0 + \beta_1 x_{ij} + u_i)}, \quad u_i \stackrel{\text{i.i.d.}}{\sim} N(0, \sigma_u^2).$$

197 Consider a generic response model, $p_r(y_{ij}, x_{ij}; \boldsymbol{\gamma}) = \Pr[R_{ij} = 1 | y_{ij}, x_{ij}, i \in s, j \in s_i; \boldsymbol{\gamma}]$.

198 The components of (3.2) can be written in this case as,

$$199 \quad E_{re} \left\{ [p_r^{-1}(y_{ij}, x_{ij}; \boldsymbol{\gamma}) - 1] \frac{\partial \log[1 - p_r(y_{ij}, x_{ij}; \boldsymbol{\gamma})]}{\partial \boldsymbol{\gamma}} \Big| x_{ij}, u_i, (i, j) \in R \right\} =$$

$$\begin{aligned}
200 \quad & p_y(x_{ij}, u_i)[p_r^{-1}(1, x_{ij}; \gamma) - 1] \frac{\partial \log[1 - p_r(1, x_{ij}; \gamma)]}{\partial \gamma} + \\
201 \quad & [1 - p_y(x_{ij}, u_i)][p_r^{-1}(0, x_{ij}; \gamma) - 1] \frac{\partial \log[1 - p_r(0, x_{ij}; \gamma)]}{\partial \gamma}; \tag{3.6}
\end{aligned}$$

$$\begin{aligned}
202 \quad & E_{re} \{ [p_r^{-1}(y_{ij}, x_{ij}; \gamma) - 1] | x_{ij}, u_i, (i, j) \in R \} = p_y(x_{ij}, u_i)[p_r^{-1}(1, x_{ij}; \gamma) - 1] + \\
203 \quad & [1 - p_y(x_{ij}, u_i)][p_r^{-1}(0, x_{ij}; \gamma) - 1]. \tag{3.7}
\end{aligned}$$

204 The random effects u_i and the logistic probabilities $p_y(x_{ij}, u_i)$ can be estimated by
205 use of the SAS procedure PROC NLMIX.

206 **Example 2.** (Sverchkov 2022): *General continuous model.*

207 In Example 1, the outcomes follow a discrete distribution. In this section, we
208 consider continuous outcomes. The proposed algorithm consists of two parts:

209 **Part 1:** Fit (estimate) the model (2.2). The output of this part (input for Part 2)
210 contains the model parameter estimates, the estimated random effects, \hat{u}_i , and
211 for each $(i, j) \in R$, estimates of $p_y^{(l)}(\mathbf{x}_{ij}, \hat{u}_i) = P_R(a_l \leq y_{ij} < a_{l+1} | \mathbf{x}_{ij}, \hat{u}_i, (i, j) \in R)$,
212 $l = 0, \dots, L + 2$; $a_0 = -\infty$, $a_{L+2} = \infty$, $a_l = \min(y_{ij}) + (l-1) \frac{\max(y_{ij}) - \min(y_{ij})}{L}$,
213 $l = 1, \dots, L + 1$. The max and min are over all the observed values y_{ij} .

214 **Part 2:** Approximate the expectations in (3.2) similarly to (3.6) and (3.7):

$$\begin{aligned}
215 \quad & E_{re} \left\{ [p_r^{-1}(y_{ij}, \mathbf{x}_{ij}; \gamma) - 1] \frac{\partial \log[1 - p_r(y_{ij}, \mathbf{x}_{ij}; \gamma)]}{\partial \gamma} \middle| \mathbf{x}_{ij}, u_i, (i, j) \in R \right\} \cong \\
216 \quad & \sum_{l=1}^{L+1} \hat{p}_y^{(l)}(\mathbf{x}_{ij}, \hat{u}_i) [p_r^{-1}(a_l, \mathbf{x}_{ij}; \gamma) - 1] \frac{\partial \log[1 - p_r(a_l, \mathbf{x}_{ij}; \gamma)]}{\partial \gamma}, \tag{3.8}
\end{aligned}$$

$$217 \quad E_{re} \left\{ [p_r^{-1}(y_{ij}, \mathbf{x}_{ij}; \gamma) - 1] \middle| \mathbf{x}_{ij}, u_i, (i, j) \in R \right\} \cong \sum_{l=1}^{L+1} \hat{p}_y^{(l)}(\mathbf{x}_{ij}, \hat{u}_i) [p_r^{-1}(a_l, \mathbf{x}_{ij}; \gamma) - 1], \tag{3.9}$$

218 where $p_r(a_l, \mathbf{x}_{ij}; \gamma) = \Pr[R_{ij} = 1 | y_{ij} = a_l, \mathbf{x}_{ij}, i \in s, j \in s_i; \gamma]$.

219 Substitute (3.8) and (3.9) into (3.2) and estimate γ by iteratively maximizing (3.5).

220

4. SELECTION OF A RESPONSE MODEL

221 There is no direct way to test the appropriateness of a chosen response model
222 because the outcome values, which are part of the model, are unknown for the
223 nonresponding units. If the model for the outcomes before nonresponse was
224 known, one could derive the model holding for the observed outcomes based on
225 this model and the model assumed for the responding units, and test the resulting
226 model by use of standard tests that compare the cumulative hypothesized
227 distribution of the observed data with the corresponding empirical distribution,
228 and/or by testing moments of the assumed model. See, e.g., Pfeffermann and
229 Landsman (2011) and Pfeffermann and Sikov (2011). However, in the approach
230 described in Section 3, we start with a model fitted to the observed outcomes,
231 which does not include the response model and therefore, we cannot use a similar
232 strategy.

233 When following the approach proposed in Section 3, the likelihood (3.4) suggests
234 at least two procedures for the selection of the response model in SAE under
235 NMAR nonresponse. **1-** Compare different models based on information criteria
236 such as the Akaike information criterion, $AIC = -2l(\gamma) + 2 \dim(\gamma)$, or Schwarz
237 information criterion, $BIC = -2l(\gamma) + \dim(\gamma) \log(n)$, $n = \sum_{i \in S} n_i$; **2-** test a saturated
238 versus a nested model based on the likelihood ratio test. In Section 5 we illustrate
239 via a simulation study how the likelihood (3.4) can be used for the application of
240 these selection procedures.

241

5. SIMULATION STUDY

242 5.1 Simulation set-up

243 We start by defining the sample model before nonresponse because as stated in
244 Section 4, our approach for estimating the response model is based on fitting a
245 model to the observed outcomes, which does not include the response model. For
246 convenience, we assume noninformative sampling of areas and within the areas,
247 such that the sample model before nonresponse is the same as the population
248 model. Note that although the sampling design defines the observed model (2.2),

249 once this model is estimated, the sampling design does not affect the estimation
250 of the response probabilities in section 3.

251 The simulation study consists of the following steps:

252 Generate auxiliary values x_{ij} , $i = 1, \dots, 100$, $j = 1, \dots, 20$ from a Uniform(0,2)
253 distribution. Next, generate sample values from the small area model,

$$254 \quad y_{ij} | x_{ij}, u_i \sim N(x_{ij} + u_i, 1), \quad i = 1, \dots, 100, \quad j = 1, \dots, 20; \quad u_i \sim N(0, 1). \quad (5.1)$$

255 Consider three unit response models (no selection of areas):

$$256 \quad p_r^{(1)}(y_{ij}, x_{ij}) = \text{logit}(-x_{ij} / 2 + 2y_{ij}),$$

$$257 \quad p_r^{(2)}(y_{ij}, x_{ij}) = \text{logit}(-x_{ij} / 2 + 2y_{ij} - 0.3y_{ij}^2),$$

$$258 \quad p_r^{(3)}(y_{ij}, x_{ij}) = \text{logit}(1.5x_{ij}).$$

259 Select 3 sets of respondents:

260 R1 uses Poisson sampling, independently between the units with response
261 probabilities $p_r^{(1)}(y_{ij}, x_{ij})$,

262 R2 is the same as R1 but with response probabilities $p_r^{(2)}(y_{ij}, x_{ij})$,

263 R3 is the same as R1 but with response probabilities $p_r^{(3)}(y_{ij}, x_{ij})$.

264 The 3 response probabilities yield similar response rates of 65 - 75 per cent.

265 The working model for the observed data for the responding units is,

$$266 \quad y_{ij} | x_{ij}, u_i \sim N(\theta_0 + \theta_1 x_{ij} + u_i, \theta_2), \quad i = 1 \dots 100, \quad j = 1 \dots 20; \quad u_i \sim N(0, \sigma_u^2). \quad (5.2)$$

267 *Remark 4.* The working model (5.2) is correct for the observed sample R3 that
268 corresponds to MAR nonresponse, but not for R1 and R2, under which the
269 nonresponse is NMAR.

270 Define three working response models:

$$271 \quad \text{M1: } p_r(y_{ij}, x_{ij}; \gamma^1) = \text{logit}(\gamma_0 + \gamma_1 x_{ij} + \gamma_2 y_{ij}),$$

$$272 \quad \text{M2: } p_r(y_{ij}, x_{ij}; \gamma^2) = \text{logit}(\gamma_0 + \gamma_1 x_{ij} + \gamma_2 y_{ij} + \gamma_3 y_{ij}^2),$$

273 M3: $p_r(y_{ij}, x_{ij}; \gamma^3) = \text{logit}(\gamma_0 + \gamma_1 x_{ij})$,

274 Note that $p_r(y_{ij}, x_{ij}; \gamma^3)$ is nested in $p_r(y_{ij}, x_{ij}; \gamma^1)$ and $p_r(y_{ij}, x_{ij}; \gamma^2)$, and

275 $p_r(y_{ij}, x_{ij}; \gamma^1)$ is nested in $p_r(y_{ij}, x_{ij}; \gamma^2)$. The response probability $p_r(y_{ij}, x_{ij}; \gamma^3)$

276 defines MAR nonresponse and hence, can be estimated by solving (3.3).

277 Estimate the unknown parameters $\theta_0, \theta_1, \theta_2$ in (5.2) by SAS Proc NMIX, and then

278 estimate γ by maximizing (3.4), as described in Section 3. The maximization was

279 carried out by use of SAS Proc NLIN under the following 9 scenarios, as defined

280 by the true response model and the assumed working response model:

281 S1: R1 set of respondents, M1 working response model.

282 S2: R1 set of respondents, M2 working response model.

283 S3: R1 set of respondents, M3 working response model.

284 S4: R2 set of respondents, M1 working response model.

285 S5: R2 set of respondents, M2 working response model.

286 S6: R2 set of respondents, M3 working response model.

287 S7: R3 set of respondents, M1 working response model.

288 S8: R3 set of respondents, M2 working response model.

289 S9: R3 set of respondents, M3 working response model.

290 Select the response model based on:

291 **1-** The Likelihood Ratio Test (LRT); test a saturated model $[I(\gamma^{**})]$ against a nested

292 model $[I(\gamma^*)]$, assuming the χ^2 distribution under the null hypothesis H_0 that the

293 nested model with a smaller number of parameters is correct. The test statistic is

294 $\lambda_{LRT} = -2[l(\hat{\gamma}^*) - l(\hat{\gamma}^{**})] \sim \chi^2_{[\dim(\gamma^{**}) - \dim(\gamma^*)]}$. Reject H_0 at the $\alpha = .05$ level.

295 **2 -** AIC selection criterion: compare the values of the AIC as obtained for the
296 corresponding two models;

297 **3 -** BIC selection criterion: compare the values of the BIC as obtained for the
298 corresponding two models.

299 Repeat the whole process independently 500 times.

300 5.2 Results

301 S1 Vs S2 (R1 – set of respondents, M1 – correct model, M2 – saturated model).

302 Note that although M2 is a saturated model, it is also correct but with an additional
303 term. The LRT selects the model M1 in 368 out of the 500 simulations. AIC selects
304 M1 in 305 out of 500 simulations, BIC selects M1 in 324 simulations.

305 S1 Vs S3 (R1 – set of respondents, M1 – correct model, M3 – incorrect nested

306 model). The LRT selects the correct model M1 in 500 out of the 500 simulations.

307 AIC and BIC likewise select M1 in all the 500 simulations.

308 S4 Vs S5 (R2 – set of respondents, M2 – correct model, M1 – incorrect nested

309 model). The LRT selects the correct model M2 in 433 out of the 500 simulations.

310 AIC and BIC select M2 in 483 simulations.

311 S4 Vs S6 (R2 - set of respondents, M2 – correct model, M3 – incorrect nested

312 model). The LRT selects the correct model M2 in 490 out of the 500 simulations.

313 AIC and BIC select the correct model in all the simulations.

314 S7 Vs S8 (R3 - set of respondents, M3 – correct model, M1 – also correct but a

315 saturated model). The LRT selects the model M3 in 241 out of 500 simulations.

316 AIC selects M3 in 225 out of 500 simulations; BIC selects M3 in 420 simulations.

317 S7 Vs S9 (R3 - set of respondents, M3 – correct model, M2 – also correct but a

318 saturated model). The LRT selects the M3 model in 241 out of 500 simulations.

319 AIC selects M3 in 361 out of 500 simulations, BIC selects M3 in 450 simulations.

320 Note that when R3 is the set of respondents and M3 is the correct model, M1 and

321 M2 also produce correct estimates of the response probabilities, although with

322 additional estimated parameters. Thus, the fact that the LRT test and the AIC

323 select the M3 model in about half of the simulations is not surprising. The use of

324 the BIC criterion performs better in these cases.

325 The results so far are summarized in table 1.

326 **Table 1.** Percentages out of 500 simulations in which each of the three selection
 327 procedures selected the correct model, for different combinations of correct
 328 (rows) and working (columns) response probability models.

	M1			M2			M3		
	LRT	AIC	BIC	LRT	AIC	BIC	LRT	AIC	BIC
R1, M1 correct	---	---	---	73.6	61.0	64.8	100	100	100
R2, M2 correct	86.6	96.6	96.6	---	---	---	98	100	100
R3, M3 correct	48.2	45	82	48.2	72.2	90	---	---	---

329

17

330 Finally, we consider the case where a working model is incorrect but might be a
 331 good approximation of the correct model: let R1 be the set of respondents such
 332 that M1 is the correct working model. Let M4 be the following working model:
 333 $p_r(y_{ij}, x_{ij}; \Upsilon^4) = \text{logit}(\gamma_0 + \gamma_1 x_{ij} + \gamma_2 y_{ij}^2 + \gamma_3 y_{ij}^3)$. Compare M4 with M1 (correct
 334 model). In this case, AIC selects the correct M1 model in 430 out of the 500
 335 simulations and BIC selects the correct model in 431 simulations.

336 Sverchkov (2013) suggested testing whether the response is NMAR or MAR by
 337 testing the significance of the corresponding estimated coefficients in the saturated
 338 response model. We applied this idea by testing the significance of the estimated
 339 coefficients $\hat{\gamma}_2$ and $\hat{\gamma}_3$ under the response models $p_r^{(1)}(y_{ij}, x_{ij})$ and $p_r^{(2)}(y_{ij}, x_{ij})$,
 340 (both assume NMAR nonresponse), when in fact the true response model is
 341 $p_r^{(3)}(y_{ij}, x_{ij})$ (MAR) or $p_r^{(1)}(y_{ij}, x_{ij})$ (NMAR), using the standard t-tests. (SAS Proc
 342 NLIN provides standard errors of the estimated coefficients.)

343 We considered two samples of respondents, R3 and R1. For R3, we found that
 344 when testing the working response model $p_r^{(1)}(y_{ij}, x_{ij})$, in 432 out of the 500
 345 simulations, $\hat{\gamma}_2$ was not significant at the 0.05 level. When testing the working
 346 response model $p_r^{(2)}(y_{ij}, x_{ij})$, in 350 out of the 500 simulations, $\hat{\gamma}_2$ was not

347 significant at the 0.05 level, and in 398 simulations $\hat{\gamma}_3$ was not significant. Recall
348 that for the respondents' sample R3, the working SAE model (5.2) for the observed
349 outcomes is correct since the response is MAR.

12
350 For the respondents' sample R1, the response model $p_r^{(1)}(y_{ij}, x_{ij})$ is correct and in
351 all the 500 simulations, the estimator $\hat{\gamma}_2$ was found significant. However, when
352 testing the response model $p_r^{(2)}(y_{ij}, x_{ij})$, in 388 simulations the estimator $\hat{\gamma}_3$ was
353 significant, even at the 0.01 level, and $\hat{\gamma}_2$ was significant in 498 simulations. This
354 result might be explained by the fact that the working outcome model (5.2) is not
355 correct when the response model is NMAR and thus, the likelihood (3.5), which
356 conditions on the estimated random effects for the estimation of the γ coefficients
357 is incorrect.

358
359

6. SUMMARY

2
360 In this paper we investigate the use of the likelihood function of the observed
361 respondents' data for selecting an appropriate response model under possible
362 NMAR nonresponse. For estimating the hypothesized model, we applied the
363 missing information principle. Despite of what seems to be a rather complex
364 estimation process, we find in our simulation study that the AIC and BIC
365 information criteria and the LRT test, when applicable, perform well for model
366 selection. Clearly, the use of other likelihood-based tests and selection criteria
367 should be investigated as well.

368
369

REFERENCES

370
371 Cepillini, R., Siniscalco, M., and Smith, C.A.B. (1955). The estimation of gene
372 frequencies in a random mating population. *Annals of Human Genetics*, **20**, 97-
373 115.
374 Orchard, T., and Woodbury, M.A. (1972). A missing information principle: theory
375 and application. *Proceedings of the 6th Berkeley Symposium on Mathematical*
376 *Statistics and Probability*, **1**, 697-715.

- 377 Pfeffermann, D. (2013). New Important Developments in Small Area Estimation.
378 *Statistical Science*, **28**, 40-68.
- 379 Pfeffermann, D., and V. Landsman (2011), "Are Private Schools Better than Public
380 Schools? Appraisal for Ireland by Methods for Observational Studies," *Annals of*
381 *Applied Statistics*, 5, 1726–1751.
- 382 Pfeffermann, D. and Sikov N. (2011). Imputation and estimation under
383 nonignorable nonresponse in household surveys with missing covariate
384 information. *Journal of Official Statistics*, **27**, 181–209.
- 385 Pfeffermann, D., and Sverchkov, M. (2007). Small-Area Estimation under
386 Informative Probability Sampling of Areas and Within Selected Areas. *Journal of*
387 *the American Statistical Association*, **102**, 1427-1439.
- 388 Pfeffermann, D., and Sverchkov, M. (2009). Inference under Informative Sampling.
389 In: *Handbook of Statistics 29B; Sample Surveys: Inference and Analysis*. Eds. D.
390 Pfeffermann and C.R. Rao. North Holland. pp. 455-487.
- 391 Pfeffermann, D. and Sverchkov, M. (2019). Multivariate small area estimation
392 under nonignorable nonresponse, *Statistical Theory and Related Fields*, **3**, pp.
393 213-223.
- 394 Rao, J.N.K., and Molina, I. (2015), *Small Area Estimation*, 2nd Edition, Wiley.
- 395 Riddles, K.M., Kim, J.K. and Im, J. (2016). A propensity-score adjustment method
396 for nonignorable nonresponse. *Journal of Survey Statistics and Methodology*, **4**,
397 215-245.
- 398 Särndal, C.E. and Swensson B. (1987). A general view of estimation for two
399 phases of selection with applications to two-phase sampling and nonresponse.
400 *International Statistical Review*, **55**, 279-294.
- 401 Sverchkov, M. (2008). A new approach to estimation of response probabilities
402 when missing data are not missing at random. *Joint Statistical Meetings*,
403 *Proceedings of the Section on Survey Research Methods*, 867-874.

- 404 Sverchkov, M. (2013). Is it MAR or NMAR? *2013 JSM Meetings, Proceedings of*
405 *the Section on Survey Methods Research*, pp. 2307-2311
- 406 Sverchkov, M. (2022). An Algorithm for Small Area Estimation under Not Missing
407 At Random Non-response. *Joint Statistical Meetings, Proceedings of the Section*
408 *on Survey Research Methods*, pp. 1735-1745.
- 409 Sverchkov, M., and Pfeffermann, D. (2004). Prediction of Finite Population Totals
410 Based on the Sample Distribution. *Survey Methodology*, **30**, 79-92.
- 411 Sverchkov, M. and Pfeffermann, D. (2018). Small area estimation under
412 informative sampling and not missing at random non-response. *Journal of Royal*
413 *Statistical Society, ser. A*, 181, Part 4, pp. 981–1008.

Peer Review

RESPONSE MODEL SELECTION IN SMALL AREA ESTIMATION UNDER NOT MISSING AT RANDOM NONRESPONSE

ORIGINALITY REPORT

36%

SIMILARITY INDEX

PRIMARY SOURCES

- 1** eprints.soton.ac.uk
Internet 914 words — 19%
- 2** www.isi2023.org
Internet 304 words — 6%
- 3** www.bls.gov
Internet 114 words — 2%
- 4** Michael Sverchkov, Danny Pfeffermann. "Small area estimation under informative sampling and not missing at random non - response", Journal of the Royal Statistical Society: Series A (Statistics in Society), 2018
Crossref 77 words — 2%
- 5** Michael Sverchkov, Danny Pfeffermann. "Small area estimation under informative sampling and not missing at random non-response", Journal of the Royal Statistical Society: Series A (Statistics in Society), 2018
Crossref 75 words — 2%
- 6** washingtonstatisticalsociety.org
Internet 36 words — 1%
- 7** www.tandfonline.com
Internet 36 words — 1%

8	xblk.ecnu.edu.cn Internet	22 words — < 1%
9	2013.isiproceedings.org Internet	17 words — < 1%
10	epdf.pub Internet	17 words — < 1%
11	Chiara Mussida, Luca Zanin. "I found a better job opportunity! Voluntary job mobility of employees and temporary contracts before and after the great recession in France, Italy and Spain", Empirical Economics, 2019 Crossref	15 words — < 1%
12	Danny Pfeffermann, Michael Sverchkov. "Multivariate small area estimation under nonignorable nonresponse", Statistical Theory and Related Fields, 2019 Crossref	12 words — < 1%
13	Anna Sikov. "A Brief Review of Approaches to Non-ignorable Non-response", International Statistical Review, 2018 Crossref	10 words — < 1%
14	www2.math.umd.edu Internet	10 words — < 1%
15	International Encyclopedia of Statistical Science, 2011. Crossref	9 words — < 1%
16	Wenchao Ma, Jimmy de la Torre. "A sequential cognitive diagnosis model for polytomous responses", British Journal of Mathematical and Statistical Psychology, 2016	9 words — < 1%

-
- 17 "Analysis of Survey Data", Wiley, 2003
Crossref 8 words — < 1%
-
- 18 Abdulhakeem A. H. Eideh. "Fitting Variance Components Model and Fixed Effects Model for One-Way Analysis of Variance to Complex Survey Data", Communications in Statistics - Theory and Methods, 2012
Crossref 8 words — < 1%
-
- 19 J.-F. Beaumont. "A new approach to weighting and inference in sample surveys", Biometrika, 09/01/2008
Crossref 8 words — < 1%
-
- 20 sit.stat.gov.pl
Internet 8 words — < 1%
-
- 21 www.statcan.gc.ca
Internet 8 words — < 1%
-
- 22 Daniela Marella, Danny Pfeffermann. "Accounting for Non - ignorable Sampling and Non - response in Statistical Matching", International Statistical Review, 2022
Crossref 7 words — < 1%
-
- 23 Danny Pfeffermann. " Bayes-based Non-Bayesian Inference on Finite Populations from Non-representative Samples: A Unified Approach *Based on S. N. Roy Memorial Lecture in the symposium. ", Calcutta Statistical Association Bulletin, 2017
Crossref 6 words — < 1%
-
- 24 Feder, Moshe, Pfeffermann, Danny. "Statistical inference under non-ignorable sampling and non-

response. An empirical likelihood approach", 'University of Southampton', 2015

Internet

EXCLUDE QUOTES ON

EXCLUDE BIBLIOGRAPHY ON

EXCLUDE SOURCES OFF

EXCLUDE MATCHES OFF