

Kernel Weighting for blending probability and non-probability survey samples

María del Mar Rueda¹, Beatriz Cobo², Jorge Luis Rueda³, Ramón Ferri-García⁴
, Luis Castro-Martín⁵

Abstract

In this paper we review some methods proposed in the literature for combining a non-probability and a probability sample with the purpose of obtaining an estimator with a smaller bias and standard error than the estimators that can be obtained using only the probability sample. We propose a new methodology based on the kernel weighting (KW) method. We discuss the properties of the new estimator when there is only selection bias and when there are coverage and selection biases. We perform an extensive simulation study to better understand the behavior of the proposed estimator.

¹Department of Statistics and Operational Research, University of Granada, Spain

²Department of Quantitative Methods for Economics and Business, University of Granada, Spain

³Department of Statistics and Operational Research, University of Granada, Spain

⁴Department of Statistics and Operational Research, University of Murcia, Spain

⁵Andalusian School of Public Health, Spain

MSC:62D05

Keywords: Kernel weighting, survey sampling, non-probability sample, coverage bias, selection bias

1. Introduction

Probability sampling methods are well established by statistical offices and researchers as one of the primary tools for data collection in surveys. This is because when controlling the sampling design, it is feasible to make valid statistical inference about large finite populations using relative small samples. There exists an extensive literature on methods for probability sampling and design-based inferences for complex surveys.

However, the deployment of probability sampling methods has become more challenging, as there has been a notorious decline in response rates [31, 21] with the subsequent increase of the survey costs. In addition, new data sources which have arisen in recent years could be considered as alternatives to survey data. Examples are large volume datasets coming from sources such as passive data or "data lakes", and web surveys that have the potential of providing more timely estimates, as well as offering easier data access and lower data collection costs than traditional probability sampling, leading to larger sample sizes. On the other hand, there are serious issues concerning the use of non-probability survey samples for estimation. The primary issue with these data sources is that the selection mechanism, which decides what individuals are eventually included in the dataset, is often unknown and may induce serious coverage and selection biases. The generalization of the results under these biases is therefore compromised.

Despite these limitations, non-probability survey designs may be particularly useful in several cases. For example, they can be used in those cases where the target population is a small subpopulation unlikely to meet sample size requirements, or when we are interested in non-demographical strata which cannot be considered in a sampling design.

Given the potential of non-probability surveys, statisticians have studied the integration or combination of data from probability and non-probability samples. Some reviews on methods of statistical data integration for finite population inference can be consulted in [2], [45], [50] or [37]. Different data integration methods, which are based on combining probability and non-probability samples, have been recently developed in the literature on survey sampling. These integration methods can be divided into three groups depending on the availability of the study variable: available in the non-probability sample only, in the probability sample only, or in both samples.

Many methods consider the first case, where the target variable has been observed in the non-probability sample only. In this situation, the probability sample plays an important role as the reference data, and can be used to increase the efficiency of the estimates through a variety of adjustment approaches to account for the selection bias in non-probability samples. However, other methods were also developed from different perspectives according to the availability of auxiliary information. Calibration [8, 14], Propensity Score Adjustment (PSA) [28, 29, 5], kernel weighting (KW) [47], Statistical Matching [38, 1], double robust estimation [24] and superpopulation modeling [46, 2] are relevant techniques to mitigate selection bias.

When the non-probability (or volunteer) survey contains auxiliary variables but no study variable, [35] shows how the use of a non-probability database can improve estimates from a probability sample and they define a class of QR predictors [42] asymptotically design-unbiased under certain conditions.

In this paper we consider the third situation posed above, where the study variables are measured in both samples. In Section 2 we review the estimation from probability and non-probability samples to introduce the notation and the framework. In Section 3 we revisit some important works in data integration for handling selection bias in our context. In Section 4 we adapt the kernel weighting method introduced in [47], to data

integration. First, we consider a situation where there are no coverage biases, and we propose a KW estimator by a linear combination of biased and unbiased estimators of a population mean. When there is coverage bias in the non-probability sample, as is usual in practice, we propose a KW estimator based on dual frame methodology. We derive conditions such that these proposed estimators are asymptotically design-unbiased. In Section 5, we use Monte Carlo simulations to compare the proposed method with several models and show that the kernel weighted estimator is a good compromise for several setups. Finally we conclude and give perspectives in 6.

2. Context and notation

Let U be the target population of size N , $U = \{1, \dots, i, \dots, N\}$. Let s_v be the set of n_v units selected from the frame U_v using a non-probability (volunteer) data collection method. Let s_r be a probability sample of size n_r selected from a frame U_r under the sampling design $d = (s_r, p_r)$ with $\pi_i > 0$ the first order inclusion probability for individual i and π_{ij} the second order probabilities for individuals i and j . Let be $d_i = 1/\pi_i$ the sampling design weight of unit i . We consider a situation in which U_r and U_v coincide with the population under study U . That is, there are no coverage biases in either the probability or the non-probability sample.

Let us denote with y_i the collected value on the unit i for the target variable y and let \mathbf{x}_i be the observed values for individual i for a vector of covariates \mathbf{x} . Both y and \mathbf{x} have been measured in both samples.

The target parameter is the population mean, $\bar{Y} = \frac{1}{N} \sum_U y_i$, that can be estimated from the probability sample using the Horvitz-Thompson estimator:

$$\bar{y}_r = \frac{1}{N} \sum_{i \in s_r} d_i y_i \quad (1)$$

and from the volunteer sample with the naive estimator:

$$\bar{y}_v = \sum_{i \in s_v} \frac{y_i}{n_v} \quad (2)$$

If there is full response in s_r , the estimator \bar{y}_r is unbiased but if the sample size is small it can lead to estimates with large sampling errors.

Let us consider the variable

$$I_{vi} = \begin{cases} 1 & i \in s_v \\ 0 & i \in U - s_v \end{cases}, \quad i = 1, \dots, N. \quad (3)$$

the estimator \bar{y}_v is biased [24] and its bias is given by

$$E_v(\bar{y}_v - \bar{Y}_N) = \frac{1}{f_v} E_v \text{Cov}(I_v, y)$$

where $E_v(\cdot)$ denotes the expectation under the selection mechanism model of the non-probability sample and $f_v = n_v/N$. Thus the mean squared error, MSE , is given by the formula

$$MSE(\bar{y}_v) = \frac{1}{f_v^2} E_v(\text{Corr}(I_v, y)^2) \text{Var}(I_v) \text{Var}(y).$$

Therefore, a non-probability sampling where $E_v\{\text{Corr}(I_v, y)\} \neq 0$ induces a certain selection bias to the results.

In the next section we will consider how we can estimate the mean population by using a data integration estimator that combine information for these two independent surveys.

3. Methodology in data integration for handling selection bias

3.1. Some previous works

Starting with the work of [11], these authors consider the problem of combining the two samples by means of a linear combination of the biased and unbiased estimators of the population mean:

$$\bar{y}_{com} = \alpha \bar{y}_r + (1 - \alpha) \bar{y}_v$$

The best estimator, in terms of efficiency, of this combination when the magnitude of the bias is known is given by:

$$\hat{y}_{EH} = \frac{\bar{y}_v \frac{\sigma_r}{n_r} + \bar{y}_r (B^2 + \frac{\sigma_v}{n_v})}{B^2 + \frac{\sigma_r}{n_r} + \frac{\sigma_v}{n_v}} \quad (4)$$

being \bar{y}_v and \bar{y}_r the sample means, with variances $\frac{\sigma_v}{n_v}$ and $\frac{\sigma_r}{n_r}$ and B the bias of \bar{y}_v .

In practice, the bias and variances have to be estimated using the information available from both samples. The bias can be estimated as the difference between the sample means of both samples. In addition, the authors calculate the maximal contribution of the non-probability sample in terms of effective sample size, the role of the non-probability sample size in approaching this limit and the roles of both sample sizes in estimating bias with enough precision. They show that a large probability sample size (1000-10000) is needed for reasonably precise estimates of the remaining bias in initially bias-adjusted non-probability sample estimators.

Other important work is due to [9]. Their proposal, based on calibration weighting, considers that auxiliary variables needed for calibration weighting must reliably differentiate between the probability sample and the non-probability sample. This calibration method has four steps:

1. Authors do a post-stratification raking calibration of s_r , using a set of demographic and geographical variables.
2. They combine the weighted s_r with the unweighted s_v . The combined sample is then weighted according to the probability sample's benchmarks from the previous

step.

3. They compare the answers from early adopter questions between the probability sample from step 1 to the answers from the blended sample from step 2.
4. They select some minimum number of early adopter questions to include in the raking due in Step 2.

Therefore, this procedure requires a good selection of early adopter questions that are included in the two surveys and that we believe will help to differentiate the samples.

Recently, [23] developed two estimators using combined data from probability sampling and non-probability sampling based on the total decomposition:

$$Y = Y_v + Y_c$$

where $Y_v = \sum_{i \in s_v} y_i = \sum_{i \in U} I_{vi} y_i$ and $Y_c = \sum_{i \in U - s_v} y_i = \sum_{i \in U} (1 - I_{vi}) y_i$. Since y is measured for all units of non-probability sampling, Y_v is known. Therefore, we only have to estimate Y_c . Authors proposed a first estimator where Y_c is estimated using the expansion estimator based on the probability sample

$$\bar{y}_{DI} = \frac{1}{N} (Y_v + \sum_{i \in s_r} d_i (1 - I_{vi}) y_i).$$

In Poisson sampling, the variance of \bar{y}_{DI} is smaller or equal to the variance of \bar{y}_r if a condition on the study variable for simple random sampling without replacement holds. When N is known, [23] propose to improve the previous estimator using the following one:

$$\bar{y}_{PDI} = \frac{1}{N} \left(Y_v + (N - n_v) \frac{\sum_{i \in s_r} d_i (1 - I_{vi}) y_i}{\sum_{i \in s_r} d_i (1 - I_{vi})} \right).$$

Authors prove that the variance of \bar{y}_{PDI} is smaller than the variance of \bar{y}_r for simple ran-

dom sampling. They also discuss how to improve the efficiency of this data integration estimator by using ratio and calibration estimation.

Other works in this matter are briefly introduced below.

[13] improve the blended calibration estimator provided by [9]. [10] develop pseudo-weights to create a representative sample using data from the non-probability sample under model assumptions that can be partially tested. With this approach, probability and non-probability samples can be blended, and the resulting sample can be treated as a probability sample with these new pseudo-weights. [49] consider a Bayesian approach for integrating a small probability sample with a non-probability sample. They show that considering informative priors based on non-probability data can reduce the variance and mean squared error of the coefficients of a linear model.

Recently, [48] do an extensive simulation study for comparing various weighting strategies where probability and non-probability samples are combined with weight normalization and raking adjustment. They apply these methods to a teen smoking behaviour survey. [36] consider the case of estimating proportions when a non-probabilistic sample and scraped data are available. They carry out a simulation study in which they evaluate the behavior of several estimators that make up the probability sample and these data sources in the Information and a Communication Technology survey. Some important works [39, 40] have appeared in which probability and non-probability samples are combined based on the propensity score adjustment technique. In the next section, we explain this technique and how it has been used by these authors.

3.2. Some estimators based on propensity score adjustment

The key concept in a non-probability survey sample is the selection mechanism. This mechanism is usually unknown and requires a suitable prediction model for the inclusion indicator variable. In this context, propensity scores, π_{vi} , can be defined as the

probability of the i -th individual of being included in the sample, $P(I_{vi} = 1)$, given the characteristics of the unit.

Let \mathbf{x} a matrix of covariates measured in s_v and also in s_r . We make the following assumption:

Assumption 1 (strong ignorability condition): the indicator variable I_v and the study variable y are conditionally independent given \mathbf{x} ; i.e. $P(I_v = 1|\mathbf{x}, y) = P(I_v = 1|\mathbf{x})$.

We assume that the selection mechanism of s_v verifies Assumption 1 and follows the model:

$$\pi_{vi} = P(I_{vi} = 1|\mathbf{x}_i) = p_i(\mathbf{x}) = m(\gamma, \mathbf{x}_i) \quad i = 1, \dots, N \quad (5)$$

where $m(\cdot)$ is a given function with second continuous derivatives with respect to γ .

We aim to estimate propensity scores using data from both samples. The maximum likelihood estimator of π_{vi} is $m(\hat{\gamma}, \mathbf{x}_i)$ where $\hat{\gamma}$ maximizes the pseudo-likelihood [7]:

$$\tilde{l}(\gamma) = \sum_{s_v} \log \frac{m(\gamma, \mathbf{x}_i)}{1 - m(\gamma, \mathbf{x}_i)} + \sum_{s_r} \frac{1}{\pi_i} \log(1 - m(\gamma, \mathbf{x}_i)). \quad (6)$$

The estimated propensities $\hat{\pi}_{vi} = m(\hat{\gamma}, \mathbf{x}_i)$ are thus used to readjust the propensity bias of the volunteer sample.

Based on these propensities, [39] define several estimators integrating the two samples.

A first estimator is calculated weighting estimators from each sample:

$$\bar{y}_{RDR1} = \alpha_1 \bar{y}_r + (1 - \alpha_1) \bar{y}_v \quad (7)$$

where $\bar{y}_v = \frac{1}{N} \sum_{s_v} y_i / q_i$ with $q_i = \frac{\pi_i \hat{\pi}_{vi}}{1 - \hat{\pi}_{vi}}$ and $\alpha_1 = \frac{(\sum_{s_r} \pi_i^{-1})(\sum_{s_v} \hat{\pi}_{vi}^{-2})}{(\sum_{s_r} \pi_i^{-1})(\sum_{s_v} \hat{\pi}_{vi}^{-2}) + (\sum_{s_r} \pi_i^{-2})(\sum_{s_v} \hat{\pi}_{vi}^{-1})}$.

For the second estimator, the authors calculate the values $p_i = \pi_i / (1 - \hat{\pi}_{vi})$ for all indi-

viduals in the joined $s = s_v \cup s_r$ and obtain a simple Horvitz-Thompson type estimator with these new weights:

$$\bar{y}_{RDR2} = \frac{1}{N} \sum_s y_i / p_i \quad (8)$$

Let \mathbf{x} be a set of auxiliary variables, related to y , whose population totals are known.

Two calibration estimators are also proposed:

$\bar{y}_{RDR3} = \frac{1}{N} (\sum_{s_v} y_i * w_{1i} + \sum_{s_r} y_i * w_{2i})$ where w_{1i} and w_{2i} are as close as possible to $1/p_i$ fulfilling $T_x = \sum_{s_v} w_{1i} \mathbf{x}_i = \sum_{s_r} w_{2i} \mathbf{x}_i$ and the estimator:

$$\bar{y}_{RDR4} = \alpha_2 \bar{y}_r + (1 - \alpha_2) \bar{y}_v \quad \text{being} \quad \alpha_2 = \frac{(\sum_r w_{1i})(\sum_v w_{2i}^2)}{(\sum_r w_{1i})(\sum_v w_{2i}^2) + (\sum_r w_{1i}^2)(\sum_v w_{2i})}.$$

[40] propose the combined estimator:

$$\bar{y}_{CPSA} = \alpha_0 \bar{y}_r + (1 - \alpha_0) \bar{y}_{IPW} \quad (9)$$

being $\bar{y}_{IPW} = \frac{1}{N} \sum_{s_v} y_i / \hat{\pi}_{vi}$, and $\alpha_0 = \frac{\hat{V}_2}{\hat{V}_1 + \hat{V}_2}$ where \hat{V}_1 and \hat{V}_2 are estimators of the variance of \bar{y}_r and the MSE of \bar{y}_{IPW} respectively. They also propose alternative methods that combine propensity score adjustment and calibration using machine learning predictive algorithms.

[3] consider a few ways on non-probability integration by combining generalized difference estimator and post-stratified calibration estimator with the inverse probability weighted estimating for estimating proportions in the survey on population by religion, native language and ethnicity in Lithuania.

The above methods can reduce bias by using propensity scores to estimate participation rates of non-probability sample units. However, they are sensitive to propensity model misspecifications and can largely increase the variance of the estimators due to extreme weights. A possible way to reduce the effect of extreme weights is the KW method [47] that uses propensity scores as a measure of similarity, and therefore is less sensitive to model misspecification while avoiding the extreme weights that may be produced in propensity score estimation. In the next section we introduce the KW approach

to create pseudoweights for the non-probability sample and propose a new method of integration based on this KW estimator.

4. Proposed estimators based on kernel weighting

The kernel weighting method was developed by [47], and is a method similar to the PSA since both consist of creating pseudoweights for the non-probability sample using auxiliary variables of a reference probability sample. However, what differentiates them is the way in which these new weights are generated, although as in PSA we will use the estimated propensities to participate in the survey. As it occurred in that case, these propensities can be estimated in different ways, even though the most commonly used one is by means of logistic regression models. However, machine learning (ML) techniques can be used in order to estimate these propensities [15].

The KW is based on using these propensities to measure the similarity between individuals based on the distributions of the auxiliary variables of the reference sample s_r and the non-probability sample s_v . These similarities will be used as weights for our estimator, after smoothing the distances using kernel functions.

The estimated probability of inclusion for individual $i \in s_v$ is obtained as

$$\hat{\pi}_{vi}^* = E_M[\delta_i^* = 1 | \mathbf{x}_i], \quad i \in s_v \cup s_r,$$

and for the individual $j \in s_r$ as:

$$\hat{\pi}_{rj}^* = E_M[\delta_j^{**} = 1 | \mathbf{x}_j], \quad j \in s_v \cup s_r$$

where, M will be one of the mentioned ML models to estimate this propensity and

$$\delta_i^* = \begin{cases} 1 & \text{for } i \in s_v \\ 0 & \text{for } i \in s_r \end{cases}$$

$$\delta_j^{**} = \begin{cases} 1 & \text{for } j \in s_r \\ 0 & \text{for } j \in s_v \end{cases}$$

Once we have these estimated propensities, we will calculate the distance between the two individuals belonging to the different samples. We define this distance as:

$$d_{ij} = \hat{\pi}_{vi}^* - \hat{\pi}_{rj}^*, \quad i \in s_v, \quad j \in s_r$$

This distance between individuals will have a value between -1 and 1. We seek to smooth these values, which is why we use a kernel function centered at zero. There are many alternative kernel functions that can be used (normal function, standard normal, triangular, etc.), see [43]. The closer this distance is to zero, the more similar the individuals are with respect to their auxiliary variables (propensities are estimated depending on the values of the auxiliary variables). Moreover, the more similar the individuals are, the greater the proportion that the KW will assign to the original weight of the reference sample d_{kj} to the i unit of the volunteer sample. This proportion is called the kernel weight, whose expression is as follows:

$$k_{ij} = \frac{K\{d_{ij}/h\}}{\sum_{i \in s_v} K\{d_{ij}/h\}}, \quad i \in s_v, \quad j \in s_r$$

where $K\{\cdot\}$ is a zero-centered kernel function [12], and h is the bandwidth corresponding to that kernel function. In addition:

$$\sum_{i \in s_v} k_{ij} = 1, \quad k_{ij} \in [0, 1]$$

The larger the value of the kernel weight k_{ij} is, the more similar the propensities will be among individuals $i \in s_v$ and $j \in s_r$.

Once we have the kernel weights, the pseudo weights KW can be calculated, w_i^{KW} for $i \in s_v$ which are the sum of the weights of the reference sample d_j , where $j \in s_r$, weighted by the kernel weights k_{ij} for the unit $i \in s_v$:

$$w_i^{KW} = \sum_{j \in s_r} d_j k_{ij}, \quad i \in s_v, \quad j \in s_r$$

Therefore a KW estimator for the population mean is:

$$\bar{y}_{KW} = \frac{1}{N} \sum_{i \in s_v} w_i^{KW} y_i$$

where $\sum_{i \in s_v} w_i^{KW} = \sum_{j \in s_r} d_j$, because of $\sum_{i \in s_v} k_{ij} = 1$.

The KW estimator is consistent if certain regularity conditions are met (see Appendix 1).

[22] improve the KW method by pairing it with Machine Learning. They consider conditional random forests, model-based recursive partitioning, gradient tree boosting and model-based boosting for estimating the propensities and constructing pseudo-weights.

4.1. Blending the samples with KW

First, we consider the situation where there is no coverage bias (U_r and U_v are equivalent to the population under study U). In this situation we propose a class of estimators based on both samples:

$$\bar{y}_C = \alpha \bar{y}_r + (1 - \alpha) \bar{y}_{KW} \tag{10}$$

where α is a nonnegative constant such that $0 \leq \alpha \leq 1$.

We study the asymptotic properties of the proposed estimator under the framework of [20] in which the properties of estimators are established under a given sequence of populations and a corresponding sequence of random sampling designs.

Theorem 1. Under assumption given in Appendix 1, the proposed estimator $\bar{y}_C \rightarrow Y$ in probability as $N \rightarrow \infty$, $n_v \rightarrow \infty$, $n_r \rightarrow \infty$ with $\frac{n_v}{N} = O(1)$ and $\frac{n_r}{N} = O(1)$.

Proof.

Assumptions 1a and 2a give sufficient conditions for the Horvitz-Thompson estimator \bar{y}_R to be consistent [20]. Under assumptions 2a-2c [47] proves that $\bar{y}_{KW} \rightarrow \bar{Y}$ in probability. Being a linear combination of consistent estimators, the proposed estimator is also consistent and converges to \bar{Y} .

Now, we consider the problem of how select the α parameter. A simple selection for α is to weight each estimator by the weight that sample has in the total sample so that $\alpha_n = n_r/(n_r + n_v)$. We propose other solution based on the idea of minimizing the MSE of the estimator:

$$\bar{y}_{CO} = \frac{\widehat{MSE}(\bar{y}_{KW})}{\widehat{MSE}(\bar{y}_{KW}) + \hat{V}(\bar{y}_r)} \bar{y}_r + \frac{\hat{V}(\bar{y}_r)}{\widehat{MSE}(\bar{y}_{KW}) + \hat{V}(\bar{y}_r)} \bar{y}_{KW} \quad (11)$$

being $\hat{V}(\bar{y}_r)$ the Horvitz-Thompson estimator of $V(\bar{y}_r)$ and $\widehat{MSE}(\bar{y}_{KW})$ an estimator for the $MSE(\bar{y}_{KW})$.

An estimator for the variance of \bar{y}_{KW} can be obtained by using the jackknife method proposed by [47] and the bias of this estimator can be estimated by $\bar{y}_r - \bar{y}_{KW}$.

4.2. Blending the samples with coverage bias

Web and social media surveys usually have a significant under-coverage bias. Thus, we consider now a more realistic situation where there is also under-coverage bias in the

non-probability sample. We will consider that U_r covers the entire finite population but the frame U_v be incomplete ($U_v \subset U$). The population of interest, U , may be divided into two mutually exclusive domains, $ab = U_v$ and $a = U \cap U_v^c$. Units in s_r can be divided as $s_r = s_{ra} \cup s_{rab}$, where $s_{ra} = s_r \cap a$ and $s_{rab} = s_r \cap (ab)$.

Following Hartley's idea [19], we can obtain a combined estimator of \bar{Y} by weighting the estimators obtained from each sample:

$$\bar{y}_H(\eta) = \frac{1}{N}(\hat{Y}_a + \eta \hat{Y}_{ab} + (1 - \eta) \hat{Y}_{KW}) \quad (12)$$

where $\hat{Y}_a = \sum_{i \in s_{ra}} d_i y_i$, $\hat{Y}_{ab} = \sum_{i \in s_{rab}} d_i y_i$ and $\hat{Y}_{KW} = \sum_{i \in s_v} w_i^{KW} y_i$ and $0 < \eta < 1$

Now, we denote as:

$$d_i^\circ = \begin{cases} d_i & \text{if } i \in s_{ra} \\ \eta d_i & \text{if } i \in s_{rab} \\ (1 - \eta) w_i^{KW} & \text{if } i \in s_v \end{cases} \quad (13)$$

then

$$\bar{y}_H(\eta) = \frac{1}{N} \sum_{i \in s} d_i^\circ y_i.$$

Theorem 2. Under the regularity conditions given in [47] for the sampling design and the propensity scores, the Hartley estimator $\bar{y}_H(\eta)$ is asymptotically unbiased for \bar{Y} .

Proof.

Since each domain is estimated by its Horvitz-Thompson estimator, $\hat{Y}_a + \eta \hat{Y}_{ab}$ is an unbiased estimator of $\sum_{i \in a} y_i + \eta \sum_{i \in ab} y_i$, for a given η . Under the regularity conditions given in [47] the estimator \hat{Y}_{KW} is asymptotically unbiased for $Y_{ab} = \sum_{i \in ab} y_i$, thus the estimator $\bar{y}_H(\eta)$ is asymptotically unbiased for \bar{Y} .

As U_r and U_v are sampled independently, the asymptotic variance of $\bar{y}_H(\eta)$ is given by

$$V(\bar{y}_H(\eta)) = \frac{1}{N^2} (V(\hat{Y}_a + \eta \hat{Y}_{ab}) + V((1 - \eta) \hat{Y}_{KW})) = \frac{1}{N^2} (V(\hat{Y}_a) + \eta^2 V(\hat{Y}_{ab}) + (1 - \eta)^2 V(\hat{Y}_{KW})) \quad (14)$$

where $V(\hat{Y}_a)$ and $V(\hat{Y}_{ab})$ are computed under the sampling design $d = (s_r, p_r)$ and $V(\hat{Y}_{KW})$ under the propensity model π_v .

The choice of the value for η is an important issue. For a fixed value of η , the estimator is simple to implement and gives internal consistency given that the same set of adjusted weights is used for all variables. The value of $\eta = 0.5$ is frequently used in dual frame estimation [34]. The value of η that minimizes the asymptotic variance in 14 is:

$$\eta_o = \frac{MSE(\hat{Y}_{KW}) - cov(\hat{Y}_a, \hat{Y}_{ab})}{V(\hat{Y}_{ab}) + MSE(\hat{Y}_{KW})} \quad (15)$$

This value depends on unknown population variances and covariances. By substituting the variances and MSE for its sample based estimators we obtain an estimator that we denote by $\bar{y}_H(opt)$. We note that these modified weights are random variable and their variability needs to be accounted for in standard errors of estimators.

A simple multiplicity estimator can be also obtained extending the multiplicity-adjusted methodology proposed in [34]. We denote by m_i the number of frames in which every unit is included. Thus $\bar{Y} = \frac{1}{N} (\sum_{U_v} y_i m_i^{-1} + \sum_{U_r} y_i m_i^{-1})$ and we can propose a single frame estimator:

$$\bar{y}_{MA} = \frac{1}{N} \left(\sum_{i \in s_v} y_i w_i^{KW} m_i^{-1} + \sum_{i \in s_r} y_i d_i m_i^{-1} \right). \quad (16)$$

It is easy to see that the estimator \hat{Y}_{MA} coincides with the estimator $\hat{Y}_H(\eta)$ for $\eta = 0.5$.

5. Simulation studies

We have conducted a simulation study to compare the efficiency of some of the proposed estimators based on KW. We are interested in comparing those estimators with some alternative estimators defined in Section 3, in the effect of the machine learning algorithm used in KW, in the effect of the kernel function used in the construction of KW pseudo-weights and also in the effect of considering coverage bias. In order to illustrate that the superiority of some estimators compared to others depends on the data, we define different setups based on different artificial populations and different sampling strategies.

5.1. Populations and setups

We consider a finite population of size $N = 500000$. The variables of interest were designed with the objective of having various types of relationships with the covariates and the propensities. We consider 8 auxiliary variables x , 2 variable of interest y and a variable π_{vi} which indicates the probability of being included in the non-probability sample. All of them were simulated as follows:

1. The covariates x_1, x_3, x_5 and x_7 followed a Bernoulli distribution with $p = 0.5$, and x_2, x_4, x_6 and x_8 followed Normal distributions with standard deviation of one and a mean parameter of 0 or 2, depending on the value of the previous Bernoulli variable. That is to say, in order to calculate x_2 we relied on the variable x_1 and if this variable was equal to 1, then the mean would be 2, or if the variable was equal to 0, then the mean would be 0. The same procedure was followed for the rest of the variables. The propensity models were fitted using all of the 8 auxiliary variables.

2. The target variables were created in order to have different relationships with the covariates and the propensities were simulated according to the formulas:

$$y_{1i} = N(8, 2) + 3(x_{5i} = 1) + 5\pi_i, \quad i \in U$$

$$y_{2i} = \begin{cases} 1 & \text{if } y_{1i} > \text{Median} \\ 0 & \text{if } y_{1i} \leq \text{Median} \end{cases}, \quad i \in U \quad (17)$$

3. The non-probability samples were drawn with a Poisson sampling design where the probability depends on variables x_5, x_6, x_7 y x_8 as:

$$\ln\left(\frac{\pi_{vi}}{1 - \pi_{vi}}\right) = -0.5 + 2.5(x_{5i} = 1) + \sqrt{2\pi}x_{6i}x_{8i} - 2.5(x_{7i} = 1), \quad i \in U. \quad (18)$$

We considered three setups. In the first setup the probability sample was drawn by simple random sampling without replacement (SRSWOR) from the full population; in the second setup the probability sample was drawn with stratified random sampling by the auxiliar variable x_7 and considering an allocation by strata of 1/3 and 2/3; in the third setup, the probability sample was selected with Midzuno sampling where the probabilities were proportional to a variable following a Normal distribution with a mean parameter dependent on the value of the auxiliar variable x_7 and a standard deviation of 0.5.

The aim of the described selection mechanism was to create weights with large variability. As a result, the mean propensity is 0.7050, with a standard deviation of 0.3792, and thus a coefficient of variation of 0.5379. The histogram of propensities $\pi_{vi}, i \in U$, is provided in figure 1.

5.2. The simulation procedure

The first simulation study evaluates the performance of some estimators for \bar{Y} there is selection bias in the estimates. We focused on the proposed estimator discussed in the

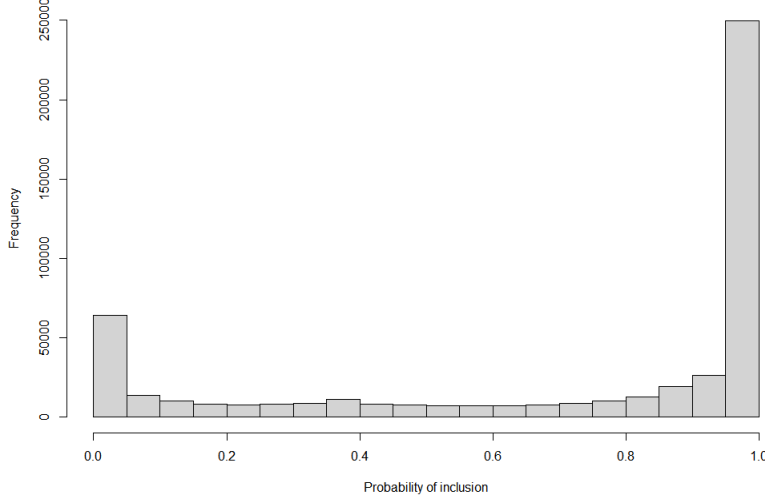


Figure 1. *Histogram of the population propensities*

paper, \bar{y}_{CO} , and we compared it with others estimators based on propensities. As a reference estimator we have considered the naive estimator that weights the estimators simply by their sizes $\bar{y}_{REF} = \frac{n_r}{N_r}\bar{y}_r + \frac{n_v}{N_v}\bar{y}_v$. We also evaluate the estimators \bar{y}_{RDR1} (7), \bar{y}_{RDR2} (8) and \bar{y}_{CPSA} (9) that do not use calibration.

We considered the XGBoost [6] algorithm among several machine learning approaches for estimating the propensities in all estimators. This algorithm builds decision trees ensembles that optimize an objective function via Gradient Tree Boosting [18]. Literature shows that PSA with Gradient Boosting Machines provides better results than other machine learning approaches [26, 27, 33, 32, 16, 41]. The method depends on several hyperparameters for a proper functioning and in order to avoid overfitting. We have considered the following hyperparameters: the number of trees forming the ensemble (50, 100 or 150), the weight shrinkage applied after each boosting step (0.3 or 0.4), the maximum number of splits that each tree can contain (1, 2 or 3), the proportion of variables used in each step (0.6 or 0.8) and the proportion of data used in each step (0.5, 0.75 or 1).

For each setup we select 500 probability samples of size $n_r = 250$ and 500 non-probability samples of sizes $n_v = 500; 1000; 2000$. We compute the Monte Carlo relative bias of the estimators:

$$|RB| = \frac{1}{B} \sum_{i=1}^B \frac{|\bar{y}_i - \bar{Y}|}{\bar{Y}} * 100 \quad (19)$$

and the Monte Carlo root mean square relative error (RMSRE):

$$RMSRE = \sqrt{\frac{1}{B} \sum_{i=1}^B \left(\frac{\bar{y}_i - \bar{Y}}{\bar{Y}} \right)^2} * 100 \quad (20)$$

where B is the number of iterations, \bar{y}_i is an estimate of \bar{Y} computed for the i -th sample.

We also examine the behaviour of variance estimators. We consider the jackknife method used in [47] to account for all sources of variability. The performance of a variance estimator along with the point estimator \bar{y}_i is assessed by the length of the intervals obtained at 95% confidence level and their real coverage.

The simulation study has been carried out using the statistical software R, and for its implementation we have needed the use of specific packages of the area, such as NonProbEst [4], KWML [22], sampling [44] and caret [25].

5.3. Results

Tables 1 and 2 contain the simulation results for y_1 and y_2 respectively for the three setups considering different sample sizes. In all setups, as expected, the proposed estimator with gradient boosting and kernel weighting (\bar{y}_{CO}) provides lower values of both $|RB|$ and RMSRE. The second best estimator is \bar{y}_{CPSA} , which obtains results similar to the first and with the rest of the estimators we obtain higher values of the $|RB|$ and RMSRE. It is also observed that the behavior pattern in terms of reduction $|RB|$ and RMSRE is similar in the three sample designs considered for the probabilistic sample.

Tables 3 and 4 show the real coverages and lengths of the corresponding 95% confi-

Table 1. Monte Carlo bias and root mean square relative error. Variable y_1

	$n_r = 250, n_v = 500$		$n_r = 250, n_v = 1000$		$n_r = 250, n_v = 2000$	
	RB	RMSRE	RB	RMSRE	RB	RMSRE
Simple random sampling without replacement						
\bar{y}_{REF}	4,7718	4,8468	4,7317	4,7947	4,8013	4,8718
\bar{y}_{RDR1}	3,0811	3,2305	2,7362	2,8952	2,7696	2,9516
\bar{y}_{RDR2}	3,2464	3,3888	2,8882	3,0445	2,9073	3,0881
\bar{y}_{CPSA}	1,2507	1,5540	1,1968	1,5116	1,3407	1,6633
\bar{y}_{CO}	1,1766	1,4611	1,1348	1,4401	1,2750	1,5757
Stratified sampling						
\bar{y}_{REF}	4,8002	4,8801	4,8636	4,9416	4,7296	4,7880
\bar{y}_{RDR1}	2,9984	3,1900	2,9125	3,1139	2,7517	2,9095
\bar{y}_{RDR2}	3,6513	3,7879	3,6009	3,7463	3,4347	3,5462
\bar{y}_{CPSA}	1,4478	1,7984	1,5946	2,0034	1,3261	1,6646
\bar{y}_{CO}	1,2242	1,5209	1,3219	1,6714	1,1624	1,4314
Midzuno sampling						
\bar{y}_{REF}	4,7708	4,8447	4,7662	4,8273	4,7351	4,7924
\bar{y}_{RDR1}	3,1003	3,2574	2,8011	2,9466	2,7663	2,9116
\bar{y}_{RDR2}	3,3814	3,5202	3,1216	3,2496	3,0686	3,1975
\bar{y}_{CPSA}	1,2194	1,5259	1,2608	1,5543	1,2386	1,5733
\bar{y}_{CO}	1,0998	1,3930	1,1406	1,4115	1,1243	1,4246

dence intervals. The coverage of intervals based on estimators \bar{y}_{REF} , \bar{y}_{RDR1} and \bar{y}_{RDR2} are very low, as expected, due to the bias in the estimates. On the contrary, the proposed estimator \bar{y}_{CO} and \bar{y}_{CPSA} have good performance, having the intervals a real coverage close to the nominal coverage. With respect to the length of the intervals, as we expected, the \bar{y}_{CO} estimator is the one with the shortest length for all types of sampling considered, sample sizes and type of variable. The KW is intended to reduce variance and indeed it succeeds for these scenarios and variables.

Table 2. Monte Carlo bias and root mean square relative error. Variable y_2

	$n_r = 250, n_v = 500$		$n_r = 250, n_v = 1000$		$n_r = 250, n_v = 2000$	
	RB	RMSRE	RB	RMSRE	RB	RMSRE
Simple random sampling without replacement						
\bar{y}_{REF}	16,2289	16,6534	16,1988	16,5467	16,3959	16,7748
\bar{y}_{RDR1}	10,0245	10,7692	9,1303	9,8089	9,1881	9,9711
\bar{y}_{RDR2}	10,6498	11,3665	9,5999	10,2814	9,6045	10,3978
\bar{y}_{CPSA}	5,1880	6,4063	5,1357	6,3999	5,7588	7,1760
\bar{y}_{CO}	4,5377	5,6647	4,6808	5,9199	5,3433	6,6418
Stratified sampling						
\bar{y}_{REF}	16,3173	16,7383	16,7025	17,1332	16,2218	16,5372
\bar{y}_{RDR1}	11,0094	11,7336	11,0115	11,7715	10,5184	11,0677
\bar{y}_{RDR2}	12,8187	13,4127	12,8210	13,4465	12,3257	12,7850
\bar{y}_{CPSA}	5,6469	7,1102	5,8664	7,6132	5,1256	6,4080
\bar{y}_{CO}	5,1189	6,4439	5,1977	6,7040	4,6117	5,6837
Midzuno sampling						
\bar{y}_{REF}	16,4209	16,8289	16,2477	16,5812	16,6900	17,0297
\bar{y}_{RDR1}	10,7378	11,4370	9,8655	10,5116	10,1815	10,8242
\bar{y}_{RDR2}	11,6344	12,2710	10,7709	11,3824	11,0198	11,6323
\bar{y}_{CPSA}	5,2458	6,6523	5,0519	6,2055	5,6324	7,0043
\bar{y}_{CO}	4,7062	5,9025	4,4902	5,5831	4,9655	6,1634

5.4. Influence of the machine learning method

In the previous simulation we used gradient boosting machine as a machine learning method, but different methods can be used. In this case we are going to make a comparison of the most used machine learning methods to see if the results are influenced by them. Specifically, we are going to compare neural networks (NNET), K-nearest neighbors (K) and logistic regression (LR) with respect to gradient boosting machine for qualitative and quantitative variables y_1 and y_2 considering the three types of sampling and for the different sample sizes. The results obtained in the comparative study can be

Table 3. Confidence intervals' real coverage and length. Variable y_1

	$n_r = 250, n_v = 500$		$n_r = 250, n_v = 1000$		$n_r = 250, n_v = 2000$	
	Coverage	Length	Coverage	Length	Coverage	Length
Simple random sampling without replacement						
\bar{y}_{REF}	0,0000	0,4811	0,0000	0,4534	0,0000	0,4378
\bar{y}_{RDR1}	0,1080	0,5171	0,1680	0,4922	0,1640	0,4745
\bar{y}_{RDR2}	0,0880	0,5220	0,1460	0,5050	0,1460	0,4904
\bar{y}_{CPSA}	0,9560	0,8534	0,9620	0,8539	0,9180	0,8546
\bar{y}_{CO}	0,9620	0,8111	0,9600	0,8073	0,9280	0,7809
Stratified sampling						
\bar{y}_{REF}	0,0020	0,5028	0,0020	0,4766	0,0000	0,4633
\bar{y}_{RDR1}	0,1900	0,5523	0,1780	0,5249	0,1780	0,5093
\bar{y}_{RDR2}	0,0540	0,5333	0,0520	0,5068	0,0300	0,4935
\bar{y}_{CPSA}	0,9440	0,9393	0,8980	0,9476	0,9540	0,9479
\bar{y}_{CO}	0,9580	0,8440	0,9060	0,8215	0,9520	0,7877
Midzuno sampling						
\bar{y}_{REF}	0,0000	0,4883	0,0000	0,4615	0,0000	0,4454
\bar{y}_{RDR1}	0,1240	0,5288	0,1460	0,5051	0,1400	0,4869
\bar{y}_{RDR2}	0,0900	0,5279	0,0900	0,5087	0,0940	0,4925
\bar{y}_{CPSA}	0,9580	0,8855	0,9620	0,8871	0,9560	0,8861
\bar{y}_{CO}	0,9520	0,8204	0,9640	0,8082	0,9500	0,7685

seen in the tables 5, 6, 7 and 8.

When comparing the $|RB|$ and the RMSRE values for y_1 for all sample sizes (table 5), we can see that in simple random sampling and Midzuno sampling the smallest values are found for \bar{y}_{CO} , in the case of stratified sampling, the smallest values are found in \bar{Y}_{CO-K} . For y_2 (table 6) the results obtained for the gradient boosting machine and K-nearest neighbors method are similar if we compare the $|RB|$ and the RMSRE values. When looking at the tables 7 and 8 for y_1 it can be observed that the greatest coverage (0.91-0.97) obtained is given in the case of the gradient boosting machine and K-nearest

Table 4. Confidence intervals' real coverage and length. Variable y_2

	$n_r = 250, n_v = 500$		$n_r = 250, n_v = 1000$		$n_r = 250, n_v = 2000$	
	Coverage	Length	Coverage	Length	Coverage	Length
Simple random sampling without replacement						
\bar{y}_{REF}	0,0080	0,0756	0,0060	0,0695	0,0020	0,0662
\bar{y}_{RDR1}	0,2760	0,0767	0,2860	0,0696	0,2360	0,0658
\bar{y}_{RDR2}	0,2420	0,0774	0,2520	0,0712	0,2320	0,0677
\bar{y}_{CPSA}	0,9680	0,1302	0,9540	0,1304	0,9300	0,1308
\bar{y}_{CO}	0,9500	0,1182	0,9540	0,1186	0,9040	0,1157
Stratified sampling						
\bar{y}_{REF}	0,0180	0,0794	0,0020	0,0735	0,0020	0,0705
\bar{y}_{RDR1}	0,1980	0,0785	0,1740	0,0720	0,1320	0,0686
\bar{y}_{RDR2}	0,1080	0,0781	0,0840	0,0714	0,0520	0,0679
\bar{y}_{CPSA}	0,9440	0,1387	0,9240	0,1395	0,9760	0,1395
\bar{y}_{CO}	0,9320	0,1259	0,9160	0,1206	0,9440	0,1142
Midzuno sampling						
\bar{y}_{REF}	0,010	0,077	0,002	0,071	0,002	0,068
\bar{y}_{RDR1}	0,232	0,077	0,232	0,071	0,162	0,067
\bar{y}_{RDR2}	0,168	0,078	0,178	0,072	0,126	0,068
\bar{y}_{CPSA}	0,950	0,133	0,988	0,134	0,958	0,134
\bar{y}_{CO}	0,950	0,122	0,960	0,118	0,924	0,115

neighbors methods. For y_2 the K-nearest neighbors method obtains the greatest coverage (0.93-0.96). With respect to the length of the confidence interval, gradient boosting machine obtains the smallest values and logistic regression model obtains the largest. The performance of the logistic regression was to be expected since the propensities do not depend on all the covariates and there is an error in the propensity model specification.

Table 5. Monte Carlo bias and root mean square relative error of estimators changing the ML method. Variable y_1

	$n_r = 250, n_v = 500$		$n_r = 250, n_v = 1000$		$n_r = 250, n_v = 2000$	
	RB	RMSRE	RB	RMSRE	RB	RMSRE
Simple random sampling without replacement						
\bar{y}_{CO}	1,1564	1,4063	1,1615	1,4275	1,2597	1,5996
$\bar{Y}_{CO-NNET}$	1,1649	1,4220	1,2432	1,5206	1,3170	1,6755
\bar{Y}_{CO-K}	1,1654	1,4182	1,1654	1,4384	1,2703	1,6103
\bar{Y}_{CO-LR}	1,1970	1,4680	1,2788	1,5679	1,3389	1,6945
Stratified sampling						
\bar{y}_{CO}	1,2497	1,5468	1,2606	1,5948	1,2566	1,5782
$\bar{Y}_{CO-NNET}$	1,3894	1,7134	1,3793	1,7729	1,4737	1,8288
\bar{Y}_{CO-K}	1,2338	1,5270	1,2501	1,5821	1,2395	1,5500
\bar{Y}_{CO-LR}	1,4669	1,8142	1,4774	1,8911	1,5567	1,9234
Midzuno sampling						
\bar{y}_{CO}	1,2538	1,5665	1,1914	1,4779	1,3072	1,6153
$\bar{Y}_{CO-NNET}$	1,3311	1,6648	1,2774	1,6051	1,4898	1,8843
\bar{Y}_{CO-K}	1,2721	1,5918	1,2034	1,4948	1,3365	1,6580
\bar{Y}_{CO-LR}	1,3823	1,7318	1,3128	1,6493	1,5287	1,9292

5.5. Influence of the kernel function

In the previous simulations we used the triangular distribution as kernel function in the construction of KW pseudo-weights, but different distributions can be used. In this case we are going to make a comparison of the distribution implemented in the R package Boosted Kernel Weighting [22] to see if the results are influenced by them. Specifically, we are going to compare triangular, standard normal (SN) and truncated standard normal (TSN) for qualitative and quantitative variables y_1 and y_2 considering the three types of sampling and for the different sample sizes. The results obtained in the comparative study can be seen in the tables 9, 10, 11 and 12.

The values of |RB| and the RMSRE are similar for the kernel functions used, so we

Table 6. Monte Carlo bias and root mean square relative error of estimators changing the ML method. Variable y_2

	$n_r = 250, n_v = 500$		$n_r = 250, n_v = 1000$		$n_r = 250, n_v = 2000$	
	RB	RMSRE	RB	RMSRE	RB	RMSRE
Simple random sampling without replacement						
\bar{y}_{CO}	5,1436	6,2637	4,9318	6,1052	5,1150	6,3776
$\bar{Y}_{CO-NNET}$	5,5098	6,7600	5,3148	6,5543	5,6183	6,8396
\bar{Y}_{CO-K}	5,1070	6,2783	5,0230	6,2403	5,0572	6,3344
\bar{Y}_{CO-LR}	5,7386	7,0105	5,5578	6,9025	5,8422	7,1011
Stratified sampling						
\bar{y}_{CO}	5,0446	6,3335	5,1506	6,5656	5,1996	6,4553
$\bar{Y}_{CO-NNET}$	5,4493	6,8482	5,8697	7,4717	6,0577	7,4612
\bar{Y}_{CO-K}	4,9665	6,3122	5,2403	6,7163	5,5527	6,8372
\bar{Y}_{CO-LR}	5,5933	7,0279	5,9386	7,5076	6,1965	7,6151
Midzuno sampling						
\bar{y}_{CO}	4,7813	5,8679	4,9698	6,3090	5,1372	6,3269
$\bar{Y}_{CO-NNET}$	5,2897	6,4671	5,6832	6,9717	5,4322	6,6979
\bar{Y}_{CO-K}	4,9220	6,0782	5,1750	6,4543	5,0860	6,2919
\bar{Y}_{CO-LR}	5,5195	6,7173	5,8068	7,1237	5,5916	6,9159

can say that there is no influence of the kernel function in this study. Regarding coverage, we see that in all cases it is quite good, moving around 0.91-0.96, obtaining the shortest length of the interval in most cases in the \bar{y}_{CO} estimator.

5.6. Results under coverage bias

In order to check the behavior of the Hartley estimator $\bar{y}_H(opt)$, proposed in section 4.2, we have repeated the previous simulation but now we include a mechanism to reproduce coverage bias in our simulation. This context is compared with the same estimators considered in the first simulation.

The probability sample is selected by SRSWOR from the full population but the non-probability sample is now selected from a frame U_v created from the population U

Table 7. Confidence intervals' real coverage and length changing the ML method. Variable y_1

	$n_r = 250, n_v = 500$		$n_r = 250, n_v = 1000$		$n_r = 250, n_v = 2000$	
	Coverage	Length	Coverage	Length	Coverage	Length
Simple random sampling without replacement						
\bar{y}_{CO}	0,9740	0,8116	0,9480	0,8051	0,9240	0,7803
$\bar{Y}_{CO-NNET}$	0,9700	0,8304	0,9560	0,8387	0,9180	0,8579
\bar{Y}_{CO-K}	0,9700	0,8225	0,9600	0,8200	0,9360	0,8208
\bar{Y}_{CO-LR}	0,9700	0,8543	0,9560	0,8667	0,9160	0,8760
Stratified sampling						
\bar{y}_{CO}	0,9460	0,8452	0,9300	0,8201	0,9160	0,7842
$\bar{Y}_{CO-NNET}$	0,9260	0,9047	0,9320	0,9145	0,9260	0,9293
\bar{Y}_{CO-K}	0,9540	0,8540	0,9380	0,8488	0,9360	0,8490
\bar{Y}_{CO-LR}	0,9180	0,9361	0,9240	0,9513	0,9200	0,9631
Midzuno sampling						
\bar{y}_{CO}	0,9140	0,8229	0,9520	0,8044	0,9180	0,7698
$\bar{Y}_{CO-NNET}$	0,9220	0,8603	0,9400	0,8753	0,8920	0,8820
\bar{Y}_{CO-K}	0,9300	0,8352	0,9520	0,8272	0,9120	0,8292
\bar{Y}_{CO-LR}	0,9180	0,8934	0,9500	0,9009	0,9080	0,9106

containing only individuals whose variable $x_5 = 1$ (related to target variables).

In tables 13 and 14 values of $|\text{RB}|$ and the RMSRE can be seen for each of the considered estimators.

As expected, all the estimators considered now have greater bias than in the previous simulation. We observe that the estimators \bar{y}_{CPSA} and \bar{y}_{CO} continue to be better than the other PSA-based estimators in terms of $|\text{RB}|$ and RMSRE reduction. As expected, the estimator based on dual frames, $\bar{y}_H(\text{opt})$, is the one that produces estimates with less $|\text{RB}|$, and consequently is also able to reduce the RMSRE compared to its competitors.

Table 8. Confidence intervals' real coverage and length changing the ML method. Variable y_2

	$n_r = 250, n_v = 500$		$n_r = 250, n_v = 1000$		$n_r = 250, n_v = 2000$	
	Coverage	Length	Coverage	Length	Coverage	Length
Simple random sampling without replacement						
\bar{y}_{CO}	0,9400	0,1181	0,9380	0,1198	0,9260	0,1186
$\bar{Y}_{CO-NNET}$	0,9060	0,1236	0,9260	0,1253	0,9000	0,1255
\bar{Y}_{CO-K}	0,9540	0,1204	0,9440	0,1199	0,9380	0,1198
\bar{Y}_{CO-LR}	0,9200	0,1283	0,9160	0,1303	0,9000	0,1295
Stratified sampling						
\bar{y}_{CO}	0,9500	0,1265	0,9280	0,1203	0,8980	0,1155
$\bar{Y}_{CO-NNET}$	0,9440	0,1375	0,9300	0,1392	0,9420	0,1391
\bar{Y}_{CO-K}	0,9600	0,1279	0,9520	0,1273	0,9500	0,1270
\bar{Y}_{CO-LR}	0,9380	0,1398	0,9200	0,1403	0,9460	0,1425
Midzuno sampling						
\bar{y}_{CO}	0,9520	0,1204	0,9340	0,1204	0,9000	0,1141
$\bar{Y}_{CO-NNET}$	0,9640	0,1299	0,9140	0,1296	0,9580	0,1329
\bar{Y}_{CO-K}	0,9620	0,1233	0,9440	0,1213	0,9580	0,1218
\bar{Y}_{CO-LR}	0,9580	0,1337	0,9300	0,1345	0,9520	0,1367

6. Discussion

In the last decade, survey research has witnessed the surge of non-probability sampling as a feasible alternative to probability sampling. In theory, the superiority of probability sampling should be clear, as it has a theoretical basis in design-based inference allowing for unbiased estimation of population parameters along with the calculation of exact sampling error. However, they are very expensive and usually have small sizes. Non-probability samples can offer some advantages in that sense, as they can be deployed in many relatively inexpensive ways, but they lack an underlying mathematical theory given their usual lack of design. This is troublesome with respect to achieving accuracy and representativeness for estimates derived from such samples.

Table 9. Monte Carlo bias and root mean square relative error of estimators changing the kernel.
Variable y_1

	$n_r = 250, n_v = 500$		$n_r = 250, n_v = 1000$		$n_r = 250, n_v = 2000$	
	RB	RMSRE	RB	RMSRE	RB	RMSRE
Simple random sampling without replacement						
\bar{y}_{CO}	1,1599	1,4479	1,1437	1,4372	1,2609	1,5779
$\hat{\bar{Y}}_{CO-SN}$	1,1642	1,4487	1,1403	1,4347	1,2641	1,5767
$\hat{\bar{Y}}_{CO-TSN}$	1,1614	1,4508	1,1449	1,4373	1,2611	1,5772
Stratified sampling						
\bar{y}_{CO}	1,2451	1,5730	1,4139	1,7341	1,2102	1,4916
$\hat{\bar{Y}}_{CO-SN}$	1,2501	1,5792	1,3891	1,7031	1,2062	1,4921
$\hat{\bar{Y}}_{CO-TSN}$	1,2559	1,5973	1,3885	1,7188	1,1998	1,4885
Midzuno sampling						
\bar{y}_{CO}	1,2210	1,5402	1,2291	1,5130	1,3117	1,6308
$\hat{\bar{Y}}_{CO-SN}$	1,2203	1,5356	1,2316	1,5175	1,3081	1,6255
$\hat{\bar{Y}}_{CO-TSN}$	1,2304	1,5478	1,2312	1,5176	1,3203	1,6322

Given their potential, many efforts have been undertaken in recent years to combine both probability and nonprobability samples to produce a single inference which may be able to overcome the limitations of each method, resulting in a rich literature on data integration in finite populations. Most of this literature is based on considering a framework where the variables of interest have not been observed in the probability sample. In this paper, we have considered the problem of observed study variables in both the non-probability sample and the probability sample, in presence of auxiliary information.

Since both samples contain the same variables, we propose a methodology to combine two surveys based on probability and non-probability samples with the help of machine learning algorithms, in order to obtain reliable estimations with small variance. We have introduced a general class of estimators, based on the kernel weighting method, and

Table 10. Monte Carlo bias and root mean square relative error of estimators changing the kernel. Variable y_2

	$n_r = 250, n_v = 500$		$n_r = 250, n_v = 1000$		$n_r = 250, n_v = 2000$	
	RB	RMSRE	RB	RMSRE	RB	RMSRE
Simple random sampling without replacement						
\bar{y}_{CO}	4,6408	5,7620	4,8391	6,0696	5,0120	6,3777
$\hat{\bar{Y}}_{CO-SN}$	4,5672	5,6788	4,8318	6,0877	5,0024	6,3352
$\hat{\bar{Y}}_{CO-TSN}$	4,6274	5,7763	4,7826	6,0414	5,0399	6,4093
Stratified sampling						
\bar{y}_{CO}	5,2150	6,6265	4,9023	6,1754	5,0693	6,2978
$\hat{\bar{Y}}_{CO-SN}$	5,1991	6,6123	4,9914	6,2498	5,0636	6,3320
$\hat{\bar{Y}}_{CO-TSN}$	5,2711	6,6313	4,9883	6,2295	5,0988	6,3772
Midzuno sampling						
\bar{y}_{CO}	4,6565	5,8727	5,1223	6,2736	4,9655	6,2108
$\hat{\bar{Y}}_{CO-SN}$	4,7364	5,8957	5,2016	6,3113	5,0140	6,2627
$\hat{\bar{Y}}_{CO-TSN}$	4,6167	5,8703	5,2201	6,3752	4,9927	6,2561

studied theoretically their bias properties. Using simulations we have also compared the proposed estimators with other methods for integrating probability and non-probability samples developed in the literature in different simulation setups, both in terms of |RB| and RMSRE.

The simulation study indicates that |RB| and RMSRE of estimators can be reduced when combining the probability and the non-probability sample using the KW method proposed here in the case where there is a relationship between the variable of interest and the participation probability. We also observed that the choice of the ML method used for propensity predictions is very important and can influence the estimates obtained. However, the kernel function in the construction of KW pseudo-weights does not influence the estimates obtained. From our simulation study we also deduce that in case the sample of volunteers has a coverage bias, it is appropriate to use an estimator

Table 11. Confidence intervals' real coverage and length changing the kernel. Variable y_1

	$n_r = 250, n_v = 500$		$n_r = 250, n_v = 1000$		$n_r = 250, n_v = 2000$	
	Coverage	Length	Coverage	Length	Coverage	Length
Simple random sampling without replacement						
\bar{y}_{CO}	0,9460	0,8120	0,9620	0,8069	0,9180	0,7820
$\hat{\bar{Y}}_{CO-SN}$	0,9560	0,8142	0,9660	0,8110	0,9200	0,7898
$\hat{\bar{Y}}_{CO-TSN}$	0,9500	0,8124	0,9680	0,8103	0,9180	0,7886
Stratified sampling						
\bar{y}_{CO}	0,9460	0,8433	0,9120	0,8275	0,9480	0,7846
$\hat{\bar{Y}}_{CO-SN}$	0,9540	0,8512	0,9300	0,8310	0,9360	0,7906
$\hat{\bar{Y}}_{CO-TSN}$	0,9320	0,8428	0,9320	0,8308	0,9460	0,7931
Midzuno sampling						
\bar{y}_{CO}	0,9300	0,8209	0,9580	0,8074	0,9120	0,7770
$\hat{\bar{Y}}_{CO-SN}$	0,9420	0,8247	0,9600	0,8129	0,9100	0,7823
$\hat{\bar{Y}}_{CO-TSN}$	0,9320	0,8212	0,9580	0,8099	0,9140	0,7845

based on dual frames that allows this bias to be treated as well.

These methods can be implemented using freely available statistical packages such as R. The R code used for the simulation study and the computation of the results is available on request.

Based on our results, our advice to practitioners is that the use of probability samples remains essential to obtain reliable estimates based on an accepted theory such as sampling theory [1], but complementing the probability sample with a non-probability sample can serve as a means to reduce the errors in the estimates.

There is a lot of room for future research to improve estimation by mean integration: other similarity measures and other weighting adjustment methods such as weight smoothing for multipurpose surveys [17] can be considered. In this work only the estimation of means and totals has been considered, but the method can be applied, with

Table 12. Confidence intervals' real coverage and length changing the kernel. Variable y_2

	$n_r = 250, n_v = 500$		$n_r = 250, n_v = 1000$		$n_r = 250, n_v = 2000$	
	Coverage	Length	Coverage	Length	Coverage	Length
Simple random sampling without replacement						
\bar{y}_{CO}	0,9440	0,1178	0,9400	0,1183	0,9320	0,1171
$\hat{\bar{Y}}_{CO-SN}$	0,9580	0,1196	0,9340	0,1213	0,9300	0,1187
$\hat{\bar{Y}}_{CO-TSN}$	0,9560	0,1191	0,9360	0,1205	0,9320	0,1189
Stratified sampling						
\bar{y}_{CO}	0,9240	0,1249	0,9540	0,1208	0,9220	0,1137
$\hat{\bar{Y}}_{CO-SN}$	0,9260	0,1262	0,9440	0,1221	0,9260	0,1173
$\hat{\bar{Y}}_{CO-TSN}$	0,9380	0,1255	0,9440	0,1217	0,9160	0,1169
Midzuno sampling						
\bar{y}_{CO}	0,9520	0,1211	0,9420	0,1186	0,9200	0,1166
$\hat{\bar{Y}}_{CO-SN}$	0,9500	0,1234	0,9420	0,1217	0,9180	0,1171
$\hat{\bar{Y}}_{CO-TSN}$	0,9480	0,1219	0,9360	0,1201	0,9160	0,1172

certain adjustments, to the case of other non-linear parameters such as distribution functions or quantiles. These issues will be future research topics.

Acknowledgments

The research was partially supported by a grant from Ministerio de Educación y Ciencia (PID2019-106861RB-I00, Spain), from IMAG-Maria de Maeztu CEX2020-001105-M/AEI/10.13039/501100011033 and from FEDER/Junta de Andalucía-Consejería de Transformación Económica, Industria, Conocimiento y Universidades (FQM170-UGR20, A-SEJ-154-UGR20).

Conflict of interest

The authors declare no potential conflict of interests.

Table 13. Monte Carlo bias and root mean square relative error of estimators with coverage bias. Variable y_1

	$n_r = 250, n_v = 500$		$n_r = 250, n_v = 1000$		$n_r = 250, n_v = 2000$	
	RB	RMSRE	RB	RMSRE	RB	RMSRE
\bar{y}_{REF}	5,5410	5,6152	5,5810	5,6493	5,5539	5,6186
\bar{y}_{RDR1}	3,2794	3,4268	3,2953	3,4208	3,1750	3,3042
\bar{y}_{RDR2}	3,2326	3,4089	3,1978	3,3578	2,9993	3,1657
\bar{y}_{CPSA}	1,2669	1,5739	1,2128	1,5345	1,2198	1,5429
\bar{y}_{CO}	1,2582	1,5630	1,2087	1,5294	1,2041	1,5200
$\bar{y}_{H(opt)}$	1,1953	1,4856	1,1249	1,4261	1,1317	1,4463

Table 14. Monte Carlo bias and root mean square relative error of estimators with coverage bias. Variable y_2

	$n_r = 250, n_v = 500$		$n_r = 250, n_v = 1000$		$n_r = 250, n_v = 2000$	
	RB	RMSRE	RB	RMSRE	RB	RMSRE
\bar{y}_{REF}	19,6890	20,0845	19,6890	20,0845	19,6890	20,0845
\bar{y}_{RDR1}	12,2102	12,8284	12,2102	12,8284	12,2102	12,8284
\bar{y}_{RDR2}	12,1179	12,7941	12,1179	12,7941	12,1179	12,7941
\bar{y}_{CPSA}	5,3593	6,6229	5,3593	6,6229	5,3593	6,6229
\bar{y}_{CO}	5,3746	6,6482	5,3746	6,6482	5,3746	6,6482
$\bar{y}_{H(opt)}$	5,2580	6,4532	5,2580	6,4532	5,2580	6,4532

A. Appendix 1

Regularity conditions for the HT estimator

The first and second order probabilities verify:

$$1a) N^{-2} \sum_{i \neq j=1}^N (\pi_i \pi_j - \pi_{ij})^r = O(n^{-2r\delta})$$

$$2a) N^{-1} \sum_{i=1}^N (y_i / \pi_i - Y/n)^{2k} < M < \infty \text{ for } \delta > 0 \text{ and } r^{-1} + k^{-1} = 1$$

Regularity conditions for the KW estimator:

The kernel function $K(u)$, the bandwidth h and the sampling schemes verify:

$$2a) K(u), \int K(u)du = 1, \sup_u |K(u)| < \infty, \text{ y } \lim_{|u| \rightarrow \infty} |u| |K(u)| = 0$$

$$2b) h = h(n_v), h \rightarrow 0, \text{ pero } n_v h \rightarrow \infty \text{ en cuanto } n_v \rightarrow \infty$$

and the distributions of the estimated propensity scores in the probability and non-probability samples are interchangeable.

References

- [1] Beaumont, J. F. (2020). Are probability surveys bound to disappear for the production of official statistics?. *Survey Methodology, Statistics Canada*, 46, 1. <http://www.statcan.gc.ca/pub/12-001-x/2020001/article/00001-eng.htm>.
- [2] Buelens, B., J. Burger, and J.A. vanden Brakel, (2018). Comparing inference methods for non-probability samples. *International Statistical Review* 86(2), 322-343.
- [3] Burakauskaitė, I., Čiginas, A. Non-probability sample integration in the survey of lithuanian census. *Workshop on Survey Statistics, Tartu*, 2022. http://isi-iass.org/home/wp-content/uploads/I.Burakauskaite_A.Ciginas_presentation.pdf.
- [4] Castro, L., Ferri, R., Rueda, M.M. (2020). NonProbEst: Estimation in Nonprobability Sampling. <https://CRAN.R-project.org/package=NonProbEst>.
- [5] Castro L., Rueda M., Ferri-García R. (2022). Combining Statistical Matching and Propensity Score Adjustment for inference from non-probability surveys *Journal of Computational and Applied Mathematics*, 404, 113414.
- [6] Chen, T. and Guestrin, C. (2016). Xgboost: A scalable tree boosting system. *In Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge*

Discovery and Data Mining, San Francisco, CA, USA 785–794.

- [7] Chen, Y., Li, P., Wu, C. (2020). Doubly robust inference with nonprobability survey samples. *Journal of the American Statistical Association*, 115(532), 2011–2021.
- [8] Deville, J. C., and Särndal, C. E. (1992). Calibration estimators in survey sampling. *Journal of the American statistical Association* 87(418), 376–382.
- [9] Disogra, C., Cobb, C.L., Chan, E.K. and Dennis, J.M. (2011). Calibrating Non-Probability Internet Samples with Probability Samples Using Early Adopter Characteristics. In *Proceedings of the American Statistical Association, Section on Survey Research. Joint Statistical Meetings (JSM)*.
- [10] Elliot, M. R. (2009). Combining data from probability and non-probability samples using pseudo-weights. *Survey Practice* 2(6), 2982.
- [11] Elliott, M. and Haviland, A. (2007). Use of a Web-Based Convenience Sample to Supplement a Probability Sample. *Survey Methodology* 33, 211–215.
- [12] Epanechnikov, V. A. (1969). Non-parametric estimation of a multivariate probability density. *Theory of Probability & Its Applications* 14(1), 153–158.
- [13] Fahimi, M., Barlas, F. M., Thomas, R. K., Buttermore, N. (2015). Scientific surveys based on incomplete sampling frames and high rates of nonresponse. *Survey Practice* 8(5), 1–11.
- [14] Ferri-García, R. and Rueda, M. M. (2018). Efficiency of Propensity Score Adjustment and calibration on the estimation from non-probabilistic online surveys. *SORT-Stat. Oper. Res. T.*, 42(2), 159–162.
- [15] Ferri-García, R., Rueda, M. D. M. (2020). Propensity score adjustment using machine learning classification algorithms to control selection bias in online surveys. *PloS One*, 15(4), e0231500.
- [16] Ferri-García, R. and Rueda, M.d.M. (2020). Propensity score adjustment using

machine learning classification algorithms to control selection bias in online surveys. *PLoS ONE* 15 e0231500.

- [17] Ferri-García, R., Beaumont, J. F., Bosa, K., Charlebois, J., and Chu, K. (2022). Weight smoothing for nonprobability surveys. *TEST* 31(3), 619-643.
- [18] Friedman, J.H. (2001). Greedy function approximation: A gradient boosting machine. *Annals of statistics* 29, 1189–1232.
- [19] Hartley, H.O. (1962). Multiple frame surveys. *In Proceedings of the Social Statistics Section, American Statistical Association* 203–206.
- [20] Isaki, C.T. and Fuller, W.A. (1982). Survey design under the regression superpopulation model. *Journal of the American Statistical Association* 77(377), 89–96.
- [21] Kennedy, C., Hartig, H. (2019). Response rates in telephone surveys have resumed their decline. Pew Research Center. <https://www.pewresearch.org/fact-tank/2019/02/27/response-rates-in-telephone-surveys-have-resumed-their-decline/>
- [22] Kern, C. and Wang, L. (2022). Boosted Kernel Weighting - Using Statistical Learning to Improve Inference from Nonprobability Samples. <https://github.com/chkern/KWML>
- [23] Kim, J. K. and Tam, S.M. (2021). Data Integration by Combining Big Data and Survey Sample Data for Finite Population Inference. *International Statistical Review* 89, (2), 382401.
- [24] Kim, J. K. and Wang, Z. (2019). Sampling techniques for big data analysis. *International Statistical Review* 87, 177–191.
- [25] Kuhn, M., Wing, J., Weston, S., Williams, A., Keefer, C., Engelhardt, A., Cooper, T., Mayer, Z., Kenkel, B., R Core Team, Benesty, M., Lescarbeau, R., Ziem, A., Scrucca, L., Tang, Y., Candan, C., Hunt, T. (2022). caret: Classification and Regression Training. <https://CRAN.R-project.org/package=caret>

- [26] Lee, B.K., Lessler, J. and Stuart, E.A. (2010). Improving propensity score weighting using machine learning. *Stat. Med.* 29, 337–346.
- [27] Lee, B.K., Lessler, J. and Stuart, E.A. (2011). Weight trimming and propensity score weighting. *PLoS ONE* 6 e18174.
- [28] Lee, S. (2006). Propensity Score Adjustment as a Weighting Scheme for Volunteer Panel Web Surveys. *J. Off. Stat.* 22(2), 329-349
- [29] Lee, S., Valliant, R. (2009). Estimation for Volunteer Panel Web Surveys Using Propensity Score Adjustment and Calibration Adjustment. *Sociol. Method. Res.* 37(3), 319-343.
- [30] Lohr, S. (2011). Alternative survey sample designs: Sampling with multiple overlapping frames. *Survey Methodology* 37(2), 197-213.
- [31] Marken S. (2018). Still Listening: The State of Telephone Surveys. Gallup Methodology Blog. <https://news.gallup.com/opinion/methodology/225143/listening-state-telephone-surveys.aspx>.
- [32] McCaffrey, D. F., Griffin, B. A., Almirall, D., Slaughter, M. E., Ramchand, R., and Burgette, L. F. (2013). A tutorial on propensity score estimation for multiple treatments using generalized boosted models. *Statistics in medicine* 32(19), 3388-3414.
- [33] McCaffrey, D. F., Ridgeway, G., and Morral, A. R. (2004). Propensity score estimation with boosted regression for evaluating causal effects in observational studies. *Psychological methods* 9(4), 403.
- [34] Mecatti, F. (2007). A single frame multiplicity estimator for multiple frame surveys. *Survey methodology* 33(2), 151-157.
- [35] Medous, E., Goga, C., Ruiz-Gazen, A., Beaumont, J. F., Dessertaine, A., and Puech, P. (2022). QR Prediction for Statistical Data Integration. *TSE Working*

Paper 22, 1344

- [36] Nekrašaitė-Liegė, V., Čiginas, A., Krapavickaitė, D. (2022). Usage of non-probability sample and scraped data to estimate proportions. *Workshop on survey statistics 2022*. <https://vb.vgtu.lt/object/elaba:138918140/index.html>
- [37] Rao, J.N.K.(2020). On Making Valid Inferences by Integrating Data from Surveys and Other Sources. *Sankhya B* 83(1), 242-272 .
- [38] Rivers, D. (2007). Sampling for web surveys. *In Proceedings of the 2007 Joint Statistical Meetings, Salt Lake City, UT, USA*.
- [39] Robbins, M. W., Ghosh-Dastidar, B., and Ramchand, R. (2021). Blending probability and nonprobability samples with applications to a survey of military caregivers. *Journal of Survey Statistics and Methodology* 9(5), 1114-1145.
- [40] Rueda, M. D. M., Ferri-García, R., and Castro-Martín, L. (2022). Combining Big Data with probability survey data: a comparison of methodologies for estimation from non-probability surveys. *Padua Research Archive-Institutional Repository* 711.
- [41] Rueda, M. D. M., Pasadas-del-Amo, S., Cobo, B., Castro-Martín, L., and Ferri-García, R. (2022). Enhancing estimation methods for integrating probability and nonprobability survey samples with machine-learning techniques. An application to a Survey on the impact of the COVID-19 pandemic in Spain. *Biometrical Journal*. <https://doi.org/10.1002/bimj.202200035>
- [42] Särndal, C. E., and Wright, R. L. (1984). Cosmetic form of estimators in survey sampling. *Scandinavian Journal of Statistics* 146-156.
- [43] Servy, E., Cuesta, C. B., Marí, G. P. D., and Armida, M. L. (2006). Utilización del paquete “kernsmooth” de r para construir suavizados loess y bandas de variabilidad a datos de la encuesta de ocupación hotelera. <http://hdl.handle>.

[net/2133/8792](https://CRAN.R-project.org/package=sampling)

- [44] Tillé, Y., and Matei, A. (2021). *sampling: Survey Sampling*. <https://CRAN.R-project.org/package=sampling>
- [45] Valliant, R., (2020). Comparing alternatives for estimation from nonprobability samples. *Journal of Survey Statistics and Methodology* 8(2), 231–263.
- [46] Valliant, R., Dorfman, A. H., and Royall, R. M. (2000). *Finite population sampling and inference: A prediction approach*. New York: John Wiley.
- [47] Wang, L., Graubard, B. I., Katki, H. A., and Li, Y. (2020). Improving external validity of epidemiologic cohort analyses: A kernel weighting approach. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, <https://doi.org/10.1111/rssa.12564>.
- [48] Wenna Xi, Alice Hinton, Bo Lu, Karol Krotki, Brittney Keller-Hamilton, Amy Ferketich and Amang Sukasih (2022): Analysis of combined probability and non-probability samples: a simulation evaluation and application to a teen smoking behavior survey, *Communications in Statistics - Simulation and Computation* DOI: 10.1080/03610918.2022.2102181.
- [49] Wiśniowski, A., Sakshaug, J. W., Perez Ruiz, D. A. and Blom, A. G. (2020). Integrating Probability and Nonprobability Samples for Survey Inference. *Journal of Survey Statistics and Methodology* 8(1), 120-147. <https://doi.org/10.1093/jssam/smz051>.
- [50] Yang, S. and Kim, J.K. (2020). Statistical data integration in survey sampling: a review. *Jpn J Stat Data Sci* 3, 625–650. <https://doi.org/10.1007/s42081-020-00093-w>