

Matched Mass Imputation for Survey Data Integration

JEREMY FLOOD AND SAYED A. MOSTAFA*

Department of Mathematics & Statistics, North Carolina A&T State University, Greensboro, NC, USA

Abstract

Analysis of nonprobability survey samples has gained much attention in recent years due to their wide availability and the declining response rates within their costly probabilistic counterparts. Still, valid population inference cannot be deduced from nonprobability samples without additional information, which typically takes the form of a smaller survey sample with a shared set of covariates. In this paper, we propose the matched mass imputation (MMI) approach as a means for integrating data from probability and nonprobability samples when common covariates are present in both samples but the variable of interest is available only in the nonprobability sample. The proposed approach borrows strength from the ideas of statistical matching and mass imputation to provide robustness against potential nonignorable bias in the nonprobability sample. Specifically, MMI is a two-step approach: first, a novel application of statistical matching identifies a subset of the nonprobability sample that closely resembles the probability sample; second, mass imputation is performed using these matched units. Our empirical results, from simulations and a real data application, demonstrate the effectiveness of the MMI estimator under nearest-neighbor matching, which almost always outperformed other imputation estimators in the presence of nonignorable bias. We also explore the effectiveness of a bootstrap variance estimation procedure for the proposed MMI estimator.

Keywords *data integration; mass imputation; nonignorable missingness; nonprobability samples; statistical matching*

1 Introduction

For the past century, probability samples have been considered the gold standard by statisticians as they facilitate statistical inference about the target population. These samples are drawn using a probability sampling design, $p(\cdot)$, which ensures that every unit in the population has a positive selection probability (Lohr, 2021; Särndal et al., 2003). They are not without challenges, however; recent literature explicitly highlights the impracticality of probability samples due to factors such as design costs, non-response rates, and other limitations (Yang and Kim, 2020; Wiśniowski et al., 2020; Wang et al., 2020; Kern et al., 2021; Li et al., 2023). The recent surge in high-dimensional convenience datasets has popularized the use of nonprobability samples to maximize available statistical information; nevertheless, without an explicit sampling design, results derived from these samples are likely to suffer from estimation bias (Yang et al., 2021; National Academies of Sciences, Engineering, and Medicine, 2018). An intuitive compromise involves *data integration*, an umbrella term for techniques that combine non-probability and probability data to leverage the strengths of both. Such techniques usually require a distinct

*Corresponding author. Email: sabdelmegeed@ncat.edu.

pattern of *monotone missingness*, where a set of shared covariates, \mathbf{X} , exist in both the probability and nonprobability sample, but the target variable, Y , is available only in the latter. Specifically, with \mathcal{U} denoting the finite population of interest, let $\mathbf{A} \subseteq \mathcal{U}$ denote a probability sample of size n_A that contains sample inclusion probabilities, π_i , and covariates, \mathbf{X}_i , for all elements $i \in \mathbf{A}$. Furthermore, let $\mathbf{B} \subseteq \mathcal{U}$ denote a nonprobability sample of size n_B that is assumed to contain information on the response variable, Y_j , and covariates, \mathbf{X}_j , for all $j \in \mathbf{B}$. In these contexts, blending sample \mathbf{A} with sample \mathbf{B} is of expressed interest, but remains challenging due to varying information in each sample.

This data integration problem has received considerable attention recently, with most research focusing on estimating the finite population mean, $\mu_N = \frac{1}{N} \sum_{i \in \mathcal{U}} Y_i$, using \mathbf{A} and \mathbf{B} together. Note that, if Y_i was measured for all elements in \mathbf{A} , a design-unbiased estimate of μ could be obtained by use of Horvitz and Thompson (1952)’s (HT) mean estimator, defined as $\hat{\mu}_\pi = \frac{1}{N} \sum_{i \in \mathbf{A}} \pi_i^{-1} Y_i$. Of course, $\hat{\mu}_\pi$ cannot be calculated since the values of Y are not observed in sample \mathbf{A} , but it may be possible to fill in these missing Y values using predictions from a regression model trained on sample \mathbf{B} . This approach is suitably named *mass imputation* since it ‘imputes’ mass values of Y in sample \mathbf{A} (Kim et al., 2021; Chen et al., 2022). In a similar vein, one may consider building a model on $\mathbf{C} = \mathbf{A} \cup \mathbf{B}$ to obtain predicted values of π_j for all $j \in \mathbf{B}$, and using these values to calculate a pseudo-HT mean estimator; these procedures fall under *propensity score estimation*, which estimates an individual’s ‘propensity’ to belong in \mathbf{B} (Beaumont and Rao, 2021; Chen et al., 2020). Furthermore, *doubly-robust* mean estimators combine estimated propensities and outcome predictions to ensure consistency, provided that either the propensity model or the outcome model is correctly specified (Yang et al., 2020).

A common assumption in the above work is the *ignorability* condition, which assumes that the response variable Y_j is independent of the sample membership indicators δ_j^B given \mathbf{X}_j for all $j \in \mathbf{B}$. Similar to the *missing at random* (MAR) assumption in missing data problems, ignorability is untestable since verifying it requires access to unavailable data (Little, 1988). To address potential bias from nonignorable sampling in \mathbf{B} , *statistical matching* pairs nonprobability sample observations with their probabilistic counterparts based on a chosen distance measure (Dever, 2018; Kalay, 2021). While Rivers (2007) and Bethlehem (2016) showed that matching can reduce estimation bias, the effectiveness of their approach is limited to MAR scenarios. Moreover, Bethlehem (2016) focused only on categorical \mathbf{X} , restricting its broader applicability.

To address this gap, we propose the *matched mass imputation* (MMI) procedure, which combines statistical matching with mass imputation. The approach consists of two steps: (1) identifying the subset of $j \in \mathbf{B}$ most similar to elements in \mathbf{A} and (2) using an outcome model trained on these matched units to impute the missing Y values in \mathbf{A} . By building the outcome model only on the matched subset, MMI reduces dependence on the full nonprobability sample and offers greater robustness to nonignorable selection mechanisms.

The remainder of the paper is organized as follows. Section 2 introduces the notation, data integration setup, and key definitions. Section 3 provides an overview of existing data integration procedures, highlighting their strengths and weaknesses. In Section 4, we introduce the proposed MMI procedure. Section 5 presents results from Monte Carlo simulations, while Section 6 summarizes findings from a real data application—both comparing the performance of our proposed estimators against various competitors. We conclude with a discussion of results in Section 7.

Table 1: Data structure for the probability (**A**) and nonprobability (**B**) samples.

Sample	d	X_1	X_2	\dots	X_p	Y
A	✓	✓	✓	✓	✓	×
B	×	✓	✓	✓	✓	✓

2 Notation and Preliminaries

Let $\mathcal{U} = \{1, 2, \dots, N\}$ denote an index set for the units in a finite population of size N . Also, let **A** denote a $n_A \times (p + 1)$ probability sample drawn from \mathcal{U} with measured variables $S_A = \{d, X_1, X_2, \dots, X_p\}$, where d denotes the *design weights*, $d_i = 1/\pi_i$ for $i \in \{1, \dots, n_A\}$, and X_1, \dots, X_p the auxiliary variables (or, *covariates*). Similarly, let **B** denote a $n_B \times (p + 1)$ nonprobability sample from \mathcal{U} with measured variables $S_B = \{Y, X_1, X_2, \dots, X_p\}$, where Y denotes a response variable of interest. This data structure is illustrated in Table 1.

Given that Y and d are missing in **A** and **B**, respectively, we may estimate μ by using a survey-weighted mean of imputed \hat{Y} (for *mass imputation*), estimated $\hat{\pi}$ (for *propensity score estimation*), or some amalgamation of the two (for *doubly-robust estimation* and *statistical matching*). These methods require *positivity* and *transportability* assumptions, stated in Definitions 1 and 2, respectively.

Definition 1 (Positivity Assumption). Let δ_j^B denote the sample **B** membership indicator which takes the value 1 if $j \in \mathbf{B}$. The assumption of positivity is satisfied if $\Pr(\delta_j^B = 1 | \mathbf{X} = \mathbf{x}) > 0$ for all $j \in \mathbf{B}$ and \mathbf{x} in the support of \mathbf{X} . That is, for every value \mathbf{x} , there is a positive probability for the corresponding units to be selected in sample **B**.

Definition 2 (Transportability Assumption). Let $f(Y|\mathbf{X})$ denote the distribution of Y conditional on \mathbf{X} . The assumption of transportability is satisfied if $f(Y|\mathbf{X}, \delta_j^B = 1) = f(Y|\mathbf{X})$.

Transportability is a crucial assumption as it allows us to ‘transport’ a prediction model built on sample **B** to sample **A**. As described by Kim et al. (2021), a sufficient condition for transportability is the *ignorability condition* in Definition 3, which is essentially Rubin (1976)’s *missing at random* (MAR) assumption.

Definition 3 (Ignorability Condition). Let $\Pr(\delta_j^B = 1 | \mathbf{X})$ denote the inclusion probability of unit $j \in \mathbf{B}$ conditional on \mathbf{X} . The assumption of ignorability is satisfied if $\Pr(\delta_j^B = 1 | \mathbf{X}, Y) = \Pr(\delta_j^B = 1 | \mathbf{X})$ for all $j \in \mathbf{B}$. Otherwise, the sample is said to possess a nonignorable, or *informative*, sampling mechanism.

Under the above setup, the following section provides an overview of some existing data integration methods, namely, mass imputation, propensity score estimation, and statistical matching.

3 Overview of Existing Methods

3.1 Mass Imputation

The mass imputation literature generally falls into two categories based on the assumed prediction model: (1) parametric models, which assume a known conditional mean function, $\mathbb{E}(Y|\mathbf{X}) =$

$m(\mathbf{X}; \boldsymbol{\beta})$, parameterized by an unknown coefficient vector $\boldsymbol{\beta}$, and (2) nonparametric models, which directly estimate a nonparametric conditional mean function, $m(\mathbf{X})$. These two approaches are described below.

3.1.1 Parametric Mass Imputation

In the parametric context, the finite population is considered as a realization from a superpopulation model, ξ_P , given by

$$Y = m(\mathbf{X}; \boldsymbol{\beta}) + v(\mathbf{X})\epsilon, \quad (1)$$

where $m(\mathbf{X}; \boldsymbol{\beta}) = \mathbb{E}_{\xi_P}(Y|\mathbf{X})$ is a *known* function of \mathbf{X} parameterized by an unknown parameter vector $\boldsymbol{\beta}$, $v(\cdot)$ is a known, strictly positive variance function, and ϵ is an error term satisfying $\mathbb{E}_{\xi_P}(\epsilon|\mathbf{X}) = 0$ and $\mathbb{E}_{\xi_P}(\epsilon^2|\mathbf{X}) = \sigma_\epsilon^2$. If \mathcal{U} was fully observed, a natural estimator of $\boldsymbol{\beta}$ could be chosen to solve the finite population score function

$$U(\boldsymbol{\beta}) = \frac{1}{N} \sum_{u \in \mathcal{U}} \left(Y_u - m(\mathbf{X}_u; \boldsymbol{\beta}) \right) \mathbf{W}(\mathbf{X}_u; \boldsymbol{\beta}) = 0$$

for some p -dimensional function \mathbf{W} ; this estimator is henceforth denoted as $\boldsymbol{\beta}_N$ since it requires information for all N population units. Since Y is measured in sample \mathbf{B} only, $\boldsymbol{\beta}$ can be estimated by solving a sample-based score function,

$$\widehat{U}(\boldsymbol{\beta}) = \frac{1}{n_B} \sum_{j \in \mathbf{B}} \left(Y_j - m(\mathbf{X}_j; \boldsymbol{\beta}) \right) \mathbf{W}(\mathbf{X}_j; \boldsymbol{\beta}) = 0,$$

whose solution is henceforth denoted as $\widehat{\boldsymbol{\beta}}$.

Under this framework, Kim et al. (2021) proposed a *parametric mass imputation estimator*, henceforth abbreviated as PMIE, that uses predictions from a semiparametric regression model to estimate the finite population mean μ_N . Specifically, let $\hat{Y}_{i,\text{par}} = m(\mathbf{X}_i; \widehat{\boldsymbol{\beta}})$ denote the predicted value of Y_i , $i \in \mathbf{A}$. Then, the PMIE is defined as

$$\hat{\mu}_{\text{PMIE}} = \frac{1}{N} \sum_{i \in \mathbf{A}} \pi_i^{-1} \hat{Y}_{i,\text{par}},$$

where π_i^{-1} denotes the first-order sampling weight associated with unit i . Under some regulatory conditions, Kim et al. (2021) showed that $\hat{\mu}_{\text{PMIE}} = \tilde{\mu}_{\text{PMIE}} + o_p(n_B^{-1/2})$, where

$$\tilde{\mu}_{\text{PMIE}} = \frac{1}{N} \sum_{i \in \mathbf{A}} \pi_i^{-1} Y_i^* + \frac{1}{n_B} \sum_{j \in \mathbf{B}} (Y_j - Y_j^*) \mathbf{W}(\mathbf{X}_j; \boldsymbol{\beta}^*)^\top \mathbf{c},$$

with $Y_i^* = m(\mathbf{X}_i; \boldsymbol{\beta}^*)$, $\boldsymbol{\beta}^* = \text{plim } \widehat{\boldsymbol{\beta}}$, and

$$\mathbf{c} = \left[n_B^{-1} \sum_{j \in \mathbf{B}} \frac{\partial m(\mathbf{X}_j; \boldsymbol{\beta}^*)}{\partial \boldsymbol{\beta}^*} \mathbf{W}(\mathbf{X}_j; \boldsymbol{\beta}^*)^\top \right]^{-1} \frac{1}{N} \sum_{u \in \mathcal{U}} \frac{\partial m(\mathbf{X}_u; \boldsymbol{\beta}^*)}{\partial \boldsymbol{\beta}^*}.$$

From here, they show that

$$\mathbb{E}(\tilde{\mu}_{\text{PMIE}} - \mu_N) = -\mathbb{E} \left(N^{-1} \sum_{u \in \mathcal{U}} (Y_u - m(\mathbf{X}_u; \boldsymbol{\beta}^*)) \right),$$

which implies that $\hat{\mu}_{\text{PMIE}}$ is asymptotically unbiased if the model is correct and ignorability holds (i.e., $\boldsymbol{\beta}^* = \boldsymbol{\beta}$). If ignorability is not satisfied, the asymptotic bias of the PMIE will still be smaller than that of the naive estimator (mean of sample \mathbf{B}) if $\text{Var}\left(Y_u - m(\mathbf{X}_u; \hat{\boldsymbol{\beta}})\right) < \text{Var}(Y_u)$ for $u \in \mathcal{U}$, which is likely true for reasonably specified models.

3.1.2 Nonparametric Mass Imputation

Nonparametric regression techniques can be used to provide robustness against misspecification of the functional form of the conditional mean function, $m(\mathbf{X})$. Assume that \mathcal{U} is a realization from the following nonparametric superpopulation model, ξ_{NP} :

$$Y = m(\mathbf{X}) + \nu(\mathbf{X})\epsilon. \quad (2)$$

To obtain mass imputation predictions for the Y values in sample \mathbf{A} , the conditional mean function $m(\mathbf{X})$ can be estimated by training a model on sample \mathbf{B} using nonparametric regression techniques such as k-nearest neighbors (Rivers, 2007), kernel regression (Chen et al., 2022), or generalized additive models (GAMs, Chen et al., 2022). However, when the covariates \mathbf{X} are high-dimensional, k-nearest neighbors and kernel regression may suffer from the curse of dimensionality. Therefore, we focus on mass imputation with GAM as described by Chen et al. (2022).

Assume that, conditional on \mathbf{X} , Y follows an exponential family distribution with

$$g^{-1}\left(\mathbb{E}_{\xi_{\text{NP}}}(Y|\mathbf{X})\right) := g^{-1}\left(m(\mathbf{X})\right) = \sum_{k=1}^p \rho_k(X_k), \quad (3)$$

where $g^{-1}(\cdot)$ denotes an inverse link function and $\{\rho_k\}_{k=1}^p$ a sequence of unknown smooth functions. Under this framework, Chen et al. (2022) considered using regression splines to estimate ρ_k , which approximate $g^{-1}\left(\mathbb{E}_{\xi_{\text{NP}}}(Y|\mathbf{X})\right)$ with linear combinations of $(m+1)$ -order basis functions, $B_m(\cdot)$ for $m = \{1, 2, \dots, M\}$, and spline coefficients, $\boldsymbol{\gamma}_m^k$, such that

$$g^{-1}\left(\mathbb{E}_{\xi_{\text{NP}}}(Y|\mathbf{X}_i)\right) \approx \sum_{k=1}^p \sum_{m=1}^M \boldsymbol{\gamma}_m^k B_m(X_{k,i})$$

(Wood, 2017; James et al., 2013). To offset the risk of overfitting, the coefficient vector $\boldsymbol{\gamma}$ was chosen to minimize the penalized likelihood $-2 \ln L(\boldsymbol{\gamma}) + \sum_{k=1}^p \lambda_k \boldsymbol{\gamma}_k^T \mathbf{S}_k \boldsymbol{\gamma}_k$, where λ_k denotes the smoothing hyperparameter of covariate X_k , $L(\boldsymbol{\gamma}) = \prod_{j \in \mathbf{B}} f(Y_j | \mathbf{X}_j; \boldsymbol{\gamma})$ is the likelihood function, $\boldsymbol{\gamma}_k^T = [\boldsymbol{\gamma}_1^k \ \boldsymbol{\gamma}_2^k \ \dots \ \boldsymbol{\gamma}_M^k]$, and

$$\mathbf{S}_k = \begin{bmatrix} \int B_1''(X_k) B_1'(X_k) dX_k & \int B_1''(X_k) B_2'(X_k) dX_k & \dots & \int B_1''(X_k) B_M'(X_k) dX_k \\ \vdots & \vdots & \vdots & \vdots \\ \int B_M''(X_k) B_1'(X_k) dX_k & \dots & \dots & \int B_M''(X_k) B_M'(X_k) dX_k \end{bmatrix}.$$

Then, Chen et al. (2022)'s nonparametric mass imputation estimator is defined as follows:

$$\begin{aligned} \hat{\mu}_{\text{MIgam}} &:= \frac{1}{N} \sum_{i \in \mathbf{A}} \pi_i^{-1} g \left(\sum_{k=1}^p \sum_{m=1}^M \hat{\boldsymbol{\gamma}}_m^k B_m(X_{k,i}) \right) \\ &= \frac{1}{N} \sum_{i \in \mathbf{A}} \pi_i^{-1} \hat{Y}_{i,\text{gam}}. \end{aligned} \quad (4)$$

3.2 Propensity Score Estimation

An alternative to mass imputation is the propensity score approach, which estimates the missing sample inclusion probabilities π_j^B for units in sample \mathbf{B} using a propensity score model built from combined data from \mathbf{A} and \mathbf{B} . The estimated π_j^B values are then used to apply the Horvitz-Thompson estimator to the Y values in sample \mathbf{B} . Similar to Chen et al. (2020), suppose that π_j^B can be modeled as $\pi_j^B = \Pr(j \in \mathbf{B} | \mathbf{X}_j) = \pi(\mathbf{X}_j; \boldsymbol{\zeta})$, where $\boldsymbol{\zeta}$ denotes a set of unknown nuisance parameters. We may choose $\hat{\boldsymbol{\zeta}}$ to maximize $\ell^*(\boldsymbol{\zeta})$, a pseudo log-likelihood function given by

$$\ell^*(\boldsymbol{\zeta}) = \sum_{j \in \mathbf{B}} \log \left\{ \frac{\pi(\mathbf{X}_j; \boldsymbol{\zeta})}{1 - \pi(\mathbf{X}_j; \boldsymbol{\zeta})} \right\} + \sum_{i \in \mathbf{A}} \frac{1}{\pi_i^A} \log \{1 - \pi(\mathbf{X}_i; \boldsymbol{\zeta})\}. \quad (5)$$

The second term on the right-hand side of (5) incorporates information from sample \mathbf{A} to estimate the unknown population sum appearing in the corresponding term of the following log-likelihood function for sample \mathbf{B} :

$$\begin{aligned} \ell(\boldsymbol{\zeta}) &= \sum_{u \in U} \{ \delta_u^B \log \pi_u^B + (1 - \delta_u^B) \log(1 - \pi_u^B) \} \\ &= \sum_{j \in \mathbf{B}} \log \left\{ \frac{\pi(\mathbf{X}_j; \boldsymbol{\zeta})}{1 - \pi(\mathbf{X}_j; \boldsymbol{\zeta})} \right\} + \sum_{u \in U} \log \{1 - \pi(\mathbf{X}_u; \boldsymbol{\zeta})\}. \end{aligned}$$

Then, the maximum likelihood estimator (MLE) of π_j is defined as $\hat{\pi}_j = \pi(\mathbf{X}_j; \hat{\boldsymbol{\zeta}})$. The propensity score model, at its best, produces $\hat{\pi}_j \approx \pi_j$ for all $j \in \mathbf{B}$; however, if the model is misspecified, $\hat{\pi}_j$ will be substantially biased, which makes this approach unfavorable in practice (Beaumont and Rao, 2021; Wang et al., 2020; Yang et al., 2020; Lee et al., 2011). As a compromise, Chen et al. (2020) proposed a doubly-robust approach that integrates propensity scores, $\hat{\pi}$, and mass imputations, \hat{Y} , in one estimator:

$$\hat{\mu}_{\text{DR}} = \frac{1}{N} \left(\sum_{j \in \mathbf{B}} \frac{Y_j - \hat{Y}_j}{\hat{\pi}_j} + \sum_{i \in \mathbf{A}} \frac{\hat{Y}_i}{\pi_i} \right). \quad (6)$$

The attractiveness of the DR estimator stems from its flexibility, in that it is consistent for μ_N if either the propensity score model or the outcome model is correctly specified. We refer the reader to Chen et al. (2020) for the details regarding the asymptotic properties of $\hat{\mu}_{\text{DR}}$.

3.3 Statistical Matching

We conclude this section with a discussion on statistical matching, a procedure that pairs elements in sample \mathbf{A} with closely related counterparts in sample \mathbf{B} . The goal is to create a subset of \mathbf{B} , denoted as \mathbf{M} , that is less biased than \mathbf{B} itself.

For continuous \mathbf{X} , Rivers (2007) facilitated cold-deck matching by pairing each unit $i \in \mathbf{A}$ with a donor $m \in \mathbf{B}$ that satisfies $|\mathbf{X}_m - \mathbf{X}_i| \leq |\mathbf{X}_j - \mathbf{X}_i|$ across all $j \in \mathbf{B}$. Letting $\tilde{Y}_{i,R} := Y_m$ denote unit m 's corresponding Y , Rivers (2007) proposed the following mean estimator:

$$\hat{\mu}_{\text{SM:R}} = \frac{1}{N} \sum_{i \in \mathbf{A}} \pi_i^{-1} \tilde{Y}_{i,R}. \quad (7)$$

Rivers (2007) reported that, under some regularity conditions (including ignorability), $\hat{\mu}_{\text{SM:R}}$ is asymptotically unbiased for μ_N .

Bethlehem (2016) extended this work to categorical \mathbf{X} by considering (random) cold-deck matching, which facilitates imputation by randomly selecting units in \mathbf{B} with exact profiles in \mathbf{A} . To describe their procedure, suppose that there exists a set of categorical covariates \mathbf{X} shared between \mathbf{A} and \mathbf{B} that, when crossed, produce H disjoint groups (or, *strata*). Let n_h and m_h denote the number of observations in \mathbf{A} and \mathbf{B} , respectively, that lie in stratum h . Assuming that \mathbf{B} is sufficiently large (i.e., $n_h \leq m_h$ for all $h \in H$) and the design of \mathbf{A} , $p(\mathbf{A})$, is a simple random sample without replacement (SRS), Bethlehem (2016) draws n_h individuals from m_h as a SRS, and uses their naive mean to estimate μ . This yields the following cold-deck estimator:

$$\hat{\mu}_{\text{SM:B}} := \frac{1}{n_A} \sum_{i \in \mathbf{A}} \tilde{Y}_{i,\mathbf{B}} = \frac{1}{n_A} \sum_{h=1}^H \sum_{i=1}^{n_h} \tilde{Y}_{i_h}. \quad (8)$$

Similar to Rivers (2007)'s $\hat{\mu}_{\text{SM:R}}$, Bethlehem (2016)'s $\hat{\mu}_{\text{SM:B}}$ is also asymptotically unbiased for μ_N if $p(\mathbf{B})$ is ignorable. Their simulation results supported this claim, as the bias of $\hat{\mu}_{\text{SM:B}}$ was indeed similar to $\hat{\mu}_\pi$ when the missingness in sample \mathbf{B} was missing at random; for nonignorable missingness, though, the bias of $\hat{\mu}_{\text{SM:B}}$ was quite noticeable and tended to match that of the naive mean of \mathbf{B} .

Taken together, the works of Rivers (2007) and Bethlehem (2016) demonstrate the impact of cold-deck matching on reducing estimation bias. However, neither approach proves effective when $p(\mathbf{B})$ is nonignorable. In the next section, we propose the *matched mass imputation* procedure as a compromise between cold-deck matching and mass imputation, offering robustness to nonignorable selection in the design of the nonprobability sample.

4 Matched Mass Imputation (MMI)

The aforementioned literature assumes *transportability* (Definition 2) or *ignorability* (Definition 3) which, given the lack of design information, are both incompatible with convenience samples. These assumptions are nonetheless crucial components of any imputation engine, and cannot be omitted; thus, it is better to replace some haphazard \mathbf{B} with $\mathbf{M} \subset \mathbf{B}$, a set of statistical matches designed to be as close to \mathbf{A} as possible. This makes transportability (or, sufficiently, ignorability) more plausible, as it only needs to hold within a less-biased \mathbf{M} . For this reason, we consider replacing \mathbf{B} in the mass imputation framework with \mathbf{M} , a subset of \mathbf{B} that possesses maximal similarity to \mathbf{A} . Specifically, our approach uses statistical matching (detailed below) to identify the subset \mathbf{M} . Then, it fits a regression model on \mathbf{M} to obtain $\hat{Y}_{A,i}^M$ for all $i \in \mathbf{A}$ and uses the survey-weighted mean of the resulting predictions as an estimator for μ_N . The resulting matched mass imputation (MMI) estimator, $\hat{\mu}_{\text{MMI}}$, is defined as

$$\begin{aligned} \hat{\mu}_{\text{MMI}} &= \left(\sum_{i \in \mathbf{A}} \frac{1}{\pi_i} \right)^{-1} \sum_{i \in \mathbf{A}} \frac{1}{\pi_i} \hat{Y}_{A,i}^M \\ &= \hat{N}^{-1} \Pi^\top \hat{Y}_A^M, \end{aligned} \quad (9)$$

where $\Pi^\top := \left[\frac{1}{\pi_1} \quad \frac{1}{\pi_2} \quad \cdots \quad \frac{1}{\pi_1} \right]$ is a $(1 \times n_A)$ row vector of design weights and \hat{Y}_A^M is a $(n_A \times 1)$ vector of predictions. The use of $\hat{N} := \sum_{i \in \mathbf{A}} \pi_i^{-1}$ to estimate N makes the mean estimator a

Algorithm 1 Matched mass imputation.

1. Use some statistical matching model to identify and isolate observations in \mathbf{B} with the highest similarity to those in \mathbf{A} . These units now belong to \mathbf{M} , a set of statistical matches.
2. Build a regression model on \mathbf{M} , say $\hat{m}_{\mathbf{M}}(\mathbf{X})$, and use it with the covariates $\mathbf{X}_{\mathbf{A}}$ to obtain predictions for the missing Y values in sample \mathbf{A} , $\hat{Y}_{\mathbf{A}}^{\mathbf{M}}$.
3. Calculate $\hat{\mu}_{\text{MMI}} = \hat{N}^{-1} \Pi^{\top} \hat{Y}_{\mathbf{A}}^{\mathbf{M}}$.

Hájek-type estimator, which is known to be particularly efficient for unequal probability sampling designs (Hájek, 1964). The general MMI procedure is summarized in Algorithm 1, with the remainder of this section detailing the identification of the set of matches \mathbf{M} .

4.1 Statistical Matching Models

In this subsection, we introduce three popular statistical matching models that can be used for the proposed MMI estimator: exact matching, coarsened exact matching, and nearest-neighbor (NN) matching.

4.1.1 Exact Matching

The exact matching model is straightforward, as it simply requires finding the intersection of samples \mathbf{A} and \mathbf{B} . The simplicity of this approach allows us to establish several useful remarks.

Let $\mathbf{M} := \mathbf{A} \cap \mathbf{B}$, and let $\Pr(j \in \mathbf{A} \cap \mathbf{B}) = \Pr(j \in \mathbf{A}) \Pr(j \in \mathbf{B})$ (that is, $p(\mathbf{A}) \perp p(\mathbf{B})$). Under the assumption that $\mathbf{B} \subseteq \mathcal{U}$ and that the selection mechanism of \mathbf{B} involves independent draws, we establish the following set of remarks.

Remark 1. $\Pr(j \in \mathbf{A} \cap \mathbf{B}) > 0$ for all $j \in \mathbf{B}$.

Proof. By definition of a probability sample, the inclusion probabilities for each observation in the population must be known and strictly positive (Särndal et al., 2003). Since $\mathbf{B} \subseteq \mathcal{U}$, every unit in \mathbf{B} has a strictly positive chance of being included in \mathbf{A} , even if $p(\mathbf{B})$ is non-ignorable. \square

Remark 2. Let $n_{\mathbf{M}}$ denote the number of units in \mathbf{M} and $n_{B_{\text{Uq}}}$ the number of unique units in \mathbf{B} . Then, $n_{B_{\text{Uq}}} \rightarrow N \implies n_{\mathbf{M}} \rightarrow n_{\mathbf{A}}$.

Proof. Note that $\max(n_{B_{\text{Uq}}}) = N$, as no sample from \mathcal{U} can contain more information than \mathcal{U} itself; therefore, $n_{B_{\text{Uq}}} = N$ implies that $\mathbf{B} = \mathcal{U}$, which further implies that $\mathbf{M} := \mathbf{A} \cap \mathbf{B} = \mathbf{A}$. Therefore, as the number of unique elements in \mathbf{B} increases to N , the number of observations in \mathbf{M} also increases and converges to $n_{\mathbf{A}}$. \square

Remark 3. Assuming that a selection mechanism of \mathbf{B} exists and is independent for all j ,

$$n_{\mathbf{M}} \sim \text{PoissonBinomial} \left(\lambda = \sum_{u \in \mathcal{U}} \Pr(u \in \mathbf{A}) \times \Pr(u \in \mathbf{B}) \right)$$

and

$$\mathbb{E}_{\mathcal{D}} \left(\sum_{j \in \mathbf{B}} \Pr(j \in \mathbf{A}) \Pr(j \in \mathbf{B}) - n_{\mathbf{M}} \right) = 0,$$

where $\mathbb{E}_{\mathcal{D}}(\cdot)$ denotes the design-based expectation.

Proof. Let δ_j^A and δ_j^B denote the sample membership indicators for samples \mathbf{A} and \mathbf{B} , respectively. Assuming that the sampling design of \mathbf{B} involves independent draws and that $p(\mathbf{A}) \perp p(\mathbf{B})$, it follows that

$$\delta_j^A \delta_j^B \sim \text{Bernouli}\left(\Pr(j \in \mathbf{A}) \times \Pr(j \in \mathbf{B})\right)$$

and

$$n_M = \sum_{j \in \mathbf{B}} \delta_j^A \delta_j^B.$$

The sum of independent Bernoulli trials with varying probabilities of success follows a Poisson Binomial distribution, which implies that \mathbf{M} is a Poisson random sample. That is, for $k = \{1, 2, \dots, N\}$,

$$p(\mathbf{M}) = \prod_{k \in \mathbf{M}} \{\Pr(k \in \mathbf{A}) \Pr(k \in \mathbf{B})\} \prod_{k \in \mathcal{U}/\mathbf{M}} \{1 - \Pr(k \in \mathbf{A}) \Pr(k \in \mathbf{B})\},$$

$$\mathbb{E}_{\mathcal{D}}(n_M) = \mathbb{E}_{\mathcal{D}}\left(\sum_{j \in \mathbf{B}} \delta_j^A \delta_j^B\right) = \sum_{u \in \mathcal{U}} \Pr(u \in \mathbf{A}) \Pr(u \in \mathbf{B}),$$

and

$$\text{Var}_{\mathcal{D}}(n_M) = \sum_{u \in \mathcal{U}} \Pr(u \in \mathbf{A}) \Pr(u \in \mathbf{B}) - \left(\Pr(u \in \mathbf{A}) \Pr(u \in \mathbf{B})\right)^2,$$

where \mathbf{M} belongs to the set of all 2^N subsets of \mathcal{U} (Wang, 1993; Särndal et al., 2003). \square

Therefore, for sufficiently large n_B , there will always exist an $\mathbf{M} \subseteq \mathbf{B}$ such that all units in \mathbf{M} are also *exactly* in \mathbf{A} ; and since $\pi_m := \Pr(m \in \mathbf{M}) = \Pr(m \in \mathbf{B}) \Pr(m \in \mathbf{A})$, we would intuitively expect significant reductions in estimation bias as $n_{B_{Uq}} \rightarrow N$ (Rivers, 2007).

4.1.2 Coarsened Exact Matching (CEM)

One limitation of exact matching is its tendency to discard too many observations, particularly when dealing with continuous covariates. This issue can be mitigated through *coarsened exact matching* (CEM), which performs exact matching after binning (or “coarsening”) the data into discrete groups (Stuart et al., 2011). The concept is akin to that of Bethlehem (2016), as it involves (1) binning each $X_p \in \mathbf{X}$ into h_p strata, (2) crossing the binned variables into H groups, and (3) identifying the (binned) intersection between \mathbf{A} and \mathbf{B} .

Although CEM tends to discard fewer data points than traditional exact matching, this depends heavily on the choice of h_p ; setting h_p too low results in spurious estimates, while setting it too high replicates the exact matching procedure. In practice, h_p is either set to an arbitrary constant, typically between 3 and 5, or determined by Sturges (1926)’s rule, which for a dataset of size n is defined as $(1 + \log_2(n))$ (Scott, 2009).

4.1.3 Nearest-Neighbor (NN) Matching

Despite the overall simplicity of the (coarsened) exact matching procedure, several limitations remain worth discussing. For one, if $n_B \rightarrow N$, then $\hat{\mu}_B = n_B^{-1} \sum_{j \in \mathbf{B}} Y_j \rightarrow \mu_N$, which challenges the necessity of matching in the first place. We also question the construct validity of exact matching on covariates alone. A preferable alternative is *nearest-neighbor* (NN) matching, which

uses a distance measure, say $d(\cdot)$, to pair each observation in \mathbf{A} with its closest neighbor in \mathbf{B} . The procedure is straightforward: for each $i \in \mathbf{A}$, find the $j \in \mathbf{B}$ whose distance d_{ij} is the smallest in the sample. These matching methods are considered “greedy” because they do not optimize an overall objective; nevertheless, they are widely used in the matching literature due to their simplicity and tendency to create well-matched groups (Stuart et al., 2011; Stuart, 2010).

Recall the NN matching algorithm proposed by Rivers (2007), which pairs each $i \in \mathbf{A}$ with an $m \in \mathbf{B}$ such that $|\mathbf{X}_m - \mathbf{X}_i| \leq |\mathbf{X}_j - \mathbf{X}_i|$ for all $j \in \mathbf{B}$. This procedure is theorized to bring the joint distribution of \mathbf{X} in \mathbf{M} closer to that in \mathbf{A} , assuming that closely matched covariates yield, on average, similar Y values, i.e., $\mathbb{E}_\xi(Y|\mathbf{X})$ is almost surely Lipschitz continuous on the support of \mathbf{X} . However, since observations are matched based only on \mathbf{X} , any information linking \mathbf{X} to Y is effectively discarded.

To efficiently utilize all available information in \mathbf{A} and \mathbf{B} , we propose two NN matching algorithms: one for continuous (or, at least, unique) Y_B values and another that accommodates general cases including noncontinuous Y_B with ties. These algorithms are presented below.

4.1.4 NN Matching: Unique Y_B

If $\mathbb{E}(Y|\mathbf{X})$ were known a priori, one could simply pair each unit $i \in \mathbf{A}$ with a donor $m \in \mathbf{B}$ that satisfies $|\mathbb{E}(Y_i|\mathbf{X}_i) - Y_m| \leq |\mathbb{E}(Y_i|\mathbf{X}_i) - Y_j|$ for all $j \in \mathbf{B}$, assuming that there exists no $j, k \in \mathbf{B}$ such that $Y_j = Y_k$. Since $\mathbb{E}(Y|\mathbf{X})$ is nonetheless unknown, we consider replacing it with $\widehat{\mathbb{E}}_B(Y_i|\mathbf{X}_i) := \hat{Y}_{A,i}^B$, the estimated value of $Y_{A,i}$ derived from a model built on \mathbf{B} and applied to the covariates in \mathbf{A} . The corresponding MMI estimator is presented in Algorithm 2.

Algorithm 2 Matched mass imputation under NN matching (unique Y_B).

1. Obtain \hat{Y}_A^B using some outcome model trained on \mathbf{B} .
2. Match, with replacement, each unit $i \in \mathbf{A}$ with a unit $m \in \mathbf{B}$ such that

$$\left| \hat{Y}_{A,i}^B - Y_m \right| \leq \left| \hat{Y}_{A,i}^B - Y_j \right| \quad \forall j \in \mathbf{B}.$$

These units now belong to \mathbf{M}_U , a set of nearest-neighbor statistical matches from \mathbf{B} with unique (highlighted by U) Y values.

3. Build a regression model on \mathbf{M}_U to obtain $\hat{Y}_{A:U}^M$.
 4. Calculate $\hat{\mu}_{\text{MMI}_{\text{unn}}} = \hat{N}^{-1} \Pi^\top \hat{Y}_{A:U}^M = \hat{N}^{-1} \sum_{i \in \mathbf{A}} \pi_i^{-1} \hat{Y}_{A,i:U}^M$.
-

Using $\left| \hat{Y}_{A,i}^B - Y_j \right|$ in place of $|\mathbf{X}_i - \mathbf{X}_j|$ allows each covariate to be weighted by its estimated relationship with Y , effectively accounting for covariates that may be discrete, nonlinear, or of varying importance.

4.1.5 NN Matching: Tied Y_B

Note that, if the variance of $\hat{Y}_{A:U}^M$ is too small, or if there exists repeated Y in \mathbf{B} , the risk of some $j \in \mathbf{B}$ being sampled too often under Algorithm 2 becomes substantial. While the former issue cannot be addressed in practice, the latter can be avoided by assuming that Y in \mathbf{B} is unique. This assumption is rather strong and may not hold in practice (see Section 6 for an example).

To address this issue, in Algorithm 3, we consider matching unit $i \in \mathbf{A}$ to satisfy $\|\mathbf{a}_i\| - \|\mathbf{b}_m\| \leq \|\mathbf{a}_i\| - \|\mathbf{b}_j\|$ for all $j \in \mathbf{B}$, where $\|\cdot\|$ denotes the L_2 norm, $\mathbf{a}_i := \begin{bmatrix} \hat{Y}_{A,i}^B & \mathbf{X}_i \end{bmatrix}$, and

Algorithm 3 Matched mass imputation under NN matching (with ties).

1. Obtain \hat{Y}_A^B using some outcome model trained on \mathbf{B} .
2. Calculate

$$\mathcal{N}_A = \begin{bmatrix} \|\mathbf{a}_1\| \\ \|\mathbf{a}_2\| \\ \dots \\ \|\mathbf{a}_{n_A}\| \end{bmatrix} \quad \text{and} \quad \mathcal{N}_B = \begin{bmatrix} \|\mathbf{b}_1\| \\ \|\mathbf{b}_2\| \\ \dots \\ \|\mathbf{b}_{n_B}\| \end{bmatrix},$$

where $\mathbf{a}_i := [\hat{Y}_{A,i}^B \quad \mathbf{X}_i]$, $\mathbf{b}_j := [Y_{B,j} \quad \mathbf{X}_j]$, and $\|\cdot\|$ denotes the L_2 norm.

3. Match, with replacement, each unit $i \in \mathbf{A}$ with a $m \in \mathbf{B}$ such that

$$|\mathcal{N}_{A,i} - \mathcal{N}_{B,m}| \leq |\mathcal{N}_{A,i} - \mathcal{N}_{B,j}| \quad \forall j \in \mathbf{B}.$$

These units now belong to \mathbf{M}_R , a set of nearest-neighbor statistical matches from \mathbf{B} with possibly repeated (highlighted by R) Y .

4. Build a regression model on \mathbf{M}_R to obtain $\hat{Y}_{A:R}^M$.
 5. Calculate $\hat{\mu}_{\text{MMI}rnn} = \hat{N}^{-1} \Pi^\top \hat{Y}_{A:R}^M = \hat{N}^{-1} \sum_{i \in \mathbf{A}} \pi_i^{-1} \hat{Y}_{A,i:R}^M$.
-

$\mathbf{b}_j := [Y_{B,j} \quad \mathbf{X}_j]$. Assuming that \mathbf{X} is numeric and unique, as in Rivers (2007), the requirement for unique Y_B becomes moot. Another potential benefit of explicitly including \mathbf{X} in the matching procedure is that it may reduce reliance on $\hat{\mathbb{E}}_B(Y_i | \mathbf{X}_i)$, thus protecting against potential bias in sample \mathbf{B} . However, a key drawback is that noisy covariates in \mathbf{X} can degrade the matching quality. Additionally, the computational burden of calculating the L_2 norm increases when \mathbf{X} is high-dimensional.

Despite their differences, both Algorithms 2 and 3 seek to identify an $\mathbf{M} \subset \mathbf{B}$ with maximal similarity to \mathbf{A} , in an attempt to align the distributions of $\delta^B | Y, \mathbf{X}$ and $\delta^A | Y, \mathbf{X}$, thus providing robustness to nonignorability in the design of \mathbf{B} , and consequently, protecting against possible violations of the transportability assumption. Given this rationale, MMI estimators are expected to outperform mass imputation estimators when the assumption of ignorability is violated. However, under ignorability, mass imputation estimators may be preferable, as MMI may needlessly discard observations from sample \mathbf{B} . Since the ignorability assumption is untestable, practitioners should use their judgment about the nonprobability sample's potential deviation from ignorability when choosing between MMI estimators and standard mass imputation. The performance of the MMI estimators from Algorithms 2 and 3 is thoroughly assessed through Monte Carlo simulations in Section 5 and a real data application in Section 6.

4.2 Variance Estimation

We conclude this section by describing a bootstrap variance estimator for the sampling variance of $\hat{\mu}_{\text{MMI}}$. Following Kim et al. (2021), we first generate L sets of replicate weights for sample \mathbf{A} and draw L bootstrap samples (of size n_B) with replacement from sample \mathbf{B} . We then use each weight-resample pair to compute the MMI estimator using (coarsened) exact or NN matching (say, $\hat{\mu}_l^l$, $l = 1, 2, \dots, L$). Using both the replicate estimates ($\hat{\mu}_l^l$) and the MMI estimator from

the original \mathbf{A} and \mathbf{B} sets (say, $\hat{\mu}_I$), we calculate the following bootstrap variance estimator:

$$\widehat{\text{Var}}_{\text{Boot}}(\hat{\mu}_I) = \frac{1}{L} \sum_{l=1}^L (\hat{\mu}_I^l - \hat{\mu}_I)^2. \quad (10)$$

The estimator in (10) can be used to construct approximately $(1 - \alpha) \times 100\%$ confidence intervals for μ_N as follows:

$$\hat{\mu}_I \pm z_{\alpha/2} \times \sqrt{\widehat{\text{Var}}_{\text{Boot}}(\hat{\mu}_I)},$$

where $z_{\alpha/2}$ is the z -score corresponding to the $(1 - \alpha/2)$ percentile under the standard normal distribution. We explore the behavior of this bootstrap variance estimator and the corresponding confidence intervals in the following section.

5 Simulation Study

We conducted a Monte Carlo study to contrast the performance of our MMI estimators with various competitors. We also examined the performance of the bootstrap variance estimator described above.

5.1 Simulation Settings

To thoroughly assess the performance of the proposed MMI estimators, we generated finite populations of size $N = 100,000$ from five superpopulation models:

- Model ξ_1 :

$$Y = \sum_{a=1}^2 4X_a + \sum_{b=3}^4 2X_b + \epsilon,$$

where $X_1, X_2 \sim N(\mu = 2, \sigma = 1)$, $X_3, X_4 \sim N(\mu = 4, \sigma = 1)$, and $\epsilon \sim N(\mu = 0, \sigma = 3)$.

- Model ξ_2 :

$$Y = \sum_{a=1}^2 4X_a + \sum_{b=3}^4 2X_b + \sum_{c \in \{1,3\}} (X_c + X_{c+1})^2 + \epsilon,$$

where $X_1, X_2 \sim \text{Uniform}(0, 4)$, $X_3, X_4 \sim \text{Uniform}(4, 8)$, and $\epsilon \sim N(\mu = 0, \sigma = 17.5)$.

- Model ξ_3 (Maia et al., 2021):

$$Y = X_1 + 2.121X_2^2 + 2\mathbb{1}(|X_3| > \bar{X}_3) + 2.619 \log(|X_1|)|X_3| \\ + 2.682X_2X_4 + 6\mathbb{1}(|X_5| > \bar{X}_5) + 1.392e^{X_6} + \epsilon,$$

where $\mathbb{1}(\cdot)$ denotes the usual indicator function, $X_1, \dots, X_6 \sim \text{Uniform}(-5, 5)$, and $\epsilon \sim N(\mu = 0, \sigma = \sqrt{1.8})$.

- Model ξ_4 :

$$Y = 4X_1^2 + 2X_2^2X_3e^{-|X_4|} + 1.5X_6 \log(|X_5|) - 12X_5 \log(|X_6|) + \sin X_3^2 \log(X_3^2) \\ + \sin(X_3X_4 - X_5) - \frac{1}{e^{1/X_4}} + X_2^2X_3^2 - X_3^2X_1^2 \log(|X_2 - X_5|) + \sin(X_4X_5^2) + \epsilon,$$

where $X_1, X_3, X_5 \sim \text{Uniform}(-5, 5)$, $X_2, X_4, X_6 \sim \text{Uniform}(5, 10)$, and $\epsilon \sim N(\mu = 300, \sigma = 20\sqrt{30})$.

- Model ξ_5 : Same as model ξ_4 , except $\epsilon \sim N(\mu = 750, \sigma = 50\sqrt{30})$.

From each finite population, a simple random sample without replacement of size $n_A = 1,000$ was used to obtain sample **A**. For sample **B** with $n_B = [n_A \ 10n_A \ 20n_A]$, stratified simple random sampling was used as follows: first, all units $u \in \mathcal{U}$ with $Y_u < \text{Median}(Y)$ were placed in stratum I, and the remaining units were placed in stratum II. Then, to allow an informative $p(\mathbf{B})$, we selected $r \times n_B$ and $(1 - r) \times n_B$ observations from each respective stratum, where $r = (0.50 \ 0.30 \ 0.15 \ 0.05 \ 0.01)$. Note that when $r = 0.50$, both strata are sampled at equal rates, resulting in a noninformative design for **B**. As r deviates from 0.50, the design of **B** becomes increasingly informative, making violations of the ignorability condition more likely. Following these steps, samples **A** and **B** were drawn $n_{\text{sim}} = 1000$ times and used to calculate the following sets of estimators:

- Existing estimators
 - $\hat{\mu}_B$: The naive mean of **B**.
 - $\hat{\mu}_{\text{MIgam}}$: The MIgam with natural cubic splines (see Eq. 4).
 - $\hat{\mu}_{\text{DRgam}}$: The doubly-robust estimator using logistic propensity scores and GAMs (GAM was fitted using R's `lm()` and `bs()` base functions with 10 degrees of freedom) with natural cubic splines (see Eq. 6).
 - $\hat{\mu}_{\text{SMexact}}$, $\hat{\mu}_{\text{SMcem}}$, $\hat{\mu}_{\text{SMunn}}$, and $\hat{\mu}_{\text{SMrnn}}$: The Hájek estimator on exact, coarsened exact (using Sturges (1926)'s rule), and nearest-neighbor matches (using Algorithms 2 and 3, respectively; see Eq. 7 and 8). Exact and CEM models were built using Stuart et al. (2011)'s `MatchIt` package in R.
- Proposed estimators
 - $\hat{\mu}_{\text{MMIexact}}$, $\hat{\mu}_{\text{MMIcem}}$, $\hat{\mu}_{\text{MMIunn}}$, and $\hat{\mu}_{\text{MMIrn}}$: Matched mass imputation estimators using GAMs trained on exact, coarsened exact, and nearest-neighbor matches (see Eq. 9).

For the sake of space, the evaluation of the variance estimator and corresponding confidence intervals was restricted to model ξ_1 with $n_A = 1,000$, $n_B = 20n_A$, and $r = (.50 \ .30 \ .15)$. $L = 500$ sets of replicate weights were generated from sample **A**, and corresponding samples with replacement (of size n_B) were generated from sample **B**; these were jointly used to calculate each of the aforementioned $\hat{\mu}$, whose values were subsequently used to calculate $\widehat{\text{Var}}_{\text{Boot}}$ using Eq. (10). This process was repeated $n_{\text{sim}} = 500$ times to calculate the following performance metrics.

5.2 Performance Metrics

We evaluate the performance of the various mean estimators described above by the root mean squared error ratio (RMSER):

$$\text{RMSER}(\hat{\mu}) = \frac{\text{RMSE}(\hat{\mu})}{\text{RMSE}(\hat{\mu}_\pi)}, \quad (11)$$

and the absolute bias ratio (ABR):

$$\text{ABR}(\hat{\mu}) = \left| \frac{\text{Bias}(\hat{\mu})}{\text{Bias}(\hat{\mu}_\pi)} \right|,$$

where $\hat{\mu}_\pi$ denotes the ‘gold-standard’ Horvitz-Thompson (HT) estimator from sample **A**. The performance of the bootstrap variance estimator is assessed using the following statistics:

- CL: The percentage of instances for which μ_N was captured in the 90% confidence interval.
- LL (UL): The mean of the confidence intervals' lower (upper) limit.
- $\hat{\mu}_I$: The Monte Carlo mean of $\hat{\mu}_I$.
- \bar{d} : The average width of the confidence intervals.
- $\widehat{\text{Var}}_{\text{Boot}}(\hat{\mu}_I)$: The Monte Carlo mean of $\widehat{\text{Var}}_{\text{Boot}}(\hat{\mu}_I)$.
- RB: The percent relative bias of $\widehat{\text{Var}}_{\text{Boot}}$ relative to its Monte Carlo equivalent,

$$\widehat{\text{Var}}_{\text{MC}}(\hat{\mu}_I) = \frac{1}{n_{\text{sim}} - 1} \sum_{m=1}^{n_{\text{sim}}} \left(\hat{\mu}_{I,m} - \frac{1}{n_{\text{sim}}} \sum_{m=1}^{n_{\text{sim}}} \hat{\mu}_{I,m} \right)^2, \quad (12)$$

defined as

$$\text{RB} = \frac{\widehat{\text{Var}}_{\text{MC}} - \widehat{\text{Var}}_{\text{Boot}}}{\widehat{\text{Var}}_{\text{MC}}} \times 100.$$

5.3 Mean Estimation Results

5.3.1 RMSER

The RMSER results are presented in Figure 1 below and in Tables 1(a), 2(a), and 3(a) in the Supplementary Material. We begin by addressing the instability of $\hat{\mu}_{\text{MMIexact}}$ and $\hat{\mu}_{\text{MMIcem}}$ when $n_B = n_A$. In this setting, the expected number of exact matches shrinks to 10, which is too small to support stable GAM smoothing. While coarsened exact matches provided some robustness—particularly for models with fewer covariates (ξ_1 and ξ_2)—the only true remedy was increasing n_B . Even then, the RMSERs for these estimators were either marginally better or noticeably worse than those of $\hat{\mu}_{\text{MIgam}}$.

On the other hand, both $\hat{\mu}_{\text{MMIunn}}$ and $\hat{\mu}_{\text{MMIrmn}}$ effectively mitigated \mathbf{B} 's selection bias. Except at $r = 0.50$, i.e., when \mathbf{B} is drawn by sampling at equal rates from the two strata, the effect of nearest-neighbor matching on mass imputation is substantial, yielding consistent efficiency gains across varying degrees of bias, from slight ($r = 0.30$) to moderate ($r = 0.15$), large ($r = 0.05$), and severe ($r = 0.01$). The superiority of these MMI estimators was evident across four of the five models, with the only exception occurring under model ξ_4 when $r = 0.30$, where they exhibited slightly higher RMSER values than the standard mass imputation estimator $\hat{\mu}_{\text{MIgam}}$. These results suggest that both $\hat{\mu}_{\text{MMIunn}}$ and $\hat{\mu}_{\text{MMIrmn}}$ provide considerable robustness to nonignorable bias in the design of sample \mathbf{B} . Although the differences between the two estimators were minor, some trends emerged: $\hat{\mu}_{\text{MMIunn}}$ performed better for well-specified $\widehat{\mathbb{E}}_B(Y|\mathbf{X})$ (i.e., the linear ξ_1 and quadratic ξ_2), whereas $\hat{\mu}_{\text{MMIrmn}}$ was superior under slight to moderate model misspecification (ξ_3 and ξ_4).

The RMSERs for the cold-deck estimators were less impressive: $\hat{\mu}_{\text{SMexact}}$ and $\hat{\mu}_{\text{SMcem}}$ closely mirrored $\hat{\mu}_B$, while $\hat{\mu}_{\text{SMrmn}}$ and $\hat{\mu}_{\text{SMunn}}$ resembled $\hat{\mu}_{\text{MIgam}}$. Evaluations of alternative cold-deck estimators may yield different findings and are of interest for future research. Nevertheless, the current results overwhelmingly support matched mass imputation over its cold-deck counterparts, particularly in the context of nonignorable sample design for \mathbf{B} .

5.3.2 ABR

The ABR results are displayed in Figure 2 below and in Tables 1(b), 2(b), and 3(b) in the Supplementary Material. Performance trends closely mirrored those observed under RMSER, though

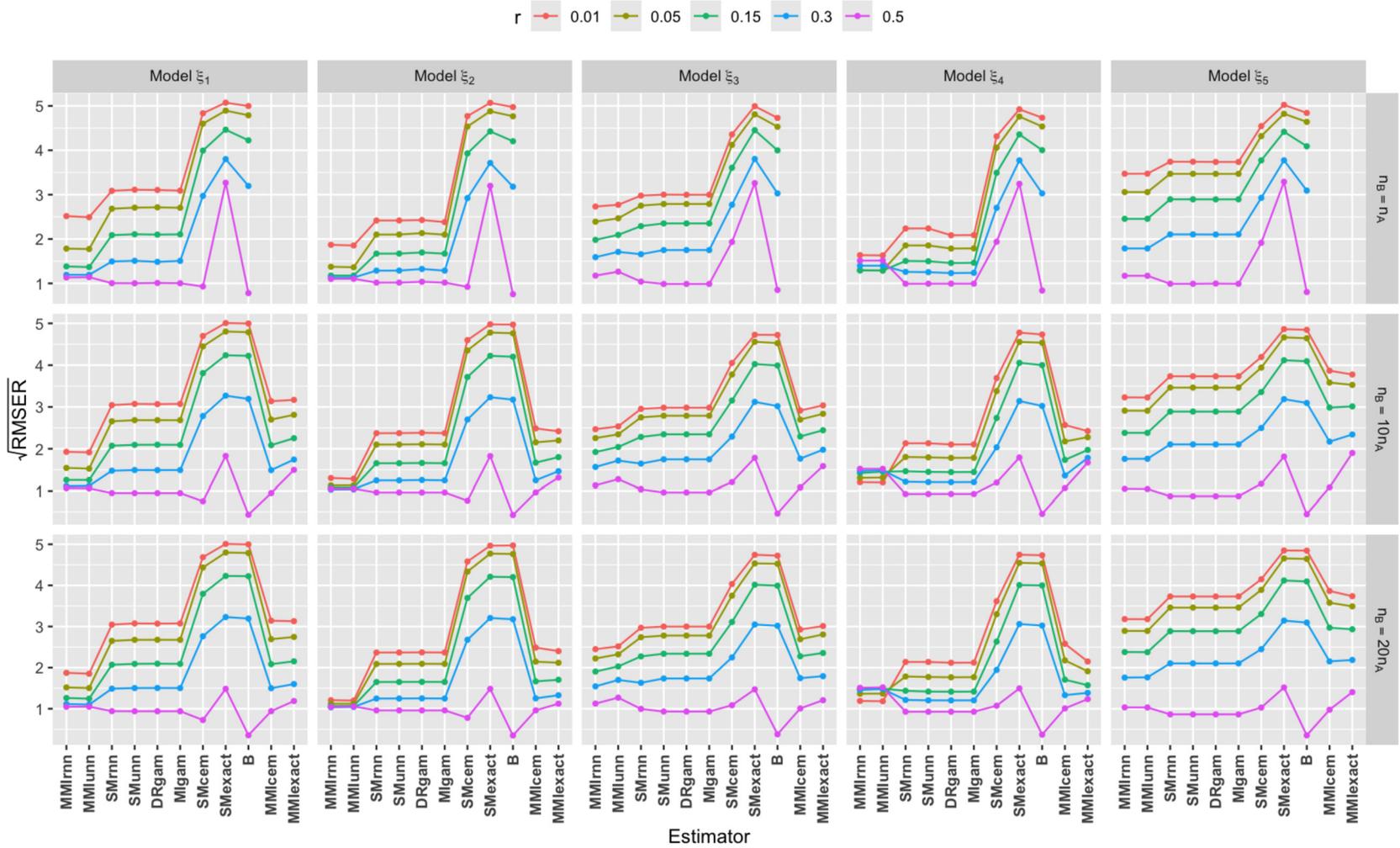


Figure 1: Square root of RMSE for eleven mean estimators. To improve readability, $\hat{\mu}_{\text{MMlcm}}$ and $\hat{\mu}_{\text{MMlexact}}$ are omitted when $n_B = n_A$ due to their extreme values.

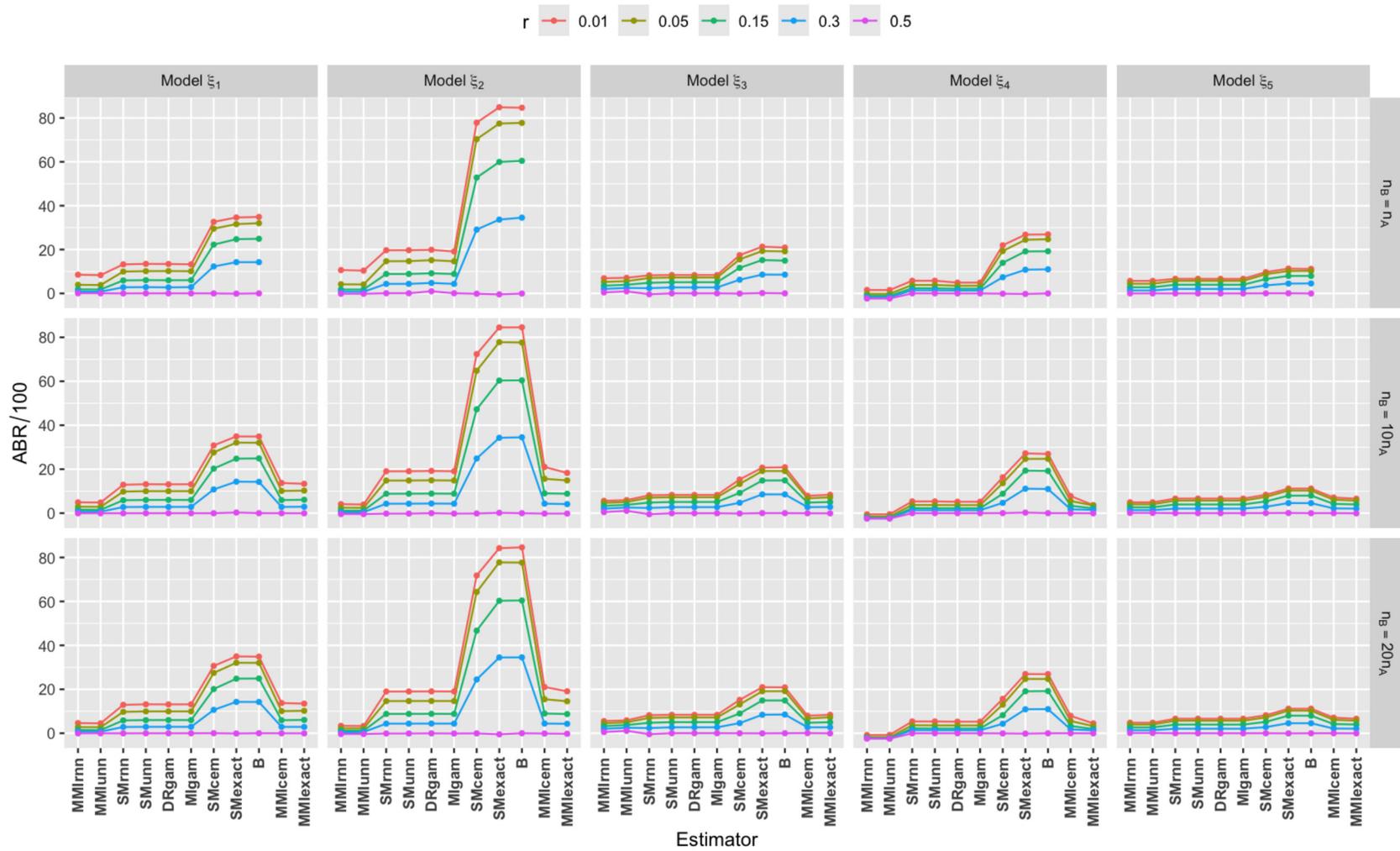


Figure 2: ABR divided by 100 for eleven mean estimators. To improve readability, $\hat{\mu}_{\text{MMlccem}}$ and $\hat{\mu}_{\text{MMllexact}}$ are omitted when $n_B = n_A$ due to their extreme values.

some additional observations warrant discussion. First, the ABRs for $\hat{\mu}_{\text{MMI}_{\text{unn}}}$ and $\hat{\mu}_{\text{MMI}_{\text{rnn}}}$ under $r \neq 0.5$ were noticeably smaller overall, suggesting that both estimators can effectively accommodate informative nonprobability samples. However, results for $r = 0.5$ highlight the consequence of needlessly discarding useful information. When \mathbf{B} is truly representative of \mathcal{U} (as is the case when $r = 0.5$), it is optimal to retain as many observations as possible. In such cases, discarding data based on Y induces missingness not at random (MNAR) in a sample that would otherwise be at least missing at random (MAR).

Additionally, we observed that both ABR and RMSEr for $\hat{\mu}_{\text{DR}_{\text{gam}}}$ were nearly identical to those of $\hat{\mu}_{\text{MI}_{\text{gam}}}$. Upon further inspection, it appears that the propensity score model yielded highly inaccurate inclusion probability estimates. The double-robust property of $\hat{\mu}_{\text{DR}_{\text{gam}}}$ mitigated this issue by effectively eliminating propensity-based influence, leaving predictions that were primarily model-driven through GAM.

5.4 Variance Estimation Results

The variance estimation results are presented in Table 2. We begin by noting a severely biased $\widehat{\text{Var}}$, which, for most mean estimators, was significantly smaller than the Monte Carlo variance $\widehat{\text{Var}}_{\text{MC}}$. These underestimated variances led to confidence intervals that were too narrow, frequently excluding μ_N at disproportionate rates. Consequently, coverage rates were significantly lower under $r = 0.30$ and $r = 0.15$, even for estimators that fared quite well in mean estimation, such as $\hat{\mu}_{\text{MMI}_{\text{unn}}}$ and $\hat{\mu}_{\text{MMI}_{\text{rnn}}}$. The only exception occurred at $r = 0.50$, where most mean estimators achieved nominal coverage or close to nominal. However, $\hat{\mu}_{\text{B}}$, $\hat{\mu}_{\text{SM}_{\text{cem}}}$, $\hat{\mu}_{\text{SM}_{\text{unn}}}$, and $\hat{\mu}_{\text{SM}_{\text{rnn}}}$ deviated from this trend: $\hat{\mu}_{\text{B}}$ produced overly conservative intervals due to large variance estimates, while the others were too liberal.

Overall, these results are unsatisfactory, as they indicate a strong reliance on MAR missingness for valid bootstrap variance estimates. If this assumption fails, the procedure proposed by Kim et al. (2021) will severely underestimate the true sampling variance. Alternative variance estimation methods, such as linearization estimators, could be explored to improve the robustness of variance estimation for MMI. We shall investigate this further in future research.

6 Real Data Application

In this section, we evaluate the performance of our proposed estimators using the National Health and Nutrition Examination Survey (NHANES) dataset published by the Centers for Disease Control and Prevention (CDC) (2015–2020). NHANES is a multistage, stratified random sample designed to be nationally representative of non-institutionalized individuals in the United States.

In this study, we estimate total cholesterol (Y , in mg/dL) in the U.S. adult population using five covariates: age (X_1), glycohemoglobin (X_2 , in %), triglycerides (X_3 , in mg/dL), direct high-density lipoprotein cholesterol (X_4 , in mg/dL), and body mass index (X_5 , in kg/m^2). The probability sample (\mathbf{A} , $n_{\text{A}} = 2,701$) consists of individuals with available data from the 2015–2016 NHANES round, while the nonprobability sample (\mathbf{B} , $n_{\text{B}} = 4,558$) comprises individuals with available data from the 2017–2020 rounds. Notably, the CDC classified observations from the 2019–2020 round as ‘nonrepresentative’ due to confidentiality concerns, requiring their merger with the 2016–2017 round. Readers interested in accessing the 2019–2020 dataset must submit a formal request through the CDC website.

Table 2: Performance statistics for the bootstrap variance estimator in Eq. (10) and corresponding 90% confidence intervals.

	Performance	$\hat{\mu}_B$	$\hat{\mu}_{MIgam}$	$\hat{\mu}_{DRgam}$	$\hat{\mu}_{SMexact}$	$\hat{\mu}_{SMcem}$	$\hat{\mu}_{SMunn}$	$\hat{\mu}_{SMrnn}$	$\hat{\mu}_{MMIexact}$	$\hat{\mu}_{MMIcem}$	$\hat{\mu}_{MMIunn}$	$\hat{\mu}_{MMIrn}$
$r = .50$	CL	99.80	90.20	90.20	89.00	61.20	14.20	14.60	92.00	90.20	86.00	86.20
	LL	31.90	31.67	31.67	31.16	31.89	31.96	31.96	31.42	31.66	31.63	31.62
	$\tilde{\mu}_I$	31.98	32.00	32.00	31.98	31.99	32.00	31.99	31.98	31.99	32.00	32.00
	UL	32.07	32.33	32.33	32.79	32.09	32.03	32.03	32.55	32.33	32.37	32.37
	\bar{d}	0.16	0.66	0.66	1.63	0.20	0.07	0.07	1.12	0.66	0.75	0.76
	\widehat{Var}_{Boot}	0.00	0.04	0.04	0.25	0.00	0.00	0.00	0.12	0.04	0.05	0.05
	RB	-270.40	-1.90	-1.80	3.10	74.00	98.80	98.70	-23.8	-1.30	17.20	18.60
$r = .30$	CL	0.00	30.60	30.60	0.00	0.00	1.80	1.40	63.40	31.40	81.60	82.60
	LL	34.14	32.10	32.10	33.44	33.56	32.39	32.37	31.85	32.10	31.70	31.71
	$\tilde{\mu}_I$	34.22	32.42	32.43	34.21	33.65	32.42	32.41	32.43	32.42	32.08	32.08
	UL	34.29	32.75	32.75	34.98	33.75	32.46	32.45	33.01	32.75	32.46	32.45
	\bar{d}	0.15	0.65	0.65	1.54	0.19	0.07	0.08	1.16	0.65	0.75	0.74
	\widehat{Var}_{Boot}	0.00	0.04	0.04	0.22	0.00	0.00	0.00	0.12	0.04	0.05	0.05
	RB	-217.10	7.70	7.70	1.60	74.00	98.80	98.70	-9.40	7.60	23.70	20.30
$r = .15$	CL	0.00	0.00	0.00	0.00	0.00	0.00	0.00	20.40	0.00	68.60	69.00
	LL	35.82	32.62	32.63	35.24	35.06	32.89	32.87	32.29	32.61	31.85	31.86
	$\tilde{\mu}_I$	35.89	32.93	32.94	35.91	35.14	32.93	32.91	32.95	32.93	32.22	32.23
	UL	35.96	33.24	33.24	36.59	35.22	32.97	32.95	33.60	33.24	32.59	32.59
	\bar{d}	0.14	0.62	0.62	1.35	0.16	0.08	0.09	1.31	0.62	0.74	0.73
	\widehat{Var}_{Boot}	0.00	0.04	0.04	0.17	0.00	0.00	0.00	0.16	0.04	0.05	0.05
	RB	-155.40	-4.30	-4.50	7.30	70.80	98.00	98.00	-16.40	-3.60	14.50	14.90

Table 3: Percent relative bias (%RB) of eleven mean estimators relative to $\hat{\mu}_\pi = 188.05$ obtained from the 2015–2016 NHANES dataset (sample **A**).

Estimator	$\hat{\mu}_{\text{MMIrn}}$	$\hat{\mu}_{\text{SMrn}}$	$\hat{\mu}_{\text{MMIcem}}$	$\hat{\mu}_{\text{MIgam}}$	$\hat{\mu}_{\text{SMcem}}$	$\hat{\mu}_{\text{B}}$	$\hat{\mu}_{\text{DRgam}}$	$\hat{\mu}_{\text{SMunn}}$	$\hat{\mu}_{\text{MMIunn}}$	$\hat{\mu}_{\text{SMexact}}$	$\hat{\mu}_{\text{MMIexact}}$
$\hat{\mu}$	186.50	186.11	185.56	185.44	184.27	179.45	167.02
%RB	0.82	1.03	1.32	1.39	2.01	4.57	11.18

Using the two samples, we computed the eleven estimators described in the previous section and their percent relative bias (%RB) defined as

$$\%RB = \frac{\hat{\mu}_\pi - \hat{\mu}}{\hat{\mu}_\pi} \times 100.$$

The results are presented in Table 3.

We begin by addressing the absence of $\hat{\mu}_{\text{MMIexact}}$, $\hat{\mu}_{\text{SMexact}}$, $\hat{\mu}_{\text{MMIunn}}$, and $\hat{\mu}_{\text{SMunn}}$. The exact match procedure failed to produce any matches, leading to the omission of its associated mean estimators from the table. Additionally, only 5.4% of Y_{B} was unique, making one-to-one matching on $\hat{Y}_{\text{A}}^{\text{B}}$ and Y_{B} infeasible. Specifically, each unit $i \in \mathbf{A}$ is matched with the first $m \in \mathbf{B}$ that satisfies the matching criteria, rendering $\hat{\mu}_{\text{MMI}}$ sensitive to the arbitrary ordering of \mathbf{B} .

Nevertheless, we observe that $\hat{\mu}_{\text{MMIrn}}$ had the lowest percent relative bias, which was substantially lower than that of $\hat{\mu}_{\text{MIgam}}$ (0.82 vs. 1.39) and even more so compared to $\hat{\mu}_{\text{B}}$ (0.82 vs. 4.57). These results, along with those presented in Section 5, further confirm that our nearest-neighbor MMI estimators provide substantial robustness to nonignorable bias in \mathbf{B} . We also note that both $\hat{\mu}_{\text{SMrn}}$ and $\hat{\mu}_{\text{MMIcem}}$ outperformed $\hat{\mu}_{\text{MIgam}}$. The bias of $\hat{\mu}_{\text{DRgam}}$ was notably higher than that of $\hat{\mu}_{\text{MIgam}}$, primarily due to severe misspecification of the propensity score model.

7 Discussion

In this paper, we addressed the problem of combining data from a probability sample (**A**) and a nonprobability sample (**B**) to estimate the finite population mean. We introduced the *matched mass imputation* (MMI) procedure, which replaces **B** in the mass imputation framework with a set of statistical matches ($\mathbf{M} \subset \mathbf{B}$). This approach shifts the burden of the ignorability assumption away from **B** and toward a less biased subset, thereby improving robustness.

The empirical simulations and real data analysis strongly support this approach, with our proposed nearest-neighbor mean estimators ($\hat{\mu}_{\text{MMIrn}}$ and $\hat{\mu}_{\text{MMIunn}}$) demonstrating the best overall performance. We also explored variance estimation using the bootstrapping procedure proposed by Kim et al. (2021). However, our results indicated a severe underestimation of sampling variance when the ignorability assumption was violated. This finding highlights the need for future research on alternative variance estimation techniques to enhance the reliability of MMI-based inference. Future studies may also explore a formal investigation of the asymptotic properties of the MMI estimators, particularly $\hat{\mu}_{\text{MMIrn}}$ and $\hat{\mu}_{\text{MMIunn}}$. This aspect was not pursued in the current work, as our primary objective was to introduce the MMI approach in its general form without restricting the discussion to a specific member of the MMI class.

Supplementary Material

The supplementary material includes the following: (1) additional simulation results showing the RMSE and ABR results in tabular format to complement the visualizations in Figures 1 and 2, (2) R code, and (3) README: a brief explanation of how to run the code.

Funding

This work of Jeremy Flood was funded by the North Carolina A&T State University Chancellor's Distinguished Fellowship, a Title III HBGI grant from the U.S. Department of Education.

References

- Beaumont JF, Rao J (2021). Pitfalls of making inferences from non-probability samples: Can data integration through probability samples provide remedies? *The Survey Statistician*, 83: 11–22.
- Bethlehem J (2016). Solving the nonresponse problem with sample matching? *Social Science Computer Review*, 34(1): 59–77. <https://doi.org/10.1177/0894439315573926>
- Centers for Disease Control and Prevention (CDC) (2015–2020). NHANES - National Health and Nutrition Examination Survey. <https://www.cdc.gov/nchs/nhanes/index.htm> (visited: 2023-10-11).
- Chen S, Yang S, Kim JK (2022). Nonparametric mass imputation for data integration. *Journal of Survey Statistics and Methodology*, 10(1): 1–24. <https://doi.org/10.1093/jssam/smaa036>
- Chen Y, Li P, Wu C (2020). Doubly robust inference with nonprobability survey samples. *Journal of the American Statistical Association*, 115(532): 2011–2021. <https://doi.org/10.1080/01621459.2019.1677241>
- Dever J (2018). Combining probability and nonprobability samples to form efficient hybrid estimates: An evaluation of the common support assumption. In: *Proceedings of the 2018 Federal Committee on Statistical Methodology (FCSM) Research Conference*, 1–15.
- Hájek J (1964). Asymptotic theory of rejective sampling with varying probabilities from a finite population. *The Annals of Mathematical Statistics*, 35(4): 1491–1523.
- Horvitz DG, Thompson DJ (1952). A generalization of sampling without replacement from a finite universe. *Journal of the American Statistical Association*, 47(260): 663–685. <https://doi.org/10.1080/01621459.1952.10483446>
- James G, Witten D, Hastie T, Tibshirani R, et al. (2013). *An Introduction to Statistical Learning*, volume 112. Springer.
- Kalay AF (2021). Double Robust Mass-Imputation with Matching Estimators. arXiv preprint: <https://arxiv.org/abs/2110.09275>.
- Kern C, Li Y, Wang L (2021). Boosted kernel weighting—using statistical learning to improve inference from nonprobability samples. *Journal of Survey Statistics and Methodology*, 9(5): 1088–1113. <https://doi.org/10.1093/jssam/smaa028>
- Kim JK, Park S, Chen Y, Wu C (2021). Combining non-probability and probability survey samples through mass imputation. *Journal of the Royal Statistical Society. Series A. Statistics in Society*, 184(3): 941–963. <https://doi.org/10.1111/rssa.12696>
- Lee BK, Lessler J, Stuart EA (2011). Weight trimming and propensity score weighting. *PLoS ONE*, 6(3): e18174. <https://doi.org/10.1371/journal.pone.0018174>

- Li Y, Fay M, Hunsberger S, Graubard BI (2023). Variable inclusion strategies for effective quota sampling and propensity modeling: An application to sars-cov-2 infection prevalence estimation. *Journal of Survey Statistics and Methodology*, 11(5): 1204–1228. <https://doi.org/10.1093/jssam/smad026>
- Little RJ (1988). A test of missing completely at random for multivariate data with missing values. *Journal of the American Statistical Association*, 83(404): 1198–1202.
- Lohr SL (2021). *Sampling: Design and Analysis*. Chapman and Hall/CRC.
- Maia M, Azevedo AR, Ara A (2021). Predictive comparison between random machines and random forests. *Journal of Data Science*, 19(4): 593–614. <https://doi.org/10.6339/21-JDS1025>
- National Academies of Sciences, Engineering, and Medicine (2018). *Federal Statistics, Multiple Data Sources, and Privacy Protection: Next Steps*. National Academies Press.
- Rivers D (2007). *Sampling for web surveys*. American Statistical Association, Alexandria, VA, 1–26.
- Rubin DB (1976). Inference and missing data. *Biometrika*, 63(3): 581–592. Publisher: Oxford University Press. <https://doi.org/10.1093/biomet/63.3.581>
- Särndal CE, Swensson B, Wretman J (2003). *Model Assisted Survey Sampling*. Springer Science & Business Media.
- Scott DW (2009). Sturges’ rule. *Wiley Interdisciplinary Reviews. Computational Statistics*, 1(3): 303–306. <https://doi.org/10.1002/wics.35>
- Stuart EA (2010). Matching methods for causal inference: A review and a look forward. *Statistical Science*, 25(1): 1. <https://doi.org/10.1214/09-STS313>
- Stuart EA, King G, Imai K, Ho D (2011). MatchIt: Nonparametric preprocessing for parametric causal inference. *Journal of Statistical Software*, 42(8): 1–28. <https://doi.org/10.18637/jss.v042.i08>
- Sturges HA (1926). The choice of a class interval. *Journal of the American Statistical Association*, 21(153): 65–66. <https://doi.org/10.1080/01621459.1926.10502161>
- Wang L, Graubard BI, Katki HA, 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*, 183(3): 1293–1311.
- Wang YH (1993). On the number of successes in independent trials. *Statistica Sinica*, 3(2): 295–312.
- Wiśniowski A, Sakshaug JW, Perez Ruiz DA, Blom AG (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>
- Wood SN (2017). *Generalized Additive Models: An Introduction with R*. CRC Press.
- Yang S, Kim JK (2020). Statistical data integration in survey sampling: A review. *Japanese Journal of Statistics and Data Science*, 3: 625–650. <https://doi.org/10.1007/s42081-020-00093-w>
- Yang S, Kim JK, Hwang Y (2021). Integration of data from probability surveys and big found data for finite population inference using mass imputation. *Survey Methodology*, 47(1): 29–58.
- Yang S, Kim JK, Song R (2020). Doubly robust inference when combining probability and non-probability samples with high dimensional data. *Journal of the Royal Statistical Society, Series B, Statistical Methodology*, 82(2): 445–465. <https://doi.org/10.1111/rssb.12354>