Use of Nonprobability Samples for Official Statistics, State of the Art

Danny Pfeffermann<sup>(1)</sup> and Michael Sverchkov<sup>(2)</sup>

Abstract

Tightened budgets, continuing decrease in response rates in traditional probability surveys and

increasing pressure by users for more timely data, has stimulated research for the use of

nonprobability sample data, such as administrative records, web scraping, mobile phone data and

voluntary internet surveys, for inference on finite population parameters like means and totals.

These data are often easier, faster and cheaper to collect than traditional probability samples.

However, a major concern with the use of this kind of data is their nonrepresentativeness due to

possible selection bias, which if not accounted for properly, could bias the inference. In this article,

we review and discuss methods considered in the literature to deal with this problem and propose

new methods, distinguishing between methods based on integration of the nonprobability sample

with an appropriate probability sample, and methods that base the inference solely on the

nonprobability sample. Empirical illustrations, based on simulated data are provided.

Key words: Empirical likelihood, Probability and nonprobability samples, Sample

integration, Selection bias.

(1) Hebrew University of Jerusalem, Israel & University of Southampton, UK.

Email: msdanny@soton.ac.uk

(2) Bureau of Labor Statistics, Washington DC, USA.

Email: Sverchkov.Michael@bls.gov

1

## 1. INTRODUCTION

Tightened budgets, continuing decrease in response rates, due in part by increased response burden in traditional probability surveys and privacy concerns, and increasing pressure by users for more timely data, has prompted research into the use of nonprobability sample data, such as administrative records, web scraping, mobile telephone data, online panels and voluntary internet surveys for inference on finite population characteristics. These data are often easier, faster and cheaper to collect than are traditional probability samples. However, a major concern with the use of this kind of data is their possible nonrepresentativeness, due to possible selection bias, which if not accounted for properly, could bias the inference. For example, house sales advertised on the internet do not represent properly all house sales. Web scraping for job vacancies does not represent all job vacancies. Data from social media do not generally represent the general public. All these examples can be considered as 'big data', but nonprobability samples do not need to be big. Baker et al.. (2013), Keiding and Louis (2016) and Elliott and Valliant (2017) discuss other potential problems with the use of nonprobability samples for inference on finite population parameters.

The basic definition of a probability sample is that every unit in the population has a positive probability of being included in the sample. Inference under the traditional randomization (design-based) distribution over all possible sample selections from a fixed target population requires that the first-order sample selection probabilities of the sampled units are known. The use of standard variance estimation procedures requires that the joint sample selection probabilities of the sampled units are also known, but these can be calculated to a desired approximation by repeated sampling from the sampling frame. (Not usually available to analysts outside National Statistical Offices-NSOs.)

By definition, nonprobability samples are not selected by use of a probability sampling scheme so no selection probabilities exist. The question arising therefore is how to draw inference from such samples, regarding the population, which they are supposed to represent. In this article, we restrict our attention to inference about target population parameters such as totals or means (proportions), which are the most common target parameters in official statistics, often published in tables.

We mention in this respect that many survey statisticians claim that traditional probability samples should be replaced by external records. Citro (2014) states that "official statistical offices need to move from the probability sample survey paradigm for the past 75 years

to a mixed mode data source paradigm for the future." Clearly, if the nonprobability sample data are timely, accurate, with good coverage and contain all the required information, there is no reason to select a corresponding probability sample.

However, this is seldom the case. Israel's population register covers all the population residing in Israel, but about 15% of the home addresses are wrong. Tax records of businesses are often obtained with a delay of up to 2 years. No administrative data are available on opinions, sentiments, detailed expenditures, and many other variables of interest. We mention also in this regard that government and private agencies are often reluctant to transfer data to NSOs, claiming data protection issues. Furthermore, the desired information is often contained in more than one file, requiring matching them, which is problematic if personal identifiers are unknown. (Requires probabilistic algorithms based on information in all the records.) Coverage of records might be different and may not apply to same time periods. Definitions and accuracy of information may differ between records. Finally, matching of different administrative data could magnify problems of data protection.

Methods considered in the literature to deal with possible non-representativeness of nonprobability (NP) samples can be divided into two classes:

- 1- Integration of the NP sample with an appropriate probability sample (PS),
- **2-** Consideration of the NP sample on its own volume of the NP sample on its own volume of the NP sample on its own volume.

REMARK 1. The methods considered in this article for inference from NP samples alone assume known population means or totals of some of the survey values, which are used for enhancing the inference. However, no detailed probability sample data are used.

In Section 2, we review several methods proposed in the literature for integration of a NP sample with an appropriate PS sample. We also present a new method. Section 3 reviews methods proposed for adjusting for selection bias of a NP sample without integration with a PS ple. In Section 4, we propose a new method for inference from a NP sample without integration. Section 5 contains simulation results illustrating the performance of our proposed method. We conclude with some summary remarks in Section 6.

#### 2. INTEGRATION OF NONPROBABILITY AND PROBABILITY SAMPLES

One of the earliest articles on this topic is by Lee (2006). The author proposes to create a pooled sample  $S_P = S_{PS} \cup S_{NP}$  from the probability sample  $S_{PS}$  and the nonprobability sample  $S_{NP}$ , assuming implicitly that the two samples do not overlap, and models the selection probability to the nonprobability sample. The  $S_{NP}$  sample is treated as a "treatment sample" in observational studies, and the  $S_{PS}$  sample is treated as the "control sample". It is assumed that every unit in the population has a positive probability to be in the  $S_{NP}$  sample, estimated by use of propensity scores,  $e(\mathbf{x}_j) = \Pr(j \in S_{NP} \mid \mathbf{x}_j; j = 1,...,n)$ , where n is the size of  $S_P$  and the  $\mathbf{x}$ -variables are assumed to be measured in both samples.

Next, the sample  $S_P$  is divided into C classes based on the ascending values of the estimated propensity scores. An adjustment factor  $f_C$  is computed for every class c as,

$$f_c = \frac{\sum_{k \in S_{PS}^c} d_{k,PS} / \sum_{k \in S_{PS}} d_{k,PS}}{\sum_{j \in S_{NP}^c} d_{j,NP} / \sum_{j \in S_{NP}} d_{j,NP}},$$
(2.1)

where  $d_{k,PS}$  and  $d_{j,NP}$  are some base weights. An adjusted weight  $d_{j,NP}^A = f_c d_{j,NP}$  is computed for every unit  $j \in S_{NP}$ .

The estimator of the target population total 
$$Y = \sum_{i \in U} Y_i$$
 is,  $\hat{Y}_{S_{NP}} = \sum_{c} \sum_{j \in S_{NP}^c} d_{j,NP}^A y_j$ .

The use of this procedure for data integration requires the existence of **x**-variables such that the assignment to  $S_{NP}$  and the target y-variable are independent given **x**,  $\Pr(j \in S_{NP} \mid \mathbf{x}_j, y_j; j \in S_P) = \Pr(j \in S_{NP} \mid \mathbf{x}_j; j \in S_P)$ . This is a limiting assumption. An extensive empirical study revealed that the use of this approach decreases (but not eliminates) the bias of inference from the  $S_{NP}$  sample, but increases the variance. See also Beaumont (2020).

Kott and Ridenhour (2024) likewise consider the use of a pooled sample  $S_P = S_{PS} \cup S_{NP}$  for inference from the nonprobability sample. The authors model the  $S_{NP}$  selection probabilities by a logistic model with covariates  $\mathbf{z}_k$  measured in both samples and for which the true population means  $T_Z$  are known or estimated from the  $S_{PS}$  sample, which

are used for calibration. The estimating equation is  $\sum_{k \in S_{NP}} [1 + \exp(\mathbf{z}_k' \mathbf{g})] \mathbf{z}_k = T_{\mathbf{Z}}(\hat{T}_{\mathbf{Z}})$ . This defines new weights  $w_k = \pi_k^{-1}[1 + \exp(\mathbf{z}_k' \hat{\mathbf{g}})]$  used for inference from the  $S_{NP}$  sample, where  $\pi_k = \Pr(k \in S_{PS})$ . When the  $S_{PS}$  sample is subject to nonresponse, the weights  $d_k = \pi_k^{-1}$  are adjusted to account for the nonresponse.

Rivers (2007) considers the case where  $\mathbf{x}$  and  $\mathbf{y}$  are measured in the  $S_{NP}$  sample but only  $\mathbf{x}$  is measured in the  $S_{PS}$  sample. The author proposes to deal with the non-representativeness of the  $S_{NP}$  sample by matching to every unit  $i \in S_{PS}$  an element k from  $S_{NP}$ , with similar values of auxiliary (matching) variables  $\mathbf{x}$ .

Denote by  $\mathbf{x}_i, i=1,...,n$ , the  $\mathbf{x}$ -vectors in  $S_{PS}$  and by  $\tilde{\mathbf{x}}_j$  the vectors in  $S_{NP}$ . The unit  $k \in S_{NP}$  satisfying  $|\tilde{\mathbf{x}}_k - \mathbf{x}_i| \leq |\tilde{\mathbf{x}}_j - \mathbf{x}_i| \ \forall j \in S_{NP}$  is chosen as the matched element for unit  $i \in S_{PS}$ , where  $|\cdot|$  is an appropriate distance. Selecting a matching element for every unit  $i \in S_{PS}$  defines a matched sample  $S_M$  of size n with y-values from the  $S_{NP}$  sample.

The proposed estimator of the population total Y is  $\hat{Y}_{SM} = \sum_{k \in S_M} w_k \tilde{y}_k$ , where  $w_k = (1/\pi_k); \pi_k = \Pr(k \in S_{PS})$  and  $\{\tilde{y}_k\}$  are the y-values measured in  $S_{NP}$ , not measured in  $S_{PS}$ . The author establishes regularity conditions under which for a scalar continuous matching variable, as  $n \to \infty$ ,  $n_{NP} \to \infty$  and  $n/n_{NP} \to 0$ ,  $(n_{NP})$  is the size of  $S_{NP}$ ,  $n^{-0.5}(\hat{Y}_{SM} - Y)/N$  converges to a normal distribution with mean zero, where N is the population size.

REMARK 2. Rather than matching one record, one can match k nearest records and select at random the matched record out of the k records, known as the kNN method. See, e.g., Conti et al.. (2008). Alternatively, a weighted mean of the y- values of the nearest records can be used for matching.

REMARK 3. The method requires a PS sample with similar  $\mathbf{X}$  values in  $S_{NP}$  and  $S_{PS}$ . It also assumes that the relationship between  $\mathbf{y}$  and  $\mathbf{x}$  in the two samples is similar or formally, that  $f_{S_{NP}}(y_i/\mathbf{x}_i) = f_{S_{PS}}(y_i/\mathbf{x}_i)$ , where  $f_{S_{NP}}(y_i/\mathbf{x}_i) = f(y_i|\mathbf{x}_i, i \in S_{NP})$  and  $f_{S_{PS}}(y_i/\mathbf{x}_i) = f(y_i|\mathbf{x}_i, i \in S_{PS})$ .

Kim & Wang (2019) propose the following procedure of integrating the data in the  $S_{PS}$  and  $S_{NP}$  samples. The authors assume that membership of the  $S_{PS}$  elements in  $S_{NP}$  is known. Let  $\delta_i = 1(0)$  if  $i \in S_{NP}$  ( $i \notin S_{NP}$ ). The  $S_{PS}$  data contains therefore the values  $\{(\mathbf{x}_i, \delta_i); i = 1, ..., n\}$ . The procedure consists of the following step:

**1-** Model  $p_i(\gamma) = \Pr(\delta_i = 1 \mid \mathbf{x}_i; \gamma)$  by use of the  $S_{ps}$  data and estimate  $\gamma$  by maximizing the "pseudo likelihood"  $l(\gamma) = \sum_{i \in S_{ps}} w_i \{ \delta_i \log p_i(\gamma) + (1 - \delta_i) \log[1 - p_i(\gamma)] \}$ .

**2-** Estimate the population total *Y* as,

$$\hat{Y}_{S_{NP}}(1) = \sum_{i \in S_{NP}} p_i^{-1}(\hat{\gamma}) y_i \text{ or } \hat{Y}_{S_{NP}}(2) = N \sum_{i \in S_{NP}} p_i^{-1}(\hat{\gamma}) y_i / \sum_{i \in S_{NP}} p_i^{-1}(\hat{\gamma})$$
(2.2)

when N is known.

The authors consider also a doubly robust estimator under the assumption of a population regression model. Consistent variance estimators are developed.

REMARK 4. The use of this method assumes that the sampling mechanism to  $S_{NP}$  is ignorable after controlling for the covariates, i.e.  $\Pr(i \in S_{NP} \mid \mathbf{x}_i, y_i) = \Pr(i \in S_{NP} \mid \mathbf{x}_i)$ , known as noninformative sampling. In addition, the assumption that membership of the  $S_{PS}$  elements in  $S_{NP}$  is known, may not hold in practice.

Chen et al. (2020) likewise assume noninformative sampling after controlling for the covariates and assume a prediction model  $y_i = m(\mathbf{x}_i) + \varepsilon_i; i = 1,...,N$  for the population units and a selection model  $\pi_i^{S_{NP}} = \Pr(i \in S_{NP} \mid \mathbf{x}_i; \gamma)$ . For the case where  $m(\mathbf{x}_i) = \mathbf{x}_i' \boldsymbol{\beta}$ , the resulting estimator of the population mean is,  $\hat{Y}_{REG} = \sum_{i \in S_{NP}} d_i^{S_{NP}} \mathbf{x}_i' \hat{\boldsymbol{\beta}}$ , where  $d_i^{S_{NP}} = 1/\pi_i^{S_{NP}}$  and  $\hat{\boldsymbol{\beta}}$  is estimated from the  $S_{NP}$  sample. This estimator is unbiased for Y.

In practice, the sample selection model is unknown, and the authors assume a parametric model  $\pi_i^{S_{NP}} = \pi(\mathbf{x}_i; \gamma)$ , which is estimated by maximizing the pseudo likelihood

$$l^*(\boldsymbol{\gamma}) = \sum_{i \in S_{NP}} \log\left[\frac{\pi(\mathbf{x}_i, \boldsymbol{\gamma})}{1 - \pi(\mathbf{x}_i, \boldsymbol{\gamma})}\right] + \sum_{i \in S_{PS}} w_i \log\left[1 - \pi(\mathbf{x}_i, \boldsymbol{\gamma})\right], \tag{2.3}$$

where  $w_{i}=1/\pi_{i}$  are the sampling weights in  $S_{PS}$  . The resulting estimator of Y is

$$\hat{Y}_{IPW} = \sum_{i \in S_{ND}} [y_i / \pi(\mathbf{x}_i.\hat{\gamma})]. \tag{2.4}$$

When estimating the mean  $\overline{Y}=(Y/N)$ , the estimator (2.4) can be divided by N when it is known or by  $\hat{N}_{S_{NP}}=\sum_{i\in S_{NP}}[1/\pi(x_i,\hat{\gamma})]$  when N is unknown. The authors prove that for the case of a logistic selection model, both estimators have an error of order  $O_P(n_{S_{NP}}^{-1/2})$ . Variance estimators are also developed, correct to order  $o(n_{S_{NP}}^{-1})$ .

REMARK 5. In a rejoinder to comments on an article by Beaumont et al. (2024a) (see below), Beaumont et al. (2024b) state that the use of the likelihood (2.3) is not efficient because the second term only uses the  $S_{PS}$  data and ignores relevant  $S_{NP}$  auxiliary data. The authors propose an improved estimator of  $\gamma$  and a sample likelihood approach that properly accounts for an overlap between the two samples, when it can be identified.

Chen et al. (2020) also consider a doubly robust estimator, defined as

$$\hat{Y}_{DR} = \sum_{i \in S_{NP}} \left[ 1 / \pi(\mathbf{x}_i.\hat{\boldsymbol{\gamma}}) \right] \left[ y_i - m_i(\mathbf{x}_i.\hat{\boldsymbol{\beta}}) \right] + \sum_{i \in S_{PS}} w_i m_i(\mathbf{x}_i.\hat{\boldsymbol{\beta}}), \qquad (2.5)$$

where  $m_i(\mathbf{x}_i, \mathbf{\beta})$  is an assumed population regression model. The estimator  $\hat{Y}_{DR}$  is shown to be consistent for Y, even if the population model or the sample selection model is misspecified. Variance estimators correct to order  $o(n_{S_{NP}}^{-1})$  are derived.

Chen et al. (2022) consider the use of the pseudo empirical likelihood for inference from nonprobability samples, defined as  $l_{PEL}(\mathbf{p}) = \sum_{i \in S_{NP}} d_i^{S_{NP}} \log(p_i)$ , where the  $p_i$ s are the EL probabilities,  $d_i^{S_{NP}} = [(1/\pi(\mathbf{x}_i.\hat{\mathbf{\gamma}})]/\hat{N}_{SNP}$  and  $\hat{N}_{S_{NP}} = \sum_{j \in S_{NP}} [(1/\pi(\mathbf{x}_i:\hat{\mathbf{\gamma}})]$ . The parameters  $\gamma$  are estimated using the likelihood (2.3) and are considered fixed in the likelihood  $l_{PEL}(\mathbf{p})$ . Maximization of the likelihood under the constraint  $\sum_{i \in S_{ND}} p_i = 1$  yields  $\hat{p}_i = d_i^{S_{NP}}$ .

The authors also consider a doubly robust estimator obtained by adding the calibration constraint  $\sum_{i \in S_{NP}} p_i[m_i(\mathbf{x}_i; \hat{\boldsymbol{\beta}})] = \overline{m}_{S_{PS}}$  where  $\overline{m}_{S_{PS}} = \hat{N}_{S_{PS}}^{-1} \sum_{i \in S_{PS}} w_i \, m_i(\mathbf{x}_i; \hat{\boldsymbol{\beta}})$ ,  $\hat{N}_{S_{PS}} = \sum_{i \in S_{PS}} w_i$ , and corresponding pseudo empirical likelihood confidence intervals, which are shown to perform generally better than the customary normal theory intervals.

We refer the readers also to a related article by Wu (2022), which contains a critical review and some extended discussions on theoretical and practical issues with inference from non-probability samples.

Beaumont et al. (2024a) likewise consider integration of  $S_{NP}$  and  $S_{PS}$  samples, again assuming that the probability of inclusion in  $S_{NP}$  only depends on the  $\mathbf{x}$  variables. The authors assume a logistic model  $p_i(\mathbf{\gamma}) = \Pr(\delta_i = 1 | \mathbf{x}_i; \mathbf{\gamma})$  for the inclusion of unit  $i \in U$  in  $S_{NP}$  and estimate  $\mathbf{\gamma}$  by solving the likelihood estimating equations  $\hat{U}(\mathbf{\gamma}) = \sum_{i \in S_{NP}} \mathbf{x}_i - \sum_{i \in S_{PS}} w_i p_i(\mathbf{\gamma}) \mathbf{x}_i = 0$ . The equations  $\hat{U}(\mathbf{\gamma})$  are design unbiased over all possible  $S_{PS}$  selections of the likelihood equations that would be obtained if the  $\mathbf{x}$ -values were known for all  $i \in U$ .

The authors develop a modified AIC criterion for stepwise selection of the  $\mathbf{x}$ -variables in the  $S_{NP}$  sample selection model  $p_i(\gamma)$ . However, a problem with the use of this criterion is that it ignores the relationship between y and the  $\mathbf{x}$ -variables. To deal with this problem, the authors extend their AIC criterion by partitioning the  $S_{NP}$  sample into homogeneous groups  $S_{NP} = S_{NP,1} \cup \ldots \cup S_{NP,G}$  based on the estimated probabilities  $p_i(\hat{\gamma})$  and a ranking method, and then assigning each unit in the  $S_{PS}$  sample to one of the groups. Let  $S_{NP,g}$  and  $S_{PS,g}$  define the  $g^{th}$  sets of units of the non-probability and probability samples, respectively. Assuming that the selection probabilities in each group are the same, the resulting estimated selection probabilities in group g are then  $\hat{p}_g = n_g^{NP} / \hat{N}_g$ , where  $n_g^{NP}$  is the size of  $S_{NP,g}$  and  $\hat{N}_g = \sum_{k \in S_{PS,g}} w_k$ . The estimator of Y is

$$\hat{Y}_{S_{NP}} = \sum_{k \in S_{NP}} \hat{w}_{k}^{NP} y_{k} = \sum_{g=1}^{G} \hat{N}_{g} \overline{y}_{S_{NP,g}}; \hat{w}_{k}^{NP} = \hat{p}_{k}^{-1} = \hat{N}_{g} / n_{g}^{NP}; \ \overline{y}_{S_{NP,g}} = \sum_{i \in S_{NP,g}} \frac{y_{i}}{n_{g}^{NP}}. \quad (2.6)$$

The variance of  $\hat{Y}_{S_{NP}}$  is estimated by an appropriate bootstrap algorithm.

REMARK 6. Rao (2021) reviews several other estimators based on data integration, distinguishing between the case where the target variable y is observed in both samples, and the case where it is only observed in the  $S_{NP}$  sample.

The common feature of all the approaches considered so far is their reliance on the assumption that the selection to the  $S_{NP}$  sample depends on known x-variables, but not on the target y-variable. (See Remark 4 above). In practice, it is likely that the selection to  $S_{NP}$  depends also on y. For example, people participating in a voluntary web survey on political tendency, may choose not to participate in the survey, depending on their

tendency. Administrative data may be missing people who do not participate in government programs, including people who do not have social security numbers, people with housing instability, or people working in the informal economy.

In addition, the  $S_{PS}$  sample used for integration with the  $S_{NP}$  sample may be subject to not missing at random (NMAR) nonresponse, in the sense that that the probability to respond depends also on the target y- variable. For example, the response of people on income may depend on their level of income. Denote by  $R_i$  the response indicator. NMAR nonresponse occurs when,

$$\Pr[R_i = 1 \mid y_i, \mathbf{x}_i, i \in s] \neq \Pr[R_i = 1 \mid \mathbf{x}_i, i \in s].$$
 (2.7)

Pfeffermann et al. (2025) consider data integration when the selection to the  $S_{NP}$  sample and the response probabilities in the  $S_{PS}$  sample depend on both y and  $\mathbf{x}$ , applying the empirical likelihood (EL) approach. It is assumed that  $\mathbf{x}$  is observed in both samples, but y is only observed in the  $S_{NP}$  sample. Let  $I_i^{PS}$  be the sample indicator for  $S_{PS}$ , taking the value 1 if unit i is sampled and 0 otherwise. For  $i \in S_{PS}$ , the sample model of  $\mathbf{x}_i$  is

$$p_{i,PS}^{X} = \Pr(\mathbf{x}_{i} \mid I_{i}^{PS} = 1) = \frac{\Pr(I_{i}^{PS} = 1 \mid \mathbf{x}_{i})}{\Pr(I_{i}^{PS} = 1)} p_{i}^{X},$$
(2.8)

where  $p_i^X = \Pr_U(\mathbf{x} = \mathbf{x}_i)$  is the probability in the population. As can be seen, under informative sampling with respect to  $\mathbf{x}$ , the sample probability  $p_{i,PS}^X$  is different from  $p_i^X$ .

Additionally, it is assumed that  $S_{PS}$  is subject to NMAR nonresponse. Let  $R_i^{PS}$  be the response indicator, taking the value 1 if sample unit  $i \in S_{PS}$  responds and 0 otherwise. Denote by  $R_{PS}$  the set of responding units in  $S_{PS}$ . Then,

$$p_{i,R_{PS}}^{X} = \Pr(\mathbf{x}_{i} \mid I_{i}^{PS} = 1, R_{i}^{PS} = 1) = \frac{\Pr(R_{i}^{PS} = 1 \mid \mathbf{x}_{i}, I_{i}^{PS} = 1)}{\Pr(R_{i}^{PS} = 1 \mid I_{i}^{PS} = 1)} p_{i,PS}^{X}.$$
(2.9)

By (2.8) and (2.9), the respondents model is a function of the true population probability, the conditional expectations of the sampling weights,  $\Pr(I_i^{PS}=1 \mid \mathbf{x}_i) = 1/E_{PS}(w_{i,PS} \mid \mathbf{x}_i)$  (Pfeffermann and Sverchkov 1999);  $w_{i,PS} = 1/\pi_{i,PS}$  are the base sampling weights in  $S_{PS}$ , and the response probabilities  $\Pr(R_i^{PS}=1 \mid \mathbf{x}_i, I_i^{PS}=1)$ . Assuming that the response is

independent of the sample selection,  $E_{PS}(w_{i,PS} \mid \mathbf{x}_i) = E_{R_{PS}}(w_{i,PS} \mid \mathbf{x}_i)$ , in which case the probabilities  $P(I_i^{PS} = 1 \mid \mathbf{x}_i)$  can be estimated by regressing  $w_{i,PS}$  against  $\mathbf{x}_i$ .

The response probabilities  $P(R_i^{PS} = 1 | \mathbf{x}_i, I_i^{PS} = 1)$  in (2.9) are unknown and need to be estimated from the available data by postulating a parametric model,

$$P(R_i^{PS} = 1 | \mathbf{x}_i, I_i^{PS} = 1, \mathbf{\rho}) = g(\mathbf{x}_i; \mathbf{\rho})$$
 (2.10)

for some known function g, (say, a logistic model), with  $\rho$  defining the model parameters. Assuming independence of the sampling and the response, the *empirical respondents'* likelihood based on  $R_{PS}$  is thus,

$$ERL_{PS}(p_i^X) = \prod_{i \in R_{PS}} p_{i,R_{PS}}^X = \prod_{i \in R_{PS}} \frac{\Pr(R_i^{PS} = 1 \mid \mathbf{x}_i, I_i^{PS} = 1)}{\Pr(R_i^{PS} = 1 \mid I_i^{PS} = 1)} p_{i,PS}^X.$$
 (2.11)

Next, consider the  $S_{NP}$  sample. Let  $I_i^{NP}$  be the sample indicator, taking the value 1 if  $i \in S_{NP}$  and 0 otherwise. Denote  $p_i^{XY} = \Pr(\mathbf{x} = \mathbf{x}_i, y = y_i)$ . For  $i \in S_{NP}$ ,

$$p_{i,NP}^{XY} = \Pr(\mathbf{x}_i, y_i \mid I_i^{NP} = 1) = \frac{\Pr(I_i^{NP} = 1 \mid \mathbf{x}_i, y_i)}{\Pr(I_i^{NP} = 1)} p_i^{XY},$$
(2.12)

where  $P(I_i^{NP}=1) = \sum_{i \in NP} P(I_i^{NP}=1 \mid \mathbf{x}_i, y_i) p_i^{XY}$ . Because no sampling weights for  $S_{NP}$  are

available, the probabilities  $P(I_i^{NP} = 1 | \mathbf{x}_i, y_i)$  need to be modelled parametrically,

$$\Pr(I_i^{NP} = 1 \mid \mathbf{x}_i, y_i; \gamma) = h(y_i, \mathbf{x}_i; \gamma)$$
(2.13)

for some known function h, with  $\gamma$  defining the model parameters. Assuming independence of the  $S_{\mathit{NP}}$  data, the *empirical likelihood* based on  $S_{\mathit{NP}}$  is

$$ESL_{NP}(p_i^{XY}) = \prod_{i \in S_{NP}} p_{i,NP}^{XY}.$$
 (2.14)

Assuming no overlap between the two samples, the *empirical likelihood* based on the data in  $S_{\mathit{NP}}$  and  $S_{\mathit{PS}}$  is,

$$EL_{R_{PS} \cup NP} = ERL_{PS}(p_i^X)ESL_{NP}(p_i^{XY}) = \prod_{i \in R_{PS}} p_{i,R_{PS}}^X \prod_{i \in NP} p_{i,NP}^{XY}.$$
 (2.15)

The unknown parameters in (2.15) are the population probabilities  $p_i^X$ ,  $p_i^{XY}$ , the sampling parameters  $\gamma$  and the response parameters  $\rho$ . The likelihood is maximized subject to constraints on the unknown probabilities and calibration constraints.

REMARK 7. The unknown probabilities  $\{p_i^X\}$  can also be estimated from the  $S_{NP}$  sample;  $\hat{p}_{i,NP}^X = \sum_{\{i: x=x_i\}} \hat{p}_{i,NP}^{XY}$ . This implies two sets of estimates of the probabilities  $\{p_i^X\}$ , which need

to be harmonized. See Marella and Pfeffermann (2023) for possible harmonization procedures. The final, integrated estimate of  $p_i^{XY}$  is  $\hat{p}_i^{XY} = \hat{p}_i^X (\hat{p}_{i,NP}^{XY} / \hat{p}_{i,NP}^X)$ , where  $\hat{p}_i^X$  is the harmonized estimator.

The population total Y can be estimated in one of the following two ways:

$$\hat{Y}_{NP}(1) = N \sum_{i \in S_{NP}} y_i \hat{p}_i^Y \; ; \; \hat{Y}_{NP}(2) = N \frac{\sum_{i \in NP} \hat{P}r^{-1} (I_i^{NP} = 1 \mid \mathbf{x}_i, y_i) y_i}{\sum_{i \in NP} \hat{P}r^{-1} (I_i^{NP} = 1 \mid \mathbf{x}_i, y_i)} \,, \tag{2.17}$$

where  $\hat{p}_i^y = \sum_{i;y=y_i} \hat{p}_i^{XY}$ . See Pfeffermann et al. (2025) for an empirical comparison of the performance of the two estimators.

REMARK 8. One of the reviewers of this article raised a concern about the model used for the selection model to the  $S_{NP}$  sample, noting that it seems difficult to obtain robustness to deviations from the model. As discussed in Section 4.3 and illustrated in Section 5, the  $S_{NP}$  model can be tested.

#### 3. INFERENCE FROM A NONPROBABILITY SAMPLE WITHOUT INTEGRATION

In Section 2, we considered methods of inference from a nonprobability sample, based on integration of the  $S_{NP}$  sample with an appropriate probability sample  $S_{PS}$ . In this section, we consider methods for adjusting the selection bias of the  $S_{NP}$  sample, without integration with a  $S_{PS}$  sample (see Remark 1).

We start with an approach based on calibration. The basic idea underlying this approach is to change some base weights,  $d_{j,NP}$  (say, based on propensity scores) to new weights  $d_{j,NP}^{cal}$ , so that when applied to a set of variables Z observed in  $S_{NP}$  and for which the true population totals are known, the  $S_{NP}$  survey estimates will equal the corresponding totals;

 $\sum_{j \in S_{NP}} d_{j,NP}^{cal} \mathbf{z}_j = \mathbf{T}_z$ , where  $\mathbf{T}_{\mathbf{z}}$  are the known population totals. (In practice, the true totals can be replaced by reliable estimates from a probability same.) See AAPOR (2010) and Baker et al. (2013) for review of methods that follow this approach, and Kott and Ridenhour (2024) reviewed in Section 2.

The success of this approach depends on the availability of calibration variables, which are highly correlated with the target y-variable (good prediction power). Lee and Valliant (2009) illustrate that combining propensity scores and calibration adjustments is more effective in reducing the bias of  $S_{NP}$  estimates than using just one of the approaches. See also Elliott and Valliant (2017).

Kim and Wang (2019) propose the use of inverse sampling to obtain a representative sample from the finite population, and hence to correct for the selection bias of the  $S_{\mathit{NP}}$  sample. The proposed inverse sampling can be viewed as a special case of two-phase sampling, where the first phase is the  $S_{\mathit{NP}}$  sample and the second phase is a subsample from the first-phase sample to correct for the selection bias.

Denote, as before, by  $\delta_i$  the indicator of whether unit  $i \in U$  is included in the  $S_{NP}$  sample. It is assumed that  $\Pr(\delta_i = 1 | y_i, \mathbf{x}_i) = \Pr(\delta_i = 1 | \mathbf{x}_i) > 0$  for all  $i \in U$ . The  $S_{NP}$  sample contains the values  $(y_i, \mathbf{x}_i)$ ,  $i \in S_{NP}$ . Denote by  $f(\mathbf{x})$  the population distribution of the  $\mathbf{x}$ -variables. If  $f(\mathbf{x})$  is known, an asymptotic unbiased estimator of  $\theta = E(Y)$  is,

$$\hat{\theta}_{S_{NP1}} = \sum_{i \in S_{NP}} \frac{f(\mathbf{x}_i)}{f(\mathbf{x}_i | \delta_i = 1)} y_i / \sum_{i \in S_{NP}} \frac{f(\mathbf{x}_i)}{f(\mathbf{x}_i | \delta_i = 1)} = \sum_{i \in S_{NP}} w_{i1} y_i.$$
(3.1)

For the more practical case where only the mean  $\overline{\mathbf{X}}_U = \sum_{i \in U} \mathbf{x}_i / N$  is known, the authors approximate  $f(\mathbf{x})$  by the function  $f_0(\mathbf{x})$ , which minimizes the Kullback–Leibler distance. The solution to the minimization distance is,

$$f_0(\mathbf{x}) = f(\mathbf{x} \mid \delta = 1) \frac{\exp(\mathbf{x}' \lambda)}{E[\exp(\mathbf{x}' \lambda \mid \delta = 1)]}, \text{ with } \lambda \text{ satisfying } \int \mathbf{x} f_0(\mathbf{x}) d\mathbf{x} = \overline{\mathbf{X}}_U.$$
 (3.2)

With this approximation, the estimator  $\hat{\theta}_{S_{NP1}}$  in (3.1) is replaced by,

$$\hat{\theta}_{S_{NP2}} = \sum_{i \in S_{NP}} w_i^* y_i ; w_i^* = \frac{\exp(\mathbf{x}_i' \hat{\lambda})}{\sum_{i \in S_{NP}} \exp(\mathbf{x}_i' \hat{\lambda})}, \text{ with } \hat{\lambda} \text{ satisfying } \sum_{i \in S_{NP}} w_i^* \mathbf{x}_i = \overline{\mathbf{X}}_U.$$
 (3.3)

Finally, the authors propose to select the second-phase sample from  $S_{NP}$  with probabilities  $\pi_{i2|i} = nw_i^*, i \in S_{NP}$  with the weights  $\{w_i^*\}$  defined by (3.3) and  $n \leq [\max_{i \in S_{NP}} \{w_i^*\}]^{-1}$ , yielding the approximately design-unbiased estimator of the  $\hat{\theta}_{S_{NP1}}$  estimator defined in (3.1),

$$\hat{\theta}_{S_{NP3}} = \sum_{i \in S_{NP}} \frac{1}{\pi_{i21}} w_i^* y_i = \frac{1}{n} \sum_{i=1}^n y_i.$$
 (3.4)

A simple estimator of the design variance of  $\,\hat{\theta}_{\scriptscriptstyle S_{NP3}}\,$  is proposed.

The two approaches considered so far assume that the selection to the  $S_{NP}$  sample is noninformative in the sense that  $\Pr(\delta_i = 1 \mid y_i, \mathbf{x}_i) = \Pr(\delta_i = 1 \mid \mathbf{x}_i) > 0$  for all  $i \in U$ . However, as discussed before, this assumption may not hold in practice, and in what follows, we consider alternative approaches aimed to deal with the case of informative sample selection.

Sayag et al. (2022) consider the following problem, underlying the computation of monthly house price indices (HPI) in many countries. A large amount of the house sales are reported several months after they occur, implying that if not accounted for, the provisional HPIs based on the on-time reported transactions are subject to large revisions, as further transactions are reported. This happens because the late-reported transactions behave differently from the transactions reported on time. This is a nice example of a nonprobability sample (the on-time reported sales), which is subject to selection bias due to late data availability of some of the sales (~40% in Israel).

To deal with this problem, the authors propose nowcasting three types of variables and adding them as input data to a hedonic regression model used for the computation of the HPI: (1)- the average characteristics of the upcoming late-reported transactions, such as the average number of rooms, the average net area size, the average age of the sold houses, etc. (2)- the average price of the late-reported transactions and (3)- the number of late-reported transactions. The three types of variables are nowcasted based on simple models fitted to data from previous months. Application of the proposed methodology shows more than 50% reduction in the magnitude of the revisions. This is a unique example of a time series of non-representative nonprobability samples for which the true population data (all the sales corresponding to a given month) become known only several months later.

Kim and Morikawa (2023) assume a non-ignorable (informative) sample selection model  $\pi_i(y_i,\mathbf{x}_i;\pmb{\phi})=\Pr(\delta_i=1\mid y_i,\mathbf{x}_i;\pmb{\phi}),$  where  $\delta_i=(1,0)$  is the  $S_{\mathit{NP}}$  sample indicator, assuming that the variables  $\mathbf{x}_i$  are known for all  $i\in U$  and  $\pi_i(y_i,\mathbf{x}_i)>0$  for all  $i\in U$ . For the case where the population model  $f(y_i\mid\mathbf{x}_i)$  is known, the authors propose estimating  $\pmb{\phi}$  by maximizing the likelihood,

$$L_{obs}(\boldsymbol{\phi}) = \prod_{i \in U} [f(y_i \mid \mathbf{x}_i) \pi(y_i, \mathbf{x}_i; \boldsymbol{\phi})]^{\delta_i} [1 - \tilde{\pi}(\mathbf{x}_i; \tilde{\boldsymbol{\phi}})]^{(1 - \delta_i)}; \tilde{\pi}(\mathbf{x}_i; \tilde{\boldsymbol{\phi}}) = \mathbb{E}[\pi(y_i, \mathbf{x}_i; \boldsymbol{\phi}) \mid \mathbf{x}_i].$$
(3.5)

However, this likelihood requires modelling the population model and the authors note that the MLE estimator obtained from (3.5) is not robust to misspecification of the model. Consequently, they develop a likelihood based on the model  $f_{SNP}(y_i | \mathbf{x}_i) = f(y_i | \mathbf{x}_i, \delta_i = 1)$ , which can be identified and estimated consistently.

Alternatively, the authors develop a methodology for estimating  $\phi$  and the population mean of the y-values by applying the Empirical Likelihood (EL) approach. For the case where the selection probabilities  $\pi_i(y_i, \mathbf{x}_i)$  are known, the authors propose estimating the  $p_i$ 's underlying the EL by maximizing the likelihood,  $l(p) = \sum_{i \in S_{NP}} \log(p_i)$ , subject to the constraints (1)-  $\sum_{i \in S_{NP}} p_i = 1$ , (2)-  $\sum_{i \in S_{NP}} p_i \pi_i(y_i, \mathbf{x}_i) = n/N$ , (3)  $\sum_{i \in S_{NP}} p_i \mathbf{x}_i = \overline{\mathbf{X}}_U$ , where n is the size of the  $S_{NP}$  sample, N is the population size and  $\overline{\mathbf{X}}_U = \sum_{i \in U} \mathbf{x}_i/N$ . The constraint (2) is referred to as a bias calibration constraint, whereas the constraint (3) is added to improve the efficiency of EL estimator.

In practice, the sample selection probabilities are unknown. The authors assume a parametric model;  $\pi_i(y_i,\mathbf{x}_i)=g(y_i,\mathbf{x}_i;\pmb{\phi})$  (say, logistic), and estimate  $\hat{\pi}_i=g(y_i,\mathbf{x}_i;\hat{\pmb{\phi}})$  by solving the estimating equations  $\sum_{i=1}^N [\frac{\delta_i}{g(y_i,\mathbf{x}_i;\pmb{\phi})}-1]\mathbf{x}_i=\mathbf{0}$ . These equations do not require knowledge of the  $\mathbf{x}$ -values for every unit in the population. By considering the estimated probabilities  $\hat{\pi}_i(y_i,\mathbf{x}_i)=g(y_i,\mathbf{x}_i;\hat{\pmb{\phi}})$  as the true inclusion probabilities, the authors maximize the constrained EL likelihood defined above with the bias calibration constraint (2) replaced by  $\sum_{i\in\mathcal{S}_{NP}}p_i\pi_i(y_i,\mathbf{x}_i;\hat{\pmb{\phi}})=N^{-1}\sum_{i=1}^N\pi_i(\mathbf{x}_i;\hat{\pmb{\phi}})$ , which does require knowledge of the population  $\mathbf{x}$ -values, yielding the estimates  $\{\hat{p}_i\}$ . The population mean of the y-values are estimated as,

$$\hat{\overline{Y}}_{EL,IPW} = \frac{1}{N} \sum_{i \in S_{NP}} \frac{y_i}{\hat{\pi}_i(y_i, \mathbf{X}_i)} \quad \hat{\overline{Y}}_{EL} = \sum_{i \in S_{NP}} \hat{p}_i y_i . \tag{3.6}$$

The authors derive asymptotic properties of their estimators and variance estimators.

This article proposes a novel approach for estimating finite population means from  $S_{NP}$  samples subject to nonignorable selection probabilities, but the assumption that the  $\mathbf{x}$ -variables are known for every unit in the population is restrictive.

REMARK 9. In Section 2, we proposed a method of inference from a  $S_{NP}$  sample alone, which likewise combines a non-ignorable sample selection model with the empirical likelihood. See Equations (2.12)-(2.14). This method does not require knowledge of the  $\bf x$  -variables for every unit in the population. See also Section 4 below.

#### 4. A NEW (OLD) APPROACH FOR INFERENCE FROM A NONPROBABILITY SAMPLE

## 4.1 Relationship between the population distribution and the $S_{\mathit{NP}}$ distribution

Following, we propose an alternative approach for inference from a nonprobability sample alone. It relies in large on Pfeffermann and Sverchkov (1999).

Denote the model holding for the target variable y in U by  $f_U(y_i \mid \mathbf{x}_i)$ . Denote the model holding for y in the  $S_{NP}$  sample by  $f_{S_{NP}}(y_i \mid \mathbf{x}_i)$ , and let  $\delta_i = 1(0)$  if  $i \in S_{NP}(i \notin S_{NP})$ . The target model is  $f_U(y_i \mid \mathbf{x}_i)$ , but observations  $\{y_i, \mathbf{x}_i\}$  are only available for  $f_{S_{NP}}(y_i \mid x_i)$ . We assume,  $\Pr(i \in S_{NP}) > 0$  for all  $i \in U$  (also assumed in the other approaches considered before). The two distributions are connected via the link function  $\Pr(\delta = 1 \mid y, \mathbf{x})$ .

$$f_{S_{NP}}(y_i \mid \mathbf{x}_i) = f(y_i \mid \mathbf{x}_i, \delta_i = 1) \stackrel{Bayes}{=} \frac{\Pr(\delta_i = 1 \mid \mathbf{x}_i, y_i) f_U(y_i \mid \mathbf{x}_i)}{\Pr(\delta_i = 1 \mid \mathbf{x}_i)}. \tag{4.1}$$

As discussed below, the relationship (4.1) enables estimating the target population distribution from the observations in  $S_{NP}$  alone. Notice that  $f_{S_{NP}}(y_i | \mathbf{x}_i) = f_U(y_i / \mathbf{x}_i)$  iff

 $\Pr(\delta_i = 1 \mid \mathbf{x}_i, y_i) = \Pr(\delta_i = 1 \mid \mathbf{x}_i) \ \forall y_i$ , in which case the model fitted based on the  $S_{NP}$  sample holds for the population data and if the  $\mathbf{x}$ -values are known for all  $i \in U$ , (or in the case of a linear population model  $\overline{\mathbf{X}}_U$  is known), inference based on the  $S_{NP}$  sample is valid. See Rao (2021) for discussion of this method under these conditions.

REMARK 10. In the first part of their article, Kim and Morikawa (2023) also assume parametric models for the population model and the sample selection probabilities, (see above), but we do not assume knowledge of the population  $\mathbf{x}$ -values. Additionally, the authors estimate the parameters underlying the sample selection model outside the likelihood, whereas we estimate them jointly with the population model parameters (see below). We utilize similar calibration constraints to the ones used by Kim and Morikawa (2023), see Equation (4.3) below. We also test the goodness of fit of the resulting model  $f_{S_{NP}}(y_i | \mathbf{x}_i)$ , see section 4.3.

The probabilities  $\Pr(\delta_i = 1 \mid \mathbf{x}_i, y_i)$  need to be modelled. They are allowed to depend on the target y variable, thus accounting for informative sample selection. They may depend also on other variables  $\mathbf{Z}$ , but we only need to model  $\Pr(\delta_i = 1 \mid \mathbf{x}_i, y_i)$ . The use of a Logistic model for  $\delta_i$  has some theoretical justification. See Lemma 1 in Pfeffermann et al. (2025) for details. When  $\mathbf{Z}$  is observed in the  $S_{NP}$  sample, we may include it among the  $\mathbf{x}$ -variables.

### 4.2 Estimation of model parameters

Unlike the use of the empirical likelihood approach, application of this approach requires specifying the population model and the model for the sample selection probabilities, which depend on unknown parameters that need to be estimated from the observations in  $S_{NP}$ . Adding parameters to (4.1), we have

$$f_{S_{NP}}(y_i \mid \mathbf{x}_i; \boldsymbol{\beta}, \boldsymbol{\phi}) = \frac{\Pr(\delta_i = 1 \mid y_i, \mathbf{x}_i; \boldsymbol{\phi}) f_U(y_i \mid \mathbf{x}_i; \boldsymbol{\beta})}{\Pr(\delta_i = 1 \mid \mathbf{x}_i; \boldsymbol{\phi}, \boldsymbol{\beta})}.$$
(4.2)

Assuming independence of the observations in  $S_{\mathit{NP}}$ , the corresponding log likelihood is

$$l_{S_{NP}}(\pmb{\phi},\pmb{\beta};y) = \sum_{i \in S_{NP}} \log f_{S_{NP}}(y_i \mid \pmb{\mathrm{x}}_i; \pmb{\beta}, \pmb{\phi}) \text{ , which we maximize subject to the constraints,}$$

$$\frac{1}{N} \sum_{i \in \mathcal{S}_{NP}} \frac{1}{\Pr(\delta_i = 1 \mid y_i, \mathbf{x}_i; \boldsymbol{\phi})} \mathbf{x}_i = \frac{1}{N} \sum_{j \in U} \mathbf{x}_j = \overline{\mathbf{X}}_U.$$
(4.3)

The constraints (4.3) are used for enhancing the estimation of the parameters ( $\beta$ , $\phi$ ). We assume throughout that  $\mathbf{x}$  contains a "1" in the first position.

REMARK 11. In the empirical study in Section 5 with continuous y, we approximated the probabilities  $\Pr(\delta_i = 1 \mid \mathbf{x}_i; \boldsymbol{\phi}, \boldsymbol{\beta})$  by Riemann's sums over 350 sub-groups of the y-values. When y is binary,

$$Pr(\delta_i = 1 \mid \mathbf{x}_i; \boldsymbol{\phi}, \boldsymbol{\beta}) = Pr(\delta_i = 1 \mid y_i = 1, \mathbf{x}_i; \boldsymbol{\phi}) Pr(y_i = 1 \mid \mathbf{x}_i; \boldsymbol{\beta})$$
$$+ Pr(\delta_i = 1 \mid y_i = 0, \mathbf{x}_i; \boldsymbol{\phi}, \boldsymbol{\beta}) Pr(y_i = 0 \mid \mathbf{x}_i; \boldsymbol{\beta}).$$

We maximized the likelihood with the constraints by use of the SAS procedure NLIN, iterating between the maximization with respect to  $\phi$  for given  $\beta$ , and the maximization of  $\beta$  for given  $\phi$ , with the "given" values defined by the estimates in the previous iteration. See Section 5 for how we estimated the population mean of the y-values in our simulations.

### 4.3 Model testing and Identifiability concerns

The application of the proposed approach assumes a model  $f_U(y_i | \mathbf{x}_i; \boldsymbol{\beta})$  for the population values and a model  $\Pr(\delta_i = 1 | y_i, \mathbf{x}_i; \boldsymbol{\phi})$  for the selection probabilities, which permits estimating the parameters  $(\boldsymbol{\phi}, \boldsymbol{\beta})$  by means of (4.2) and (4.3), using the data in  $S_{NP}$ . No direct testing of the population model or the model for the selection probabilities is possible, since no data are available from the population distribution and the y-values are unknown for units  $j \notin S_{NP}$ . However, contrary to a common perception that it is impossible to test a model fitted to the  $S_{NP}$  data, we contend this is not true. We have observations from the fitted model, so we are faced with the classical problem of testing the goodness of fit of a hypothesized model to the observed data. See Krieger and Pfeffermann (1997) and Pfeffermann and Sikov (2011) for plausible tests.

REMARK 12. Rejection of the null hypothesis implies that at least one of the two models is misspecified. See Section 5 for examples and the concluding remarks in Section 6.

A common argument in favor of the claim that the sample model cannot be tested is that it may be the case that there is more than one combination of a population model and a selection model, yielding the same model for the observed data, such that the model fitted to the data in  $S_{NP}$  is not identifiable or weekly identifiable. Pfeffermann and Landsman (2011) and Wang et al. (2014) establish conditions under which the model  $f_{S_{NP}}(y_i | \mathbf{x}_i)$  is

identifiable, with references to other related studies. See Section 5 for the identifiability conditions for the models considered.

REMARK 13. In a highly cited article, Molenberghs et al. (2008) prove and illustrate that for every NMAR model fitted to a set of data, there is a MAR counterpart providing exactly the same fit to the data. The authors note that "such a construction does not lead to a member of a conventional parametric family". A simple example for this argument is where the population model  $f_U(y \mid \mathbf{x})$  is assumed to be defined by the sample model  $f_{S_{NP}}(y_i \mid \mathbf{x}_i)$  (Eq. 4.2), and the sample inclusion probability satisfies  $\Pr(\delta_i = 1 \mid y_i, \mathbf{x}_i; \phi) = \Pr(\delta_i = 1 \mid \mathbf{x}_i; \phi)$ . Clearly,  $f_U(y_i \mid \mathbf{x}_i) = f_{S_{NP}}(y_i \mid \mathbf{x}_i)$  is a very odd population distribution. Molenberghs et al. (2008) also note that "we can make progress if attention is confined to a given parametric family, in which we put sufficiently strong prior belief". This is what we do under our proposed approach. Notice that the selection model is used to obtain valid estimates of the population model, and as shown below and illustrated in Section 5, it can be tested.

Consider first the case where y is a continuous variable. In our empirical applications, we applied the following UNIF test statistic (Krieger and Pfeffermann, 1997).

#### **Preliminaries:**

- 1- For a continuous variable Z with cumulative distribution F,  $F(z) \sim U(0,1)$ .
- 2- Under general conditions, the set of all the moments of F(z) determine the distribution.

#### **Proposed test:**

- (*i*) Compute  $T_i = F_{S_{NP}}(y_i \mid \mathbf{x}_i), i = 1,...,n$  based on the estimated coefficients  $(\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\phi}})$ .
- (ii) Compute the sample moments  $u_m = \sum_{i=1}^n T_i^m / n, m = 1,...,M$ .
- (iii) Compute the Wald test statistic based on the estimated sample moments.

For the moments of the U(0,1) distribution,  $\mu_m = E(u_m) = 1/(m+1)$ ;

 $Cov(u_m, u_l) = ml / [(m+1)(l+1)(m+l+1)n]$ . Assuming  $u' = (u_1, ..., u_m)$  is normal,

UNIF = 
$$(u - \mu)' \Sigma^{-1} (u - \mu) \sim \chi_M^2$$
, (4.4)

where  $\Sigma$  is the Variance-Covariance matrix defined by the covariances above.

REMARK 14. In the proposed test, we replace the true moments by the estimated moments. The estimators  $(\hat{\beta}, \hat{\phi})$  are obtained by MLE and under some regularity conditions, they converge almost surely (a.s.) to the true parameters  $(\beta, \phi)$  (Zacks, 1971). Then, if the true distributional function F is smooth, e.g. twice differentiable with respect to  $\beta$  and  $\phi$ ,  $F(y_i | \mathbf{x}_i, \delta_i = 1; \hat{\phi}, \hat{\beta}) \xrightarrow{a.s.} F(y_i | \mathbf{x}_i, \delta_i = 1; \phi, \beta)$ , justifying the use of the UNIF test defined by (4.4). See Figure 1 in Section 5 for a simulation illustration.

REMARK 15. In our simulation study we used M=5 moments, which was found to perform well in Krieger and Pfeffermann (1997). Notice that  $Corr^2(u_m,u_{m-1}) = (1-\frac{1}{4}m^2)$ , so higher order moments add only marginally to the power of the test.

For the case where y is binary, we apply in Section 5 the Hosmer and Lemeshow (1980, hereafter H-L) test defined as follows:

- (*i*) Sort the observed data in  $S_{NP}$  based on the estimated probabilities  $\hat{\eta}_i = \hat{\Pr}(y_i = 1 \mid \mathbf{x}_i, \delta_i = 1), \, i \in S_{NP} \,.$
- $(\pmb{ii} \ ) \ \text{Divide the sorted data into G groups of approximately equal size } \ n_{_g} \cong (n/G) \ \text{ and compute for each group } \ g \ ; \ o_{_g} \text{- the number of values } \ y=1 \ \text{and } \ \overline{\eta}_{_g} = \frac{1}{n_{_g}} \sum_{_{i \in g}} \hat{\eta}_{_i} \ . \ \text{The test statistic is,}$

$$H - L = \sum_{g=1}^{G} \frac{(o_g - n_g \overline{\eta}_g)^2}{n_o \overline{\eta}_o (1 - \overline{\eta}_o)} \sim \chi_{(G-2)}^2.$$
 (4.5)

#### 5. SIMULATION RESULTS

In this section, we present simulation results to illustrate the performance of our proposed approach, separately for the case where the target variable y is continuous, and for the case where it is binary.

## 5.1 Simulation setup with a continuous target variable- correct model

We start by repeating the same simulation study as performed by Kim and Morikawa (2023), which consists of the following steps:

**1.** Generate 5,000 population values as  $y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \varepsilon_i$ , where  $x_{1i}, x_{2i} \sim N(2,1); \ \varepsilon_i \sim N(0,1). \ \text{(The values of the } \boldsymbol{\beta} \ \text{ coefficients are in Table 1 below.)}$ 

**2.** Generate selection probabilities to the  $S_{\mathit{NP}}$  sample as,

$$\pi_{i,S_{NP}} = \Pr(\delta_i = 1 \mid y_i, \mathbf{x}_i) = \frac{\exp(\phi_0 + \phi_1 x_{1i} + \phi_2 y_i)}{1 + \exp(\phi_0 + \phi_1 x_{1i} + \phi_2 y_i)}. \text{ (The } \phi \text{ coefficients are in Table 1.)}$$

- 3. Repeat Steps 1 and 2 1,000 times, yielding an average selection rate of 50%.
- **4.** For each simulation, estimate the model parameters and the population mean  $\overline{Y}_U = \sum_{i=1}^{5,000} y_i \,/\, N$  .

#### **Estimators considered:**

 $\hat{\bar{Y}}_{U,X\,known} = \frac{1}{N} \sum_{i \in U} (\hat{\beta}_0 + \hat{\beta}_1 x_{1i} + \hat{\beta}_2 x_{2i}). \text{ The x-variables are known for every unit } i \in U \text{ , } \boldsymbol{\beta}$  is estimated by maximization of the likelihood  $l_{S_{NP}}(\boldsymbol{\phi},\boldsymbol{\beta};y) = \sum_{i \in S_{NP}} \log f_{S_{NP}}(y_i \mid \mathbf{x}_i;\boldsymbol{\beta},\boldsymbol{\phi}) \text{ , under the constraints in (4.3). Note: since the population model is linear, it suffices to know the population means of the <math>\mathbf{x}$ -variables.

$$\hat{\bar{Y}}_{U,GREG} = \frac{\sum\limits_{i \in S_{NP}} k_i y_i}{\sum\limits_{i \in S_{NP}} k_i} + \hat{\mathbf{B}}'_{pk} [\mathbf{\bar{X}}_U - \frac{\sum\limits_{i \in S_{NP}} k_i \mathbf{X}_i}{\sum\limits_{i \in S_{NP}} k_i}], \quad k_i = (1/\hat{\pi}_{i,S_{NP}}), \text{ the GREG estimator with the }$$

standard base sampling weights  $w_i = (1/\pi_i)$  replaced by  $k_i = (1/\hat{\pi}_{i,S_{NP}})$ .  $\hat{\mathbf{B}}_{pk}$  is the probability weighted estimator of  $\boldsymbol{\beta}$ , with weights  $k_i$ .

 $\hat{\overline{Y}}_{U,KM} = \sum_{i \in S_{NP}} \hat{p}_i y_i$ , the estimator of Kim and Morikawa (2023). ( $\hat{\overline{Y}}_{EL}$  in Equation 3.6).

 $\hat{Y}_{U,MAR}$ - the estimator obtained by assuming that the selection probabilities only depend on the  $\mathbf{x}$ -variables;  $\Pr(\delta_i = 1 \mid y_i, \mathbf{x}_i) = \Pr(\delta_i = 1 \mid \mathbf{x}_i)$ , where  $\mathbf{x}_i = (x_{1i}, x_{2i})'$ . We assume a logistic model, using all the population  $\mathbf{x}$ -values.

The first 2 estimators are obtained by application of our approach. The estimation of the  $\beta$ -coefficients in the first estimator is only based on the data in  $S_{NP}$ .

REMARK 16. An important question regarding the models used in this simulation study is whether the resulting sample model  $f_{S_{NP}}(y_i \mid \mathbf{x}_i; \boldsymbol{\beta}, \boldsymbol{\phi}) = \frac{\Pr(\delta_i = 1 \mid y_i, \mathbf{x}_i; \boldsymbol{\phi}) f_U(y_i \mid \mathbf{x}_i; \boldsymbol{\beta})}{\Pr(\delta_i = 1 \mid \mathbf{x}_i; \boldsymbol{\phi}, \boldsymbol{\beta})}$  is

identifiable. By identifiability we mean that there are no different pairs [ $\Pr_j(\delta_i=1|\ y_i,\mathbf{x}_i;\phi_j)$ ],  $f_{Uj}(y_i|\mathbf{x}_i;\beta_j)$ ], j=1,2 inducing the same sample model for every y and  $\mathbf{x}$ . Pfeffermann and Landsman (2011) consider sets of conditions guaranteeing the identifiability of the sample model. In particular, for the case of a normal population model and a logistic model for the sample selection probabilities, the sample model is identifiable if the  $\mathbf{x}$ -variables in the two models differ by at least one variable. Notice that in the models underlying the present simulation, the population model is a function of  $(x_{1i}, x_{2i})$ , but the selection logistic model is only a function of  $x_{1i}$ , so that the identifiability condition defined above is satisfied.

## 5.2 Results for continuous case when fitting the correct model, 1,000 simulations.

Table 1. Estimation of model coefficients under the proposed method.

	Population model coefficients			Selection model coefficients		
	$oldsymbol{eta}_0$ $oldsymbol{eta}_1$ $oldsymbol{eta}_2$			$\phi_0$	$\phi_{_{\! 1}}$	$\phi_2$
True coefficients	-4	1	1	-2	1	0.5
Mean estimators*	-3.92	0.98	0.99	-2.15	0.80	0.43
Empirical S.E.	.004	.001	.001	0.023	0.008	0.002
Mean PWR estimator	-3.88	0.96	0.99	NA	NA	NA
Empirical S.	0.006	.002	.001	NA	NA	NA

<sup>\*</sup> The mean estimators are the MLE estimators. The PvvR estimator is the probability weighted estimator with weights  $k_i = (1/\hat{\pi}_{i,S_{NR}})$ .

As can be seen, the  $\beta$  coefficients are estimated quite accurately on average. The estimators of the  $\phi$  coefficients are less accurate, but the estimators of the population mean in Table 2 still have a negligible bias with these estimators.

Table 2. Estimation of population mean. (Mean true value =-0.00), 1,000 simulations.

Method	Bias	Emp. Var X 1000	MSE X 1000 (Bootstrap estimates)*
$\hat{m{Y}}_{U,Xknown}$	-0.01	2.263	2.363 (3.36)
$\hat{ar{Y}}_{U,GREG}$	-0.02	2.423	2.823 (3.89)
$\hat{Y}_{U,KM}$	0.01	2.030	2.080 ()
$\hat{oldsymbol{Y}}_{U,MAR}$	0.25	2.106	64.606 (65.11)

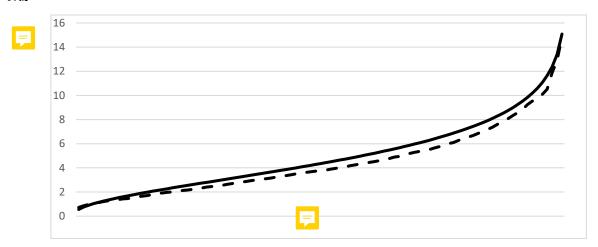
<sup>\*</sup> The bootstrap MSE estimates are based on 100 simulations with 100 bootstrap samples for each simulation.

Estimation of the population mean of the y-values is the primary target of inference in the simulation study and the first three estimators are seen to be literally unbiased. The estimator  $\hat{ar{Y}}_{\!\!U,\mathit{KM}}$  uses all the population **x**-values and performs best. The estimator  $\hat{ar{Y}}_{U,X\,known}$  likewise uses all the population  ${f x}$ -values (or  ${f x}_i$  ,  $i\in S_{NP}$  and  $\overline{m{X}}_U$  ), but the estimation of the model coefficients is only based on the  $S_{\mathit{NP}}$  sample. The estimator  $\hat{\overline{Y}}_{\!\!U\!,\!GREG}$  uses the  $S_{\scriptscriptstyle NP}$  model for estimating the  $\phi$ -coefficients and likewise performs well on average although with somewhat larger variance and MSE. The bootstrap MSE estimators are conservative with large upward bias. We selected the bootstrap samples by following the procedure proposed in Sverchkov and Pfeffermann (2004), which consists of selecting with replacement a pseudo-population from the sample with probabilities proportional to  $k_i = (1/\hat{\pi}_{i,S_{NP}})$  , and then selecting the bootstrap samples  $S_{NP}^b$  with the estimated probabilities  $\hat{\pi}_{_{i,S_{NP}}}$  obtained from the original sample. We only considered 100 simulations and 100 bootstrap samples for each simulation, which may explain the upward biases. As expected, the estimator  $\hat{\overline{Y}}_{U,MAR}$ , which assumes that the selection probabilities only depend on the x-variables has a large positive bias and extremely large MSE. Kim and Morikawa (2023) obtained similar bias and MSE figures in this case.

Overall, the use of our proposed approach seems to perform well in this part of the simulation study.

**Model testing:** As discussed In Section 4.3, our proposed approach enables testing the models assumed for the population and the sample selection probabilities. Figure 1, compares the empirical cumulative distribution of the UNIF statistic (Equation 4.4) with the corresponding  $\chi_M^2$  distribution for the case M=5.

Figure 1. Empirical cumulative distribution of UNIF (dashed curve) and under the  $\chi_M^2$  distribution (solid curve). M=5, 1,000 simulations.



We applied the UNIF test for this part of the simulation study and obtained the following results for the case of M=5 and  $\alpha=0.05$  significance level.

	Mean	Standard Deviation	Minimum	Maximum
UNIF statistic	4.64	2.93	0.45	22.80
P-value	0.53	0.28	~0	0.99
H0 not rejected	0.97	0.18	0	1

We conclude that the UNIF test performs well when testing the correct model, with non-rejection rate of 97%.

## 5.3 Application of the proposed procedure when the models are misspecified

In Section 5.2 we assume that the population model and the model for the selection probabilities are specified correctly. In this section, we consider the case where they are misspecified, using the same simulation setup as in Section 5.1.

Case 1. The population model is specified correctly, the sample selection model is misspecified.

In this case, we selected the  $S_{NP}$  sample with probabilities,

 $\pi_{i,S_{NP}} = \frac{\exp(\phi_0 + \phi_1 x_{1i} + \phi_2 y_i^2)}{1 + \exp(\phi_0 + \phi_1 x_{1i} + \phi_2 y_i^2)}, \text{ but assumed as our working model that the selection}$ 

probabilities are as in Section 5.1 (with  $y_i$  in the exponent rather than  $y_i^2$ ). The population model of y is specified correctly. The average selection rate over the 1,000 simulations is in this case 0.53, similar to what we had before.  $\phi$ 

Table 3. Estimation of model coefficients with misspecified selection probabilities, 1,000 simulations.

	Population model coefficients			Selection model coefficients		
	$eta_0$ $eta_1$ $eta_2$			$\phi_0$	$\phi_{_{\mathrm{l}}}$	$\phi_2$
True coefficients	-4	1	1	-2	1	0.5
Mean estimators	-4.66	1.14	1.14	-0.52	0.33	0.02
Empirical S.E.	.002	.001	.001	0.007	0.003	0.002
Mean PWR estimator	-4.69	1.16	1.14	NA	NA	NA
Empirical S.E.	.006	.002	.002	NA	NA	NA

Estimation of the  $\phi$ -coefficients is of little interest in this case because the selection model is misspecified, but notice the relative large bias in the estimation of the  $\beta$ -coefficients even though the population model is specified correctly. Thus, misspecifying the selection model affects the estimation of the population model.

Table 4. Estimation of population mean. (Mean true value =-0.00).

Method	Bias	Emp. Var. X 1000	MSE x 1000
$\hat{oldsymbol{Y}}_{U,Xknown}$	0.091	1.089	9.37
$\hat{oldsymbol{Y}}_{U,GREG}$	0.096	1.369	10.585
$\hat{ar{Y}}_{U,MAR}$	0.231	0.676	54.037

As can be seen, the bias, empirical variance and MSEs are much larger in this case than under the correct model (Table 2). This is not surprising since we fitted a wrong model. Here again, we applied the UNIF test for each simulation and obtained the following results.

	Mean	Standard Deviation	Minimum	Maximum
UNIF statistic	27.24	10.75	2.23	71.18
P-value	0.01	0.04	~0	0.82
H0 not rejected	0.04	0.17	0	1

For this case, the UNIF test performs well in rejecting the model fitted, with an average rejection rate of 96%.

# Case 2. The sample selection model is specified correctly, the population model is misspecified.

Here, we consider the case where the sample selection model is specified correctly (same as in Section 5.1), but the population model is misspecified. Specifically, the population values have been generated as  $y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i}^2 + \varepsilon_i$ , but the assumed working model is as in Section 5.1 (with  $x_{2i}$  instead of  $x_{2i}^2$ ). All the other model specifications are the same as in Section 5.1.

Table 5. Estimation of model coefficients with misspecified population model. 1,000 simulations.

	Population model coefficients			Selection model coefficients		
	$eta_0$ $eta_1$ $eta_2$		$\phi_0$	$\phi_{_{\! 1}}$	$\phi_2$	
True coefficients	-4	1	0.5	-2	1	0.5
Mean estimators	-5.77	0.94	2.20	-1.39	0.600	0.390
Empirical S.E.	0.013	0.002	0.002	0.040	0.011	0.007
Mean PWR estimator	-5.36	0.94	1.99	NA	NA	NA
Empirical S.E.	0.01	0.005	0.007	NA	NA	NA

As expected, the estimators of the  $\beta$ -coefficients are highly biased and so are the estimators of the  $\phi$ -coefficients. Thus, as already noted regarding Table 3, misspecification of one of the models affects the estimation of both models.

Table 6. Estimation of population mean. (Mean true value =-0.00).

Method	Bias	Emp. Var. X 1000	MSE x 1000
$\hat{oldsymbol{Y}}_{U,Xknown}$	-0.024	20.16	20.74
$\hat{oldsymbol{Y}}_{U,GREG}$	-0.010	42.03	42.11
$\hat{ar{Y}}_{U,MAR}$	-0.209	5.85	49.53

The estimators of the population mean are less biased than for the case where the sample selection model is misspecified (Table 4), but with relatively large variances, particularly for the GREG estimator. Notice that the GREG estimator depends directly on the estimated sample selection probabilities, which are highly biased (Table 5).

Application of the UNIF test yields in this case,

	Mean	Standard Deviation	Minimum	Maximum
UNIF statistic	207.12	46.49	82.12	394.47
P-value	~0.00	~0	~0	~0
H0 not rejected	0	0	0	0

The UNIF test rejects the models fitted in each of the 1,000 simulations.

## 5.4 Simulation setup with binary target variable- correct model

So far, we illustrated the performance of our proposed method for the case where the target y-variable is continuous. Following, we consider the case where y is binary. We use a similar simulation setup to the setup used for the continuous case, except that the population y-values are now generated as  $Pr(y_i = 1) = logit^{-1}(\beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i})$ , with the before. x-values generated We as again the logistic model use  $\pi_{i,S_{NP}} = \Pr(\delta_i = 1 \mid y_i, \mathbf{x}_i) = \frac{\exp(\phi_0 + \phi_1 x_{1i} + \phi_2 y_i)}{1 + \exp(\phi_0 + \phi_1 x_{1i} + \phi_2 y_i)} \quad \text{for selecting the } S_{NP}$ sample,

maximizing the likelihood under the same constraints as before.

The question arising is whether the  $S_{NP}$  model (4.2) is identifiable in this case as well. Wang et al. (2014) establish the following condition for model identifiability. The auxiliary variables  ${\bf x}$  in the population model can be decomposed as  ${\bf x}=({\bf x}_1,{\bf x}_2)$ , with the dimension of  ${\bf x}_2 \ge 1$ , such that  $\pi_{i,S_{NP}} = \Pr(\delta_i = 1 \mid y_i,{\bf x}_i) = \Pr(\delta_i = 1 \mid y_i,{\bf x}_{1i})$ , implying that the sample selection model does not depend on  ${\bf x}_2$ , given y and  ${\bf x}_1$ . This condition is

satisfied in our simulation setup. Recall that for a normal population model and logistic selection probabilities, the sample model is identifiable if the  $\mathbf{x}$  variables in the two models differ in at least one variable, a somewhat weaker condition. See Remark 16.

The results in Tables 7 and 8 below are based on 1,000 simulations with an average selection rate of 70%. The estimated value is again the true population mean (proportion) of the target y-variable.

Table 7. Estimation of model coefficients under correct models. 1,000 simulations.

	Population model coefficients			Selection model coefficients		
	$eta_0$ $eta_1$ $eta_2$			$\phi_0$	$\phi_{_{\! 1}}$	$\phi_2$
True coefficients	-4	1	1	-2	1	5
Mean estimators*	-4.40	1.18	1.01	-2.89	1.50	5.65
Empirical S.E.	0.01	0.004	0.002	0.016	0.008	0.085
PWR mean estimator	-0.24	0.20	0.16	NA	NA	NA
PWR Empirical S.E.	0.001	0.0005	0.0005	NA	NA	NA

<sup>\*</sup> The mean estimators of the model coefficients are of the MLE estimators. The PWR estimator is the probability weighted estimator with weights  $k_i = (1/\hat{\pi}_{i,S_{NP}})$ .

The MLE and PWR estimators are biased, notably the PWR estimator and the MLE estimators of the  $\phi$ - coefficients, but as can be seen in Table 8, the bias seems to have little effect on the estimation of the population mean of the target y-variable.

We consider the following estimators of the population mean:

$$\hat{\bar{Y}}_{U,H} = \sum\nolimits_{i \in S_{NP}} k_i y_i / \sum\nolimits_{i \in S_{NP}} k_i \; ; \; k_i = (1/\hat{\pi}_{i,S_{NP}}) \; .$$

$$\hat{\bar{Y}}_{U\!,E\!I} = \frac{1}{N} [\sum\nolimits_{i \in S_{NP}} y_i + \frac{N-n}{\sum\nolimits_{i \in S_{NP}} (k_i - 1)} \sum\nolimits_{i \in S_{NP}} (k_i - 1) y_i]; \text{ see Sverchkov and Pfeffermann}$$

(2004) for derivation of this estimator.

$$\hat{\bar{Y}}_{U,X \, known} = \frac{1}{N} \{ \sum_{i \in S_{NP}} y_i + \sum_{j \notin S_{NP}} \hat{E}_{S_{NP}} (y_j \mid \mathbf{x}_j) + \frac{N-n}{n} \sum_{i \in S_{NP}} \frac{k_i - 1}{\overline{k}_{S_{NP}} - 1} [y_i - \hat{E}_{S_{NP}} (y_i \mid \mathbf{x}_i)] \};$$

$$\overline{k}_{S_{NP}} = \frac{1}{n} \sum_{i \in S_{NP}} k_i$$
,  $\hat{E}_{S_{NP}}$  is the estimated expectation under the model (4.2). The estimator

when all the population  $\mathbf{x}$ 's are known. See Sverchkov and Pfeffermann (2004) for the derivation of this estimator.

$$\hat{\bar{\boldsymbol{Y}}}_{\!\!\!\boldsymbol{U}\!,\!\boldsymbol{GREG}} = \frac{\sum\limits_{i \in S_{NP}} k_i y_i}{\sum\limits_{i \in S_{NP}} k_i} + \hat{\boldsymbol{B}}_{pk}' [\boldsymbol{\overline{X}}_{\!\!\!\boldsymbol{U}} - \frac{\sum\limits_{i \in S_{NP}} k_i \boldsymbol{X}_i}{\sum\limits_{i \in S_{NP}} k_i}]; \text{ Same as when y is continuous.}$$

 $\hat{ar{Y}}_{U,M\!A\!R}=$  the estimator obtained from  $\hat{ar{Y}}_{U,E\!I}$  when replacing  $k_i$  by the MAR weight,

$$k_i^* = \frac{1 + \exp(\hat{\alpha}_0 + \hat{\alpha}_1 x_{1i} + \hat{\alpha}_2 x_{2i})}{\exp(\hat{\alpha}_0 + \hat{\alpha}_1 x_{1i} + \hat{\alpha}_2 x_{2i})}.$$

Table 8. Estimation of population mean. (Mean true value=0.5). 1,000 simulations.

Estimator	Bias	Emp. Var X1000	Emp. MSE x 1000*
			(Bootstrap estimate)*
$\hat{oldsymbol{Y}}_{U\!,H}$	-0.051	1.600	4.201 (5.60)
$\hat{oldsymbol{Y}}_{U\!,EI}$	0.001	0.009	0.010 (0.026)
$\hat{oldsymbol{Y}}_{U,Xknown}$	-0.006	0.169	0.205 (0.300)
$\hat{ar{Y}}_{U,GREG}$	-0.006	0.172	0.208 (0.309)
$\hat{oldsymbol{Y}}_{U,MAR}^{\widehat{oldsymbol{ au}}}$	0.149	0.049	22.25 (22.50)

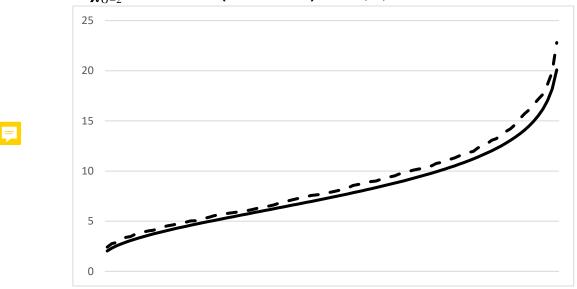
<sup>\*</sup> The bootstrap MSE estimates are based on 100 simulations with 100 bootstrap samples for each simulation.

As can be seen, all the estimators except  $\hat{Y}_{U,MAR}$  have a negligible bias, despite the bias of the estimated  $\phi$ -coefficients. Among the estimators,  $\hat{Y}_{U,EI}$  is the clear winner, with surprisingly small MSE, much lower than the MSE of  $\hat{Y}_{U,X\,known}$ . This might be due to the fact that this estimator uses the observed y's, (~70% in this case), and only predicts the sum of the unobserved y's. The estimator  $\hat{Y}_{U,X\,known}$  also uses the observed y's, but it uses

the estimated expectation under the  $S_{NP}$  model for predicting the sum of the unobserved y's. The estimator  $\hat{Y}_{UH}$  has a relatively large MSE due to its relatively larger bias.

**Model testing:** As for the continuous case, we tested the goodness of fit of our model, using in this case the Hosmer and Lemeshow (1980, H-L) test (Equation 5). Figure 2 compares the empirical cumulative distribution of the H-L statistic with the corresponding  $\chi^2_{G-2}$  distribution with G=10 groups.

Figure 2. Empirical cumulative distribution of H-L statistic (dashed curve) and under the  $\chi^2_{G-2}$  distribution (solid curve). *G*=10, 1,000 simulations.



Application of the test in the simulations with  $\alpha = 0.05$  significance level yields,

	Mean	Standard deviation	Minimum	Maximum
H-L test	8.56	4.30	1.108	30.57
p-value	0.46	0.29	~0	~1
H0 not rejected	0.934	0.248	0	1

The H-L test performs well when testing the correct model.

### 5.5 Application of proposed method for binary case with misspecified models

In Section 5.4, we assumed that the population model and the model for the selection probabilities are specified correctly. In this section we consider the case where they are misspecified, using the same simulation setup as before.

## Case 1. The population model is specified correctly, the sample selection model is misspecified.

In this case, we selected the  $S_{NP}$  sample with probabilities,  $\Pr(\delta_i = 1 \mid y_i, \mathbf{x}_i) = \frac{\exp(-2 + 5x_{1i}y_i)}{1 + \exp(-2 + 5x_{1i}y_i)} \,, \text{ but assumed as our working model the same}$ 

model as in Section 5.4. The population model of y is specified correctly. The average selection rate over the 1,000 simulations is in this case 54%.

Table 9. Estimation of model coefficients under misspecified model. 1,000 simulations.

	Population model coefficients			Selection model coefficients		
	$oldsymbol{eta}_0$	$oldsymbol{eta}_1$	$oldsymbol{eta}_2$	$\phi_0$	$\phi_{_{\mathrm{l}}}$	$\phi_2$
True coefficients	-4	1	1	NA	NA	NA
Mean estimators	-4.78	1.75	1.00	-3.1	0.96	3.1
Empirical S.E	0.01	.003	.003	0.02	0.003	0.14
PWR mean estimator	-0.16	0.25	0.14	NA	NA	NA
PWR empirical S.E	0.002	0.001	0.0004	NA	NA	NA

As can be seen, except for  $\beta_2$ , the MLE estimates of the other  $\beta$ -coefficients are biased, with larger bias than when the sample selection model was specified correctly (Table 7.)

Table 10. Estimation of population mean. (Mean true value=0.5).

Estimator	Bias	Emp. Var 1000	Emp. MSE x 1000
$\hat{m{Y}}_{U\!,H}$	0.04	1.60	3.2
$\hat{oldsymbol{Y}}_{U,EI}$	0.05	0.29	2.8
$\hat{oldsymbol{Y}}_{U\!,Xknown}$	0.11	0.53	12.6
$\hat{ar{Y}}_{U,GREG}$	0.11	0.53	12.6
$\hat{ar{Y}}_{U,MAR}$	0.29	0.17	84.3

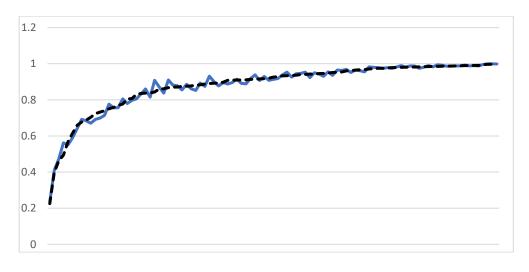
The results in Table 10 indicate that the first 2 estimators have small bias despite of the model misspecification, with smaller MSE of  $\hat{\overline{Y}}_{U,H}$ , but much larger MSEs of  $\hat{\overline{Y}}_{U,EI}$ ,  $\hat{\overline{Y}}_{U,X\;known}$  and  $\hat{\overline{Y}}_{U,GREG}$ , compared to the MSEs obtained under the correct model (Table 8). These large MSEs are clearly explained by the misspecification of the sample selection model. As before,  $\hat{\overline{Y}}_{U,MAR}$  has a large bias and an extreme MSE.

We applied the H-L test with  $\,\alpha=0.05\,$  significance level and obtained the following results.

	Mean	Standard deviation	Minimum	Maximum
H-L test	16.12	179.4	0.778	5557.8
p-value	0.43	0.30	~0	~1
H0 not rejected	0.89	0.312	0	1

It follows that the H-L test fails to reject the misspecified model in this case. In an attempt to understand this outcome, Figure 3 compares the correct  $S_{NP}$  model with true coefficients, used to select the sample with the corresponding estimated model under the misspecified model, for a simple random sample of 100 observations from the  $S_{NP}$  sample.

Figure 3. Comparison of Correct model and estimated misspecified model\*



<sup>\*</sup> Dashed curve represents the correct  $S_{\mathit{NP}}$  model, twisted curve represents the estimated  $S_{\mathit{NP}}$  model based on the misspecified model.

As can be seen, the estimated model under wrong specification yields almost perfect estimators of the correct model producing the  $S_{\it NP}$  data, which explains why the H-L test

does not reject the model. This is an example for what is known as "practical nonidentifiability" (Lee and Berger, 2001), meaning that even though the  $S_{NP}$  model is theoretically identifiable, another model may fit the data almost as well. Notice in Table 10 that the use of the misspecified working model yields two almost unbiased estimators of the true population mean.

## Case 2. The population model is misspecified, the sample selection model is specified correctly.

In this case, we used the same sample selection model as in Section 5.4 (correct specification of our working model), but we generated the population values as,  $y_i = \text{logit}(\beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i}^2)$ . As our working model we assumed the model of Section 5.4 ( $x_{2i}$ , instead of  $x_{2i}^2$ ). We used 1,000 simulations with an average selection rate of 0.73. All the other model specifications are as in Section 5.4.

Table 11. Estimation of model coefficients under misspecified model

	Population model coefficients			Selection model coefficients		
	$oldsymbol{eta}_0$	$oldsymbol{eta}_1$	$oldsymbol{eta}_2$	$\phi_0$	$\phi_{_1}$	$\phi_2$
True coefficients	-4	1	1	-2	1	5
Mean estimators	-6.07	1.34	1.62	-3.50	1.55	10.66
Empirical S.E	.015	.005	.002	.025	.01	.129
PWR mean estimator	-0.35	0.18	0.23	NA	NA	NA
PWR empirical S.E	.001	.001	.001	NA	NA	NA

As can be seen, all the estimators are highly biased, due to misspecification of the population model.

Table 12. Estimation of population mean (True mean value=0.55)

Estimator	Bias	Emp. Var 1000	Emp. MSE x 1000
$\hat{ar{Y}}_{U,H}$	-0.15	3.64	26.14
$\hat{ar{Y}}_{U,EI}$	-0.004	0.01	0.026
$\hat{ar{Y}}_{U,Xknown}$	-0.06	0.53	4.13
$\hat{ar{Y}}_{U,GREG}$	-0.05	0.58	3.08
$\hat{ar{Y}}_{U,MAR}$	0.11	0.05	12.15

All the estimators except for  $\hat{\bar{Y}}_{U,H}$  and  $\hat{\bar{Y}}_{U,MAR}$  have a negligible bias in this case, with  $\hat{\bar{Y}}_{U,EI}$  performing really well, as in the case of correct model specification (Section 5.4). On the other hand,  $\hat{\bar{Y}}_{U,X\;known}$ , although having a negligible bias, has a large MSE, even larger than the MSE of  $\hat{\bar{Y}}_{U,GREG}$ .

Application of the H-L test with  $\alpha = 0.05$  significance level yields in this case,

	Mean	Standard deviation	Minimum	Maximum
H-L test	42.5	41.3	6.09	950.0
p-value	0.002	0.02	~0	0.637
H0 not rejected	0.006	0.08	0	1

The H-L test performs well in rejecting the misspecified model.

## 6. CONCLUSION

In recent years, there is growing research on the use of NP samples for inference on population parameters, as an alternative or complement to the use of probability samples. A major problem with the use of these samples is their possible nonrepresentativeness of the corresponding target population, which if not accounted for properly, may lead to large bias in the inference process. In this article, we review and discuss several approaches proposed in the literature to deal with this problem, distinguishing between methods based on integration of the NP sample with an appropriate probability sample, and methods that

base the inference solely on the NP sample. Another distinction emphasized is between methods that assume that the selection to the NP sample depends on known auxiliary variables  $\mathbf{x}$ , but not on the target study y variable, and methods that assume that the selection depends also on y.

We also propose two additional methods for inference from a nonprobability sample, one that employs the empirical likelihood approach and one that requires specifying the population model parametrically. We discuss the conditions guaranteeing that the resulting model holding for the NP sample is identifiable, and propose simple tests for testing that the models are specified correctly. Our simulation study illustrates good performance of the proposed method and generally good performance of the test statistics.

A major problem underlying all the methods considered in this article is that they assume, at least implicitly, that every unit in the population has a positive probability to be in the NP sample. Clearly, if this is not the case, inference on the target population could be highly biased. This problem also exists with traditional probability samples when the sampling frame is not complete, known as under-coverage. When the group of units with zero probability to be included in the NP sample is known, say certain geographical areas, industries or ethnic groups, the target population should be redefined accordingly. When this is not the case, integration of the NP sample with an appropriate PS sample and the use of known population means of the true target population for calibration is a possible way to at least reduce the bias of the NP sample. This is an important topic for more research.

There are two important questions regarding our proposed method that require further investigation. The first question is how to proceed when the test statistic rejects the models defining the NP model. We do not have a clear answer to this question at this stage other than a scholarly consideration of alternative models. We mention again that the use of a logistic model for the selection probabilities has some theoretical justification, and this model is in common use.

The second related question is the choice of the  $\mathbf{x}$  variables in the models, when there are many of them. In practice, it may be the case that the analyst has a set of variables that he likes to include in the population model, which as explained in Section (4.1), defines also the variables included in the sample selection model. When this is not the case, one can use an appropriate stepwise algorithm. Beaumont et al. (2024) use a forward stepwise procedure, aimed at minimizing their proposed AIC criterion.

All the methods discussed in the present article should be considered as first attempts of inference from nonprobability samples, and more theoretical research and practical applications are required before they can be used routinely for the production of official statistics.

#### **REFERENCES**

AAPOR (2010), "Report on Online Survey Panels." available at <a href="http://poq.oxfordjournals.org/content/early/2010/10/19/poq.nfq048.full.html">http://poq.oxfordjournals.org/content/early/2010/10/19/poq.nfq048.full.html</a>

Baker, R., Brick, J. M., Bates, N. A., Battaglia, M., Couper, M. P., Dever, J. A., Gile, K. J. and Tourangeau, R. (2013). Report of the AAPOR Task Force on Non-probability Sampling. *Journal of Survey Statistics and Methodology*, 1, 90-143.

Beaumont, J.F. (2020). Are probability surveys bound to disappear for the production of official statistics? *Survey Methodology*, 46, 1-28.

Beaumont, J.F., Bosa, K., Brennan, A., Charlebois, J. and Chu, K. (2024a). Handling non-probability samples through inverse probability weighting with an application to Statistics Canada's crowdsourcing data. *Survey Methodology*, 50, 77-106.

Beaumont, J.F., Bosa, K., Brennan, A., Charlebois, J. and Chu, K. (2024b). Authors' response to comments on "Handling non-probability samples through inverse probability weighting with an application to Statistics Canada's crowdsourcing data." *Survey Methodology*, 50, 123-141.

Chen, Y., Li, P. and Wu, C. (2020). Doubly robust inference with non-probability survey samples. *Journal of the American Statistical Association*, 115, 2011-2021.

Chen, Y.; Li, P.; Rao, J.N.K and Wu, C. (2022). Pseudo empirical likelihood inference for non-probability survey samples. *Canadian Journal of Statistics*, 50, 1166-1185.

Citro, C. (2014). From multiple modes for surveys to multiple data sources for estimates. *Survey Methodology*, 40, 137-161.

Conti, P.L., Marella, D. and Scanu, M. (2008). Evaluation of matching noise for imputation techniques based on nonparametric local linear regression estimators. *Computational Statistics and Data Analysis*, 53, 354-365.

Elliott, M.R. and Valliant, R. (2017). Inference for Nonprobability Samples. *Statistical Science*, 32, 249-264.

Hosmer, D.W. and Lemeshow, S. (1980). A goodness-of-fit test for the multiple logistic regression model. *Communications in Statistics*, A10, 1043–1069.

Keiding, N. and Louis, T. A. (2016). Perils and potentials of self-selected entry to epidemiological studies and surveys. *Journal of the Royal Statistical Society*, Series A, 179, 319-376.

Kim, J.K. and Wang, Z. (2019). Sampling techniques for big data analysis. *International Statistical Review*, 87, 177-191.

Kim, J. K. and Morikawa, K. (2023). An empirical likelihood approach to reduce selection bias in voluntary samples. *Calcutta Statistical Association Bulletin*, 75, 8-27.

Kott, P. and Ridenhour, J. (2024). Calibration weighting with a blended (probability and nonprobability) sample: mean and variance estimation when errors can come from both samples. *Methods Report, RTI Press*.

Krieger, A.M. and Pfeffermann, D. (1997). Testing of Distribution Functions from complex Sample Surveys. *Journal of Official Statistics*, 13, 123-142.

Lee, J. and Berger, J. O. (2001). Semiparametric Bayesian analysis of selection models. *Journal of the American Statistical Association*, 96, 1397–1409.

Lee, S. (2006). Propensity score adjustment as a weighting scheme for volunteer panel Web survey. *Journal of Official Statistics*, 22, 329-349.

Lee, S. and Valliant, R. (2009), Estimation for volunteer panel web surveys Using propensity score adjustment and calibration adjustment, *Sociological Methods and Research*, 37, 319–343.

Marella, D. and Pfeffermann, D. (2023). Accounting for Non-ignorable Sampling and Non-response in Statistical Matching. *International Statistical Review*, 91, 269–293.

Molenberghs, G. Beunckens, C., and Kenward, M.G. (2008). Every missing not at random has a missingness at random counterpart with equal fit. *Journal of the Royal Statistical Society*, Series B, 70, 371–388.

Pfeffermann D. and Sverchkov M. (1999). Parametric and semiparametric estimation of regression models fitted to survey data. *Sankhyā*, *Series B*, 61, 166–186.

Pfeffermann, D. and Landsman, A. (2011). Are private schools really better than public schools? Assessment by methods for observational studies. *Annals of Applied Statistics*, 5, 1726–1751.

Pfeffermann, D. and Sikov, A. (2011). Imputation and estimation under non ignorable nonresponse in household surveys with missing covariate information. *Journal of Official Statistics*, 27, 181–209.

Pfeffermann, D., Marella, D. and Summa, D. (2025). Matching of a Non-probability sample with a probability sample affected by nonignorable sampling and nonresponse. Submitted for publication.

Pfeffermann, D., Preminger, A. and Sikov, S. (2025). Statistical inference under nonignorable sampling and nonresponse - an empirical likelihood approach. (Under revision).

Qin J., Leung, D. and Shao J. (2022). Estimation with survey data under non-ignorable nonresponse or informative sampling. *Journal of the American Statistical Association*, 97, 193–200.

Rao, J.N.K. (2021). On making valid inferences by integrating data from surveys and other sources. *Sankhyā*, *Series B*, 83, 242-272.

Rivers, D. (2007). Sampling from web surveys. *Proceedings of the Survey Research Methods Section*, American Statistical Association, Alexandria, VA.

Sayag, D., Ben-Hur, D. and Pfeffermann, D. (2022). Reducing revisions in hedonic house price indices by the use of nowcasts. *International Journal of Forecasting*, 38, 253–266.

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

Wang, S., Shao, J. and Kim, J. K. (2014), An Instrument Variable Approach for Identification and Estimation with Nonignorable Nonresponse. *Statistica Sinica*, 24, 1097–1116.

Wu, C. (2022). Statistical inference with non-probability survey samples. *Survey Methodology*, 48, 283-311.

Zacks, S. (1971) *The theory of statistical inference*. (Wiley series in probability and mathematical statistics)