

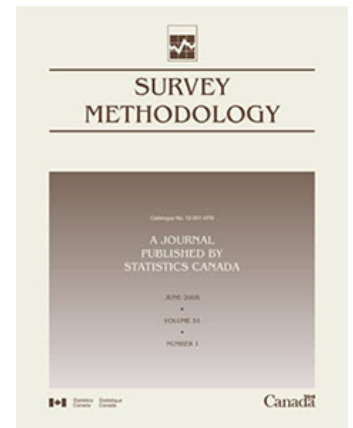
Catalogue no. 12-001-X
ISSN 1492-0921

Survey Methodology

Fully synthetic data for complex surveys

by Shirley Mathur, Yajuan Si and Jerome P. Reiter

Release date: December 20, 2024



Statistics
Canada

Statistique
Canada

Canada

How to obtain more information

For information about this product or the wide range of services and data available from Statistics Canada, visit our website, www.statcan.gc.ca.

You can also contact us by

Email at infostats@statcan.gc.ca

Telephone, from Monday to Friday, 8:30 a.m. to 4:30 p.m., at the following numbers:

- | | |
|---|----------------|
| • Statistical Information Service | 1-800-263-1136 |
| • National telecommunications device for the hearing impaired | 1-800-363-7629 |
| • Fax line | 1-514-283-9350 |

Standards of service to the public

Statistics Canada is committed to serving its clients in a prompt, reliable and courteous manner. To this end, Statistics Canada has developed standards of service that its employees observe. To obtain a copy of these service standards, please contact Statistics Canada toll-free at 1-800-263-1136. The service standards are also published on www.statcan.gc.ca under "Contact us" > "[Standards of service to the public.](#)"

Note of appreciation

Canada owes the success of its statistical system to a long-standing partnership between Statistics Canada, the citizens of Canada, its businesses, governments and other institutions. Accurate and timely statistical information could not be produced without their continued co-operation and goodwill.

Published by authority of the Minister responsible for Statistics Canada

© His Majesty the King in Right of Canada, as represented by the Minister of Industry, 2024

Use of this publication is governed by the Statistics Canada [Open Licence Agreement](#).

An [HTML version](#) is also available.

Cette publication est aussi disponible en français.

Fully synthetic data for complex surveys

Shirley Mathur, Yajuan Si and Jerome P. Reiter¹

Abstract

When seeking to release public use files for confidential data, statistical agencies can generate fully synthetic data. We propose an approach for making fully synthetic data from surveys collected with complex sampling designs. Our approach adheres to the general strategy proposed by Rubin (1993). Specifically, we generate pseudo-populations by applying the weighted finite population Bayesian bootstrap to account for survey weights, take simple random samples from those pseudo-populations, estimate synthesis models using these simple random samples, and release simulated data drawn from the models as public use files. To facilitate variance estimation, we use the framework of multiple imputation with two data generation strategies. In the first, we generate multiple data sets from each simple random sample. In the second, we generate a single synthetic data set from each simple random sample. We present multiple imputation combining rules for each setting. We illustrate the repeated sampling properties of the combining rules via simulation studies, including comparisons with synthetic data generation based on pseudo-likelihood methods. We apply the proposed methods to a subset of data from the American Community Survey.

Key Words: Bootstrap; Confidentiality; Disclosure; Privacy; Weights.

1. Introduction

Many national statistics agencies, survey organizations, and researchers – henceforth all called agencies – disseminate microdata, i.e., data on individual units, to the public. Wide dissemination of microdata greatly benefits society, enabling broad subsets of the research community to access and analyze the collected data (Reiter, 2009). Often, however, agencies cannot release microdata as collected, because doing so could reveal survey respondents’ identities or values of sensitive attributes, thereby failing to satisfy ethical or legal requirements to protect data subjects’ confidentiality (Reiter and Raghunathan, 2007).

To manage these risks, several agencies have implemented or are considering synthetic data approaches, as first proposed by Rubin (1993). In this approach, the agency (i) randomly and independently samples units from the sampling frame to comprise each synthetic data set, (ii) imputes the unknown data values for units in the synthetic samples using models fit with the original survey data, and (iii) releases multiple versions of these data sets to the public. These are called fully synthetic data sets (Drechsler, 2011; Raghunathan, 2021). Releasing fully synthetic data can preserve confidentiality, since identification of units and their sensitive data can be difficult when the released data are not actual, collected values (Reiter and Drechsler, 2010). Methods for inferences from these multiply-imputed data files have been developed for a variety of statistical inference tasks (Raghunathan, Reiter and Rubin, 2003; Reiter, 2002, 2005a,b; Drechsler and Reiter, 2010; Si and Reiter, 2011).

While prominent applications of fully synthetic data exist for censuses or administrative data (e.g., Kinney, Reiter, Reznick, Miranda, Jarmin and Abowd, 2011), many research data sets are based on surveys

1. Shirley Mathur, Department of Statistics, B-313 Padelford Hall, University of Washington, Seattle, WA 98195-4322; Yajuan Si, Survey Research Center, Institute for Social Research, University of Michigan, Rm 4014, 426 Thompson St., Ann Arbor, MI 48104. E-mail: yajuan@umich.edu; Jerome P. Reiter, Department of Statistical Science, 214a Old Chemistry Building, Duke University, Durham, NC 27708-0251.

collected with sampling designs that use unequal probabilities of selection. Previous research on multiple imputation for missing data suggests that imputation models should account for the survey design features, such as stratification, clustering, and survey weights (Reiter, Raghunathan and Kinney, 2006). Similarly, when using multiple imputation for synthetic data, the models also should account for the survey design (Mitra and Reiter, 2006; Fienberg, 2010; Kim, Drechsler and Thompson, 2021). The key challenge is properly incorporating weights in the synthesis models, which relates to the long-standing debate about the role of survey weights in model-based inferences (Pfeffermann, 1993, 2011; Little, 2004).

Researchers have proposed a variety of approaches for generating fully synthetic data in complex surveys. The suggestion in early work (Rubin, 1993; Raghunathan et al., 2003; Reiter, 2002) was to take a Bayesian finite population inference approach, in which the agency (i) builds predictive models for the survey variables conditional on design features like stratum/cluster indicators or size measures, which are assumed known by the agency for every unit in the population, (ii) imputes the missing survey variables for the nonsampled units in the population, and (iii) takes a simple random sample from the completed population to release as one synthetic data set. A related approach uses the weighted finite population Bayesian bootstrap (WFPBB) (Dong, Elliott and Raghunathan, 2014), in which the agency generates completed populations by replicating individuals from the confidential data in proportion to their survey weights and then releases the completed populations, forgoing the step of simple random sampling. More recently, it has been suggested to build synthetic data models that account for the sampling design directly, so that they estimate the joint distribution of the population data. For example, the agency can use a pseudo-likelihood approach (Pfeffermann, 1993; Savitsky and Toth, 2016), in which each individual's contribution to the likelihood function of a synthesis model is raised to a power that is a function of the survey weights (Kim et al., 2021). Departing from the proposal of Rubin (1993), a completely different approach is to create and attach new weights to synthetic data records simulated from models that are agnostic to the survey weights (United Nations Economic Commission for Europe, 2022). Here, the goal is to allow users to use weighted estimates that scale up to the finite population. The new weights can be created by treating the survey weights as a variable in the synthesis, so that the agency specifies a predictive model for the weights. The simulated weights may be adjusted by raking or calibration before inclusion in the released file.

Each of these methods has its potential drawbacks. The Bayesian finite population inference approach, while theoretically principled, requires completing full populations, which can be cumbersome, and the availability of design variables for all records in the population, which may not be the case in some surveys. The WFPBB releases (multiple copies of) individuals' genuine data records, which creates obvious disclosure risks. Pseudo-likelihood approaches may not estimate sampling variability correctly (Williams and Savitsky, 2021), and it is not clear how easily they can be implemented with machine learning synthesizers like classification and regression trees (Reiter, 2005c), which are commonly used in practical synthetic data projects (Raab, Nowok and Dibben, 2018). With synthesized weights, secondary analysts are expected to use the simulated weights to approximate design-based inference. This approximation does not have a theoretical basis; as such, it is unclear whether the synthetic weights approach facilitates accurate inferences in general.

In this article, we propose an approach to generate fully synthetic data from complex samples in the spirit of the original proposal of Rubin (1993), i.e., the agency releases simple random samples that do not require users to perform survey-weighted analyses with the synthetic data. To do so, we build on the WFPBB approach of Dong et al. (2014) by first creating pseudo-populations that account for the survey weights. We then take simple random samples (SRSs) from each pseudo-population, estimate synthesis models from each SRS, and generate draws from these models to create multiply-imputed, fully synthetic public use files. The latter step provides confidentiality protection, as the agency is not releasing genuine records. We consider two processes for the last step of generating the synthetic data. In *SynRep-R*, we generate multiple synthetic data sets from each SRS. In *SynRep-I*, we generate one synthetic data set from each SRS. *SynRep-R* releases more data sets than *SynRep-I*, which can result in reduced variances. However, the additional data sets can increase the overhead for the agency and secondary analysts, and they provide additional information for adversaries seeking disclosures. For both approaches, we derive multiple imputation combining rules that enable the estimation of variances. Using simulation studies, we illustrate the repeated sampling performances of the combining rules and compare them to fully synthetic data generated while disregarding the sampling design entirely. We also compare them against approaches that use synthesis models estimated with weighted pseudo-likelihoods (Kim et al., 2021). Finally, we illustrate the proposed methods using a subset of the American Community Survey (ACS) data. Code for the simulation studies and the ACS illustration is available at <https://github.com/yajuansi-sophie/SynRep>.

The remainder of the article is organized as follows. Section 2 describes the two synthetic data generation processes in detail and presents the new combining rules. Section 3 presents the simulation studies. Section 4 presents the illustration with the ACS data. Section 5 suggests topics for future research.

2. Proposed methods for generating fully synthetic survey data

Let \mathcal{D} be a probability sample of size n randomly drawn from a finite population comprising N units. For $i = 1, \dots, N$, let π_i be the selection probability for unit i , and let $w_i = 1/\pi_i$ be the unit's survey weight. Here, we are agnostic as to whether w_i is potentially adjusted, e.g., for normalization, calibration or nonresponse, although in our simulation studies we use pure design weights. For $i = 1, \dots, N$, let Y_i be the $p \times 1$ vector of survey variables. Hence, $\mathcal{D} = \{(w_i, Y_i) : i = 1, \dots, n\}$. For simplicity of exposition, we suppose that $p = 1$, so that Y_i is a scalar. *SynRep-R* and *SynRep-I*, and their corresponding inferential methods, can be used with multivariate survey data as well.

In Section 2.1, we describe the processes of generating synthetic data. In Section 2.2, we describe the inferential methods. As mentioned in Section 1 and following the proposal in Rubin (1993), we take as a goal allowing secondary users to analyze the released data sets as if they were simple random samples from the population.

2.1 Data generation process

Figure 2.1 and Figure 2.2 display the processes of generating synthetic data for *SynRep-R* and *SynRep-I*, respectively. We now describe these steps in detail.

Figure 2.1 Process for generating synthetic data with multiple data sets per simple random sample (SRS), which we call *SynRep-R*.

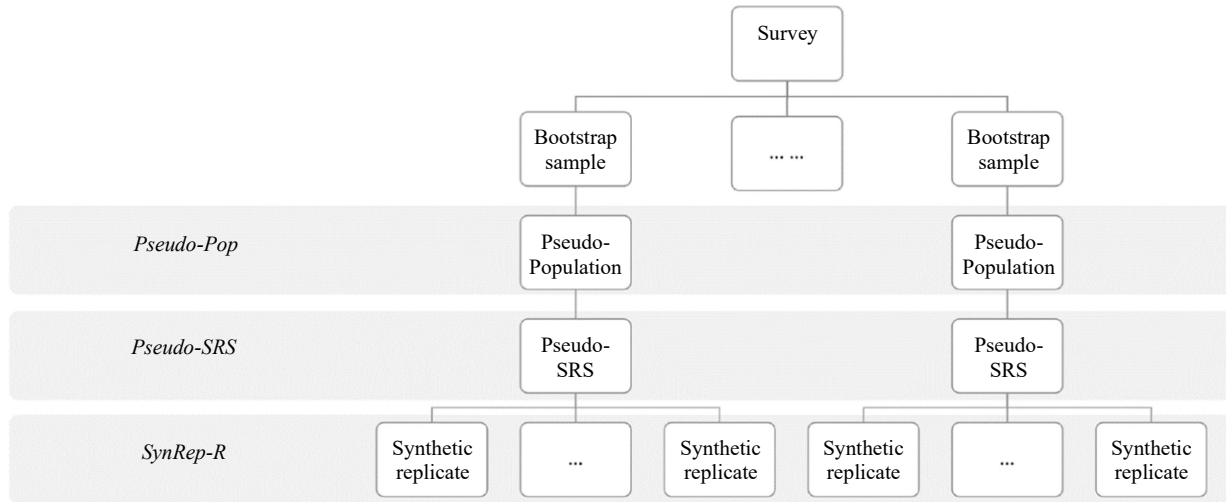
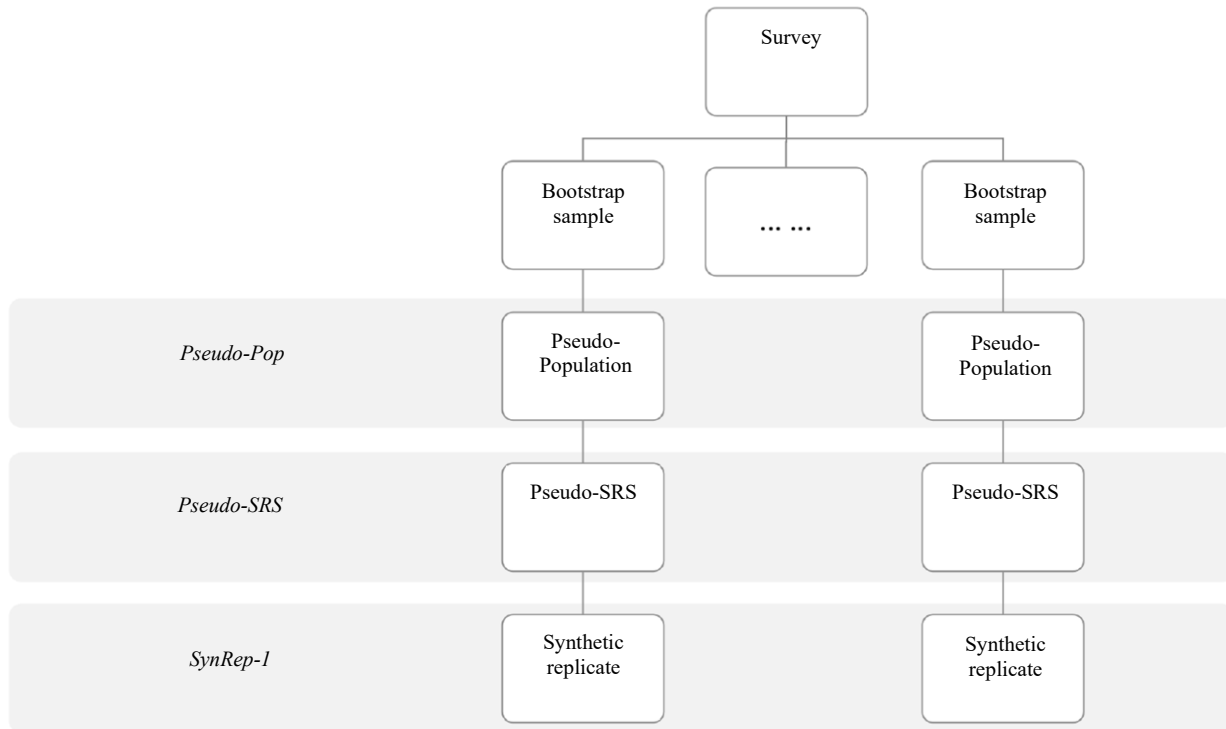


Figure 2.2 Process for generating synthetic data with one data set per simple random sample (SRS), which we call *SynRep-I*.



In either process, the first step is to generate pseudo-populations using the WFPBB (Dong et al., 2014). The WFPBB generates pseudo-populations by “undoing” the complex sampling design and accounting for the sampling weights. The idea is to draw from the posterior predictive distribution of non-observed data (Y_{nob}) given the observed data (Y_{obs}) and the survey weights, i.e., drawing from $P(Y_{\text{nob}} | Y_{\text{obs}}, w_1, \dots, w_n)$. This distribution supposes that the population is comprised of the unique values of $Y_i \in \mathcal{D}$, and that the corresponding counts for each value in the population follow a multinomial distribution. With a non-informative Dirichlet prior distribution on the multinomial probabilities, the Pólya distribution can be used to draw the predictive samples in place of the Dirichlet-multinomial distribution.

With this in mind, the process of generating the synthetic data is described below.

1. **Resample via Bayesian bootstrap:** To inject sufficient sampling variability, using the data from the “parent” sample \mathcal{D} , we generate M samples, $(\mathcal{S}^{(1)}, \dots, \mathcal{S}^{(M)})$, each of size n using independent Bayesian bootstraps (Rubin, 1981). For each $\mathcal{S}^{(m)}$ and for $i = 1, \dots, n$, let $w_i^{(m)} = cw_i r_i^{(m)}$, where $r_i^{(m)}$ is the number of times that element i from \mathcal{D} appears in $\mathcal{S}^{(m)}$. The c is a normalizing constant to ensure that the new weights sum to the population size N . Thus, in each $\mathcal{S}^{(m)}$, for $i = 1, \dots, n$, we create $w_i^{(m)} = (Nw_i r_i^{(m)}) / (\sum_k w_k r_k^{(m)})$.
2. **Use the WFPBB to make pseudo-populations:** For each $\mathcal{S}^{(m)}$, we construct an initial Pólya urn using the set of $\{Y_i, w_i^{(m)}\}$. We then draw $N - n$ units using probabilities $(p_1^{(m)}, \dots, p_n^{(m)})$ determined from

$$p_i^{(m)} = \frac{w_i^{(m)} - 1 + l_{i,k-1}^{(m)}(N - n)/n}{N - n + (k - 1)(N - n)/n}, \tag{2.1}$$

for the k th draw, $k \in \{1, \dots, N - n\}$, where $l_{i,k-1}^{(m)}$ is the number of bootstrap selections of Y_i among the elements present in the urn at the $k - 1$ draw. The $N - n$ draws combined with the data in $\mathcal{S}^{(m)}$ comprise one pseudo-population, $\mathcal{P}^{(m)}$. We repeat this for $m = 1, \dots, M$ to create $\mathcal{P}_{\text{pseudo}} = \{\mathcal{P}^{(m)} : m = 1, \dots, M\}$. When N is very large, we can save memory and computational costs by creating a pseudo-population that is large enough to be practically the same for inference as a population of size N , which we operationalize by generating $50n$ rather than $N - n$ records.

3. **Draw SRS from each pseudo-population:** For $m = 1, \dots, M$, take a simple random sample $\mathcal{D}^{(m)}$ of size n from $\mathcal{P}^{(m)}$. Let $\mathcal{D}_{\text{srs}} = \{\mathcal{D}^{(m)} : m = 1, \dots, M\}$.
4. **Generate synthetic data replicates:** For $m = 1, \dots, M$, estimate a synthesis model using $\mathcal{D}^{(m)}$, and draw from the predictive distributions to form synthetic data replicates using either Step 4a or Step 4b.
 - 4a. *SynRep-R:* For $m = 1, \dots, M$, draw $R > 1$ synthetic replicates $\mathcal{D}_{\text{syn}}^{(m,r)}$ of size n , where $r = 1, \dots, R$, using each $\mathcal{D}^{(m)}$. We release $\mathcal{D}_{\text{syn}} = \{\mathcal{D}_{\text{syn}}^{(m,r)} : m = 1, \dots, M; r = 1, \dots, R\}$ including indicators of which m each $\mathcal{D}_{\text{syn}}^{(m,r)}$ belongs to.
 - 4b. *SynRep-I:* For $m = 1, \dots, M$, draw one synthetic data sample $\mathcal{D}_{\text{syn}}^{(m)}$ of size n from each $\mathcal{D}^{(m)}$. Release $\mathcal{D}_{\text{syn}} = \{\mathcal{D}_{\text{syn}}^{(m)} : m = 1, \dots, M\}$.

The synthesis model for each $\mathcal{D}^{(m)}$ can utilize plug-in values of model parameters, e.g., their maximum likelihood estimates. It is not necessary to use posterior distributions at this stage of the process (Reiter and Kinney, 2012).

As these two processes for generating synthetic data differ from those of Raghunathan et al. (2003), as well as from other synthetic data scenarios such as those of Reiter (2003, 2004), we require new methods for inferences, to which we now turn.

2.2 Inferences for *SynRep-R* and *SynRep-1*

To derive the inferential methods, we follow the general strategy of multiple imputation (Rubin, 1987) and use a Bayesian inference approach. For any population quantity Q , such as the population mean $Q \equiv \bar{Y}$, we seek the posterior distribution $P(Q | \mathcal{D}_{\text{syn}})$. Following Raghunathan et al. (2003), we compute the following integral based upon each level of the data synthesis process from Figure 2.1 or Figure 2.2.

$$P(Q | \mathcal{D}_{\text{syn}}) = \iiint P(Q | \mathcal{D}_{\text{syn}}, \mathcal{D}_{\text{srs}}, \mathcal{P}_{\text{pseudo}}, \mathcal{D}) P(\mathcal{D} | \mathcal{D}_{\text{syn}}, \mathcal{D}_{\text{srs}}, \mathcal{P}_{\text{pseudo}}) P(\mathcal{P}_{\text{pseudo}} | \mathcal{D}_{\text{syn}}, \mathcal{D}_{\text{srs}}) P(\mathcal{D}_{\text{srs}} | \mathcal{D}_{\text{syn}}) d\mathcal{D} d\mathcal{P}_{\text{pseudo}} d\mathcal{D}_{\text{srs}}. \quad (2.2)$$

When we condition on \mathcal{D} , the values of $(\mathcal{D}_{\text{syn}}, \mathcal{D}_{\text{srs}}, \mathcal{P}_{\text{pseudo}})$ do not provide any additional information about Q . Thus, we can simplify $P(Q | \mathcal{D}_{\text{syn}}, \mathcal{D}_{\text{srs}}, \mathcal{P}_{\text{pseudo}}, \mathcal{D}) = P(Q | \mathcal{D})$. When we condition on $\mathcal{P}_{\text{pseudo}}$, the values of $(\mathcal{D}_{\text{rep}}, \mathcal{D}_{\text{syn}})$ provide no additional information about \mathcal{D} . Thus, we simplify $P(\mathcal{D} | \mathcal{D}_{\text{syn}}, \mathcal{D}_{\text{srs}}, \mathcal{P}_{\text{pseudo}}) = P(\mathcal{D} | \mathcal{P}_{\text{pseudo}})$. When we condition on \mathcal{D}_{srs} , the value of \mathcal{D}_{syn} provides no information about $\mathcal{P}_{\text{pseudo}}$. Hence, $P(\mathcal{P}_{\text{pseudo}} | \mathcal{D}_{\text{syn}}, \mathcal{D}_{\text{srs}}) = P(\mathcal{P}_{\text{pseudo}} | \mathcal{D}_{\text{srs}})$. With some re-arrangement to aid interpretation, we re-express (2.2) as

$$P(Q | \mathcal{D}_{\text{syn}}) = \int \left[\int \left[\int P(Q | \mathcal{D}) P(\mathcal{D} | \mathcal{P}_{\text{pseudo}}) d\mathcal{D} \right] P(\mathcal{P}_{\text{pseudo}} | \mathcal{D}_{\text{srs}}) d\mathcal{P}_{\text{pseudo}} \right] P(\mathcal{D}_{\text{srs}} | \mathcal{D}_{\text{syn}}) d\mathcal{D}_{\text{srs}}. \quad (2.3)$$

We begin with $P(Q | \mathcal{P}_{\text{pseudo}}) = \int P(Q | \mathcal{D}) P(\mathcal{D} | \mathcal{P}_{\text{pseudo}}) d\mathcal{D}$. We assume that, for large M , this is approximately a normal distribution. This should be reasonable in large samples, which are typical in settings where agencies want to release public use data. We only require the posterior distribution of Q to be normal, not the distribution of the survey variables themselves; indeed, the underlying data can be categorical. We note that the inferential methods are not intended for quantities like medians or other quantiles; inferential methods for such quantities is a topic for additional research.

We only require means and variances to characterize normal sampling distributions. Thus, we focus on estimating the distributions of the first two moments. For $m = 1, \dots, M$, let $Q^{(m)}$ be the computed value of Q if we had access to $\mathcal{P}^{(m)}$. Rubin (1987) shows that

$$(Q | \mathcal{P}_{\text{pseudo}}) \sim t_{M-1}(\bar{Q}, (1 + M^{-1}) B), \quad (2.4)$$

where $\bar{Q} = \sum_m Q^{(m)} / M$ and $B = \sum_m (Q^{(m)} - \bar{Q})^2 / (M - 1)$. Here $t_\nu(\mu, \sigma^2)$ denotes a t -distribution with ν degrees of freedom, location μ , and variance σ^2 . In the derivations, for convenience we approximate the t -distribution in (2.4) as a normal distribution, which should be reasonable for somewhat large M .

We next turn to $P(\mathcal{P}_{\text{pseudo}} | \mathcal{D}_{\text{srs}})$. Here, we only need $P(\bar{Q}, B | \mathcal{D}_{\text{srs}})$. For $m = 1, \dots, M$, let $q^{(m)}$ be the estimate of $Q^{(m)}$ and $v^{(m)}$ be the estimate of the sampling variance associated with $q^{(m)}$; we could compute these if we had access to $\mathcal{D}^{(m)}$. We assume that $\{q^{(m)}, v^{(m)} : m = 1, \dots, M\}$ are valid in the following sense.

- 1) For each m , $q^{(m)}$ is approximately unbiased for $Q^{(m)}$ and asymptotically normally distributed, with respect to repeated sampling from the pseudo-population $\mathcal{P}^{(m)}$ with sampling variance $V^{(m)}$. That is, we have $(q^{(m)} | \mathcal{P}^{(m)}) \sim N(Q^{(m)}, V^{(m)})$.
- 2) The sampling variance estimate $v^{(m)}$ is approximately unbiased for $V^{(m)}$, and the sampling variability in $v^{(m)}$ is negligible. That is, $(v^{(m)} | \mathcal{P}^{(m)}) \approx V^{(m)}$.
- 3) The variation in $V^{(m)}$ across the M pseudo-populations is negligible; that is, $V^{(m)} \approx V \approx \bar{v}$, where $\bar{v} = \sum_m v^{(m)} / M$.

Using standard Bayesian arguments based on these sampling distributions, it follows that

$$(Q^{(m)} | q^{(m)}, \bar{v}) \sim N(q^{(m)}, \bar{v}) \tag{2.5}$$

$$(\bar{Q} | \bar{q}, \bar{v}) \sim N(\bar{q}, \bar{v}/M), \tag{2.6}$$

where $\bar{q} = \sum_m q^{(m)} / M$.

To obtain the distribution of $(Q | \mathcal{D}_{\text{srs}})$, we integrate the distribution in (2.4), which we approximate as a normal distribution, with respect to the distributions of \bar{Q} and B . We only need the first two moments since the resulting distribution is a normal distribution. We have

$$E(Q | \mathcal{D}_{\text{srs}}) = E(E(Q | \bar{Q}) | \mathcal{D}_{\text{srs}}) = E(\bar{Q} | \mathcal{D}_{\text{srs}}) = \bar{q}. \tag{2.7}$$

We also have

$$\begin{aligned} \text{Var}(Q | \mathcal{D}_{\text{srs}}) &= E(\text{Var}(Q | \mathcal{P}_{\text{pseudo}}) | \mathcal{D}_{\text{srs}}) + \text{Var}(E(Q | \mathcal{P}_{\text{pseudo}}) | \mathcal{D}_{\text{srs}}) \\ &= (1 + M^{-1}) E(B | \mathcal{D}_{\text{srs}}) + \bar{v}/M. \end{aligned} \tag{2.8}$$

This is the variance estimator in Raghunathan et al. (2003), which analysts would use if the agency releases \mathcal{D}_{srs} as the public use files. However, since we take an additional step of replacing each $\mathcal{D}^{(m)}$ with simulated values, we need to average over the distributions of (\bar{q}, \bar{v}, B) . The result depends on whether we use *SynRep-R* or *SynRep-I*, as we now describe.

2.2.1 Derivation with *SynRep-R*

For each $\mathcal{D}_{\text{syn}}^{(m,r)}$, let $q_{\text{syn}}^{(m,r)}$ be the point estimate of Q , and let $v_{\text{syn}}^{(m,r)}$ be the estimate of the variance associated with $q_{\text{syn}}^{(m,r)}$. The analyst computes $q_{\text{syn}}^{(m,r)}$ and $v_{\text{syn}}^{(m,r)}$ acting as if $\mathcal{D}_{\text{syn}}^{(m,r)}$ is the collected data obtained

via a simple random sample of size n taken from the population. The analyst needs to compute the following quantities.

$$\bar{q}_{\text{syn}}^{(m)} = \sum_{r=1}^R q_{\text{syn}}^{(m,r)} / R \quad (2.9)$$

$$\bar{q}_{\text{syn}} = \sum_{m=1}^M \bar{q}_{\text{syn}}^{(m)} / M \quad (2.10)$$

$$b_{\text{syn}} = \sum_{m=1}^M (\bar{q}_{\text{syn}}^{(m)} - \bar{q}_{\text{syn}})^2 / (M - 1) \quad (2.11)$$

$$w_{\text{syn}}^{(m)} = \sum_{r=1}^R (q_{\text{syn}}^{(m,r)} - \bar{q}_{\text{syn}}^{(m)})^2 / (R - 1) \quad (2.12)$$

$$\bar{w}_{\text{syn}} = \sum_{m=1}^M w_{\text{syn}}^{(m)} / M \quad (2.13)$$

$$\bar{v}_{\text{syn}} = \sum_{m=1}^M \sum_{r=1}^R v_{\text{syn}}^{(m,r)} / MR. \quad (2.14)$$

We now complete the derivation of the posterior distribution for $(Q | \mathcal{D}_{\text{syn}})$ in the *SynRep-R* approach. To do so, we assume large-sample normal approximations for the sampling distributions of the point estimates. Specifically, for all (m, r) , we assume that

$$q_{\text{syn}}^{(m,r)} \sim N(q^{(m)}, W^{(m)}), \quad (2.15)$$

where $W^{(m)}$ is the sampling variance for $q_{\text{syn}}^{(m,r)}$ over draws of synthetic data from $\mathcal{D}^{(m)}$. The normality should be reasonable when n is large. Assuming diffuse prior distributions and conditioning on $W^{(m)}$, we have

$$(q^{(m)} | \mathcal{D}_{\text{syn}}^{(m,1)}, \dots, \mathcal{D}_{\text{syn}}^{(m,R)}, W^{(m)}) \sim N(\bar{q}_{\text{syn}}^{(m)}, W^{(m)} / R) \quad (2.16)$$

$$(\bar{q} | \mathcal{D}_{\text{syn}}, \bar{W}) \sim N(\bar{q}_{\text{syn}}, \bar{W} / MR), \quad (2.17)$$

where $\bar{W} = \sum_m W^{(m)} / M$.

Having now determined distributions for the point estimators, we put everything together for the posterior distribution of Q . Since all the components are normal distributions, $P(Q | \mathcal{D}_{\text{syn}})$ is a normal distribution. Thus, for the expectation, we use (2.7) and (2.17) to obtain

$$E(Q | \mathcal{D}_{\text{syn}}) = (E(Q | \mathcal{D}_{\text{srs}}) | \mathcal{D}_{\text{syn}}) = E(\bar{q} | \mathcal{D}_{\text{syn}}) = \bar{q}_{\text{syn}}. \quad (2.18)$$

For the variance, we first write the variance in terms of (B, \bar{v}, \bar{W}) and then plug in point estimates of these terms. To emphasize the use of (B, \bar{v}, \bar{W}) , we write

$$\begin{aligned}\text{Var}(Q | \mathcal{D}_{\text{syn}}, B, \bar{v}_M, \bar{W}) &= E(((1 + M^{-1})B + \bar{v}/M) | \mathcal{D}_{\text{syn}}, B, \bar{v}, \bar{W}) + \text{Var}(\bar{q} | \mathcal{D}_{\text{syn}}, B, \bar{v}, \bar{W}) \\ &= (1 + M^{-1})B + \bar{v}/M + \bar{W}/MR.\end{aligned}\quad (2.19)$$

We now define the estimates for (B, \bar{v}, \bar{W}) , which we plug into (2.19). For \bar{v} , we assume that $\bar{v}_{\text{syn}} \approx \bar{v}$. This assumption follows from the rationale in Raghunathan et al. (2003), who argue this is the case when the synthetic data are generated from the same underlying distribution as the data used to fit the models.

For \bar{W} , we note that (2.15) implies that, for $m = 1, \dots, M$,

$$\frac{(R-1)w_{\text{syn}}^{(m)}}{W^{(m)}} \sim \chi_{R-1}^2. \quad (2.20)$$

We further assume that each $W^{(m)} \approx \bar{W}$. This assumption is in line with a similar assumption provided in Reiter (2004) regarding the variability of posterior variances. Essentially, as stated in Reiter (2004), this assumption stems from the observation that variability amongst posterior variances is generally smaller in magnitude than variability in posterior expectations. With this assumption and utilizing (2.20), we have

$$\sum_{m=1}^M \frac{(R-1)w_{\text{syn}}^{(m)}}{\bar{W}} \sim \chi_{M(R-1)}^2. \quad (2.21)$$

Thus, we have

$$E\left(\sum_{m=1}^M \frac{(R-1)w_{\text{syn}}^{(m)}}{\bar{W}}\right) = M(R-1). \quad (2.22)$$

Utilizing a methods of moments approach to approximate \bar{W} , we obtain $\bar{W} \approx \bar{w}_{\text{syn}}$.

For approximating B , we note that the sampling distribution of a randomly generated $\bar{q}_{\text{syn}}^{(m)}$ over all steps in the data generation process is $N(Q, B + \bar{v} + \bar{W}/R)$. Using this fact, we have

$$\frac{\sum_{m=1}^M (\bar{q}_{\text{syn}}^{(m)} - \bar{q}_{\text{syn}})^2}{B + \bar{v} + \bar{W}/R} \sim \chi_{M-1}^2, \quad (2.23)$$

so that

$$E\left(\frac{\sum_{m=1}^M (\bar{q}_{\text{syn}}^{(m)} - \bar{q}_{\text{syn}})^2}{B + \bar{v} + \bar{W}/R}\right) = M-1. \quad (2.24)$$

Using a method of moments approach and the definition of b_{syn} in (2.11), and the plug-in estimate \bar{w}_{syn} for \bar{W} , we have $b_{\text{syn}} \approx B + \bar{v}_{\text{syn}} + \bar{w}_{\text{syn}}/R$, so that $B \approx b_{\text{syn}} - \bar{v}_{\text{syn}} - \bar{w}_{\text{syn}}/R$.

Putting all together, we can approximate $\text{Var}(Q | \mathcal{D}_{\text{syn}})$ with the estimate T_r , where

$$\begin{aligned}
 T_r &= (1 + M^{-1}) \left(b_{\text{syn}} - \bar{v}_{\text{syn}} - \bar{w}_{\text{syn}}/R \right) + \bar{v}_{\text{syn}}/M + \bar{w}_{\text{syn}}/MR \\
 &= (1 + M^{-1}) b_{\text{syn}} - \bar{v}_{\text{syn}} - \bar{w}_{\text{syn}}/R.
 \end{aligned}
 \tag{2.25}$$

We compute approximate 95% intervals for Q as $\bar{q}_{\text{syn}} \pm t_{0.975, M-1} \sqrt{T_r}$. The t -distribution is a simple approximation based on the degrees of freedom in (2.4). As with the variance estimator in Raghunathan et al. (2003), the estimate T_r can be negative, particularly for small M . As an *ad hoc* adjustment when $T_r < 0$, we recommend replacing B with \bar{v} in (2.19) and using $T_r^* = (1 + 2/M) \bar{v}_{\text{syn}} + \bar{w}_{\text{syn}}/MR$.

2.2.2 Derivation with *SynRep-1*

With large M and R , *SynRep-R* results in many synthetic data sets, which may be undesirable from the perspective of the agency and secondary data analysts. Instead, agencies may want to use *SynRep-1*. To obtain inferences for Q in this setting, we leverage the methodology of Raab et al. (2018), who observed that when the source data come from a simple random sample, as is the case for each $\mathcal{D}^{(m)}$, we can obtain valid variance estimates with single implicates with adjustments of the combining rules. We now describe this derivation.

For $m = 1, \dots, M$, let $q_{\text{syn}}^{(m)}$ be the point estimate of Q computed using $\mathcal{D}_{\text{syn}}^{(m)}$, and let $v_{\text{syn}}^{(m)}$ be the estimated variance associated with $q_{\text{syn}}^{(m)}$. The analyst computes each $(q_{\text{syn}}^{(m)}, v_{\text{syn}}^{(m)})$ by acting as if $\mathcal{D}_{\text{syn}}^{(m)}$ is a SRS of size n from the population. We require the following quantities for inferences. To economize on notation, we re-use some of the notation introduced in Section 2.2.1.

$$\bar{q}_{\text{syn}} = \sum_{m=1}^M q_{\text{syn}}^{(m)} / M
 \tag{2.26}$$

$$b_{\text{syn}} = \sum_{m=1}^M (q_{\text{syn}}^{(m)} - \bar{q}_{\text{syn}})^2 / (M - 1)
 \tag{2.27}$$

$$\bar{v}_{\text{syn}} = \sum_{m=1}^M v_{\text{syn}}^{(m)} / M.
 \tag{2.28}$$

The pairs of equations (2.26) and (2.10), (2.27) and (2.11), and (2.28) and (2.14) can be viewed as equivalent when $R = 1$.

To complete the derivation for *SynRep-1*, we follow the logic in Raab et al. (2018) and assume that $q_{\text{syn}}^{(m)} \sim N(q^{(m)}, V^{(m)})$. Assuming $V^{(m)} \approx \bar{v}$ for all m , we have

$$(q^{(m)} | \mathcal{D}_{\text{syn}}^{(m)}, \bar{v}) \sim N(q_{\text{syn}}^{(m)}, \bar{v})
 \tag{2.29}$$

$$(\bar{q} | \mathcal{D}_{\text{syn}}, \bar{v}) \sim N(\bar{q}_{\text{syn}}, \bar{v}/M).
 \tag{2.30}$$

We note, however, that one should not assume that $B \approx \bar{v}$ as well. As \mathcal{D} is a complex sample, it yields sampling variances that could differ from the simple random sampling variances associated with \mathcal{D}_{srs} .

Since all the components are approximately normal distributions, $P(Q|\mathcal{D}_{\text{syn}})$ also is approximately a normal distribution. For its expectation, we use (2.7) and (2.30) to obtain

$$E(Q|\mathcal{D}_{\text{syn}}) = E(E(Q|\mathcal{D}_{\text{srs}})|\mathcal{D}_{\text{syn}}) = E(\bar{q}|\mathcal{D}_{\text{syn}}) = \bar{q}_{\text{syn}}. \tag{2.31}$$

For its variance, as with *SynRep-R*, we write the variance in terms of (B, \bar{v}) and then plug in point estimates of these terms. We have

$$\begin{aligned} \text{Var}(Q|\mathcal{D}_{\text{syn}}, B, \bar{v}) &= E((1 + M^{-1})B + \bar{v}/M | \mathcal{D}_{\text{syn}}, B, \bar{v}) + \text{Var}(\bar{q} | \mathcal{D}_{\text{syn}}, B, \bar{v}) \\ &= (1 + M^{-1})B + \bar{v}/M + \bar{v}/M = (1 + M^{-1})B + 2\bar{v}/M. \end{aligned} \tag{2.32}$$

We now define the estimates for (B, \bar{v}) to plug into (2.32). For \bar{v} , we use \bar{v}_{syn} defined in (2.28). This should be reasonable since we are replacing the entire set of each $\mathcal{D}^{(m)}$ with synthetic values. To approximate B , we note that the sampling distribution of a randomly generated $q_{\text{syn}}^{(m)}$ over all steps in the data generation process is $N(Q, B + 2\bar{v})$. Using this fact, we have

$$\frac{\sum_{m=1}^M (q_{\text{syn}}^{(m)} - \bar{q}_{\text{syn}})^2}{B + 2\bar{v}} \sim \chi_{M-1}^2, \tag{2.33}$$

so that

$$E\left(\frac{\sum_{m=1}^M (q_{\text{syn}}^{(m)} - \bar{q}_{\text{syn}})^2}{B + 2\bar{v}}\right) = M - 1. \tag{2.34}$$

Using a method of moments approach and the definition of b_{syn} in (2.27), we have $b_{\text{syn}} \approx B + 2\bar{v}_{\text{syn}}$, so that $B \approx b_{\text{syn}} - 2\bar{v}_{\text{syn}}$.

Thus, we can approximate $\text{Var}(Q|\mathcal{D}_{\text{syn}})$ with the estimate T_m , where

$$T_m = (1 + M^{-1})b_{\text{syn}} - 2\bar{v}_{\text{syn}}. \tag{2.35}$$

We compute approximate 95% intervals for Q as $\bar{q}_{\text{syn}} \pm t_{0.975, M-1} \sqrt{T_m}$. When $T_m < 0$, as an *ad hoc* variance estimate we replace B by \bar{v} in (2.32) and use $T_m^* = (1 + 3/M)\bar{v}_{\text{syn}}$.

3. Simulation studies

In this section, we present simulation studies to illustrate the repeated sampling properties of the inferential methods in Section 2.2 for various finite population quantities.

3.1 Simulation design

We construct a finite population based on data from the Public Use Microdata Sample of the 2021 American Community Survey (United States Bureau of the Census, 2021). This file comprises 3,252,599

individuals, which we treat as a population of size N . The file also has person-level weights (named “PWGTP” in the data file). We do not treat these as survey weights, per se; rather, we treat them as size variables x_i , where $i = 1, \dots, N$, for use in probability proportional to size (PPS) sampling. We also use these constructed size measures to generate two survey variables, (y_{i1}, y_{i2}) , where $i = 1, \dots, N$. Specifically, we let each y_{i1} be a binary variable sampled from a Bernoulli distribution with probability $\Pr(y_{i1} = 1) = \exp(-7 + 2 \log x_i) / (1 + \exp(-7 + 2 \log x_i))$. We let each y_{i2} be a continuous variable sampled from a normal distribution with mean $20 + 50y_{i1}$ and standard deviation 50. We estimate the finite population proportion $\bar{Y}_1 = \sum_{i=1}^N y_{i1} / N \approx 0.765$; the finite population mean $\bar{Y}_2 = \sum_{i=1}^N y_{i2} / N \approx 58.2$; and, the finite population regression coefficient of Y_1 in the linear regression of Y_2 on Y_1 , which is $\beta \approx 50$.

From this population, we sample \mathcal{D} using a PPS sample of size $n = 500$ survey units, setting $\pi_i = nx_i / \sum_{i=1}^N x_i$ and using the function “ppss” in the R package “pps” (Gambino, 2021). Under this PPS sampling design, we expect that unweighted inferences using \mathcal{D} should be badly biased for (\bar{Y}_1, \bar{Y}_2) but perhaps not so for β . We repeat the sampling process to create 1,000 independent realizations of \mathcal{D} .

For each \mathcal{D} , we implement *SynRep-R* and *SynRep-I* with various (M, R) . Specifically, we examine $(M = 4, R = 5)$, $(M = 10, R = 5)$, $(M = 50, R = 5)$, $(M = 10, R = 10)$, $(M = 10, R = 25)$, and $(M = 10, R = 50)$. The choice of R only affects *SynRep-R*. We implement the WFPBB using the “polyapost” package in R (Meeden, Lazar and Geyer, 2020), creating pseudo-populations $(\mathcal{P}^{(1)}, \dots, \mathcal{P}^{(M)})$ each comprising 25,000 individuals. From each $\mathcal{P}^{(m)}$ where $m = 1, \dots, M$, we take a simple random sample of size n to make a corresponding $\mathcal{D}^{(m)}$. To make each synthetic data replicate stemming from each $\mathcal{D}^{(m)}$, we sample n synthetic values for Y_1 using a Bernoulli distribution with probability set to the empirical proportion of Y_1 in $\mathcal{D}^{(m)}$. We sample the corresponding synthetic values of Y_2 from normal distributions with means equal to the predicted values from the regression of Y_2 on Y_1 , computed using the synthetic values of Y_1 and the unbiased estimates of the coefficients computed with $\mathcal{D}^{(m)}$, and variance equal to the unbiased estimate of the regression variance computed with $\mathcal{D}^{(m)}$.

To assist in evaluating the repeated sampling performances of *SynRep-I* and *SynRep-R*, we also use results computed with $\mathcal{P}_{\text{pseudo}}$ and \mathcal{D}_{srs} . Specifically, in each of the 1,000 simulation runs, we define *Pseudo-Pop* as the procedure that uses a point estimator of \bar{Q} and variance estimator of $(1 + 1/M)B$ computed with the WFPBB-generated pseudo-populations $(\mathcal{P}^{(1)}, \dots, \mathcal{P}^{(M)})$. We define *Pseudo-SRS* as the procedure that uses a point estimator of \bar{q} and variance estimator of Raghunathan et al. (2003) computed with $(\mathcal{D}^{(1)}, \dots, \mathcal{D}^{(M)})$. As a comparison against what happens if we disregard the sampling design entirely, we define *SRSsyn* as the procedure that generates synthetic data by using (i) the unweighted sample proportion for Y_1 as the Bernoulli probability to generate n synthetic values of Y_1 and (ii) the unweighted estimates of parameters in the regression of Y_2 on Y_1 as the parameters of the normal distribution to generate the corresponding n synthetic values of Y_2 .

We also evaluate the repeated sampling performances of pseudo-likelihood approaches to making fully synthetic data. For each synthesis model, i.e., the Bernoulli and linear regression models, we start with a likelihood function defined as the product of the contributions from each individual in \mathcal{D} . We create the

pseudo-likelihood by raising each individual's contribution to a power defined by the individual's survey weight. We use these weighted pseudo-likelihoods to estimate synthesis model parameters. We implement this approach using the software *Stan* (Stan Development Team, 2024), which can generate posterior samples of model parameters based on user-specified likelihood functions. We run *Stan* to create four chains of 4,000 iterations and discard the first 2,000 iterations as burn-in. We randomly sample one of the resulting draws and use its parameter values in the Bernoulli and linear regression models to generate the synthetic data. We repeat this process M times and apply the inference rules in Raghunathan et al. (2003). We call this method *Wtreg*. We note that Kim et al. (2021) use the variance estimator in (2.8) from Raghunathan et al. (2003) with $\bar{v} = 0$. Kim et al. (2021) release synthetic populations (where $\bar{v} = 0$) rather than synthetic samples (where $\bar{v} > 0$).

We also consider a modification of *Wtreg* to address potential underestimation of variability in the parameter draws. We call this method *Wtreg-Boot*. First, we take a bootstrap sample of size n from \mathcal{D} . We construct the pseudo-likelihood functions using the bootstrapped data and the calibrated survey weight for each resampled individual. Using this pseudo-likelihood function, we then generate and analyze synthetic data following the steps described for *Wtreg*.

Finally, we define *Direct* as using the unweighted sample mean and standard deviation from \mathcal{D} , i.e., ignoring the survey weights, and *HT* as using the Horvitz and Thompson (1952) estimator and its estimated variance using \mathcal{D} . We use these latter two procedures to assess the importance of accounting for the sampling design in inferences with \mathcal{D} .

Let superscript s index the results from simulation run s , where $s = 1, \dots, 1,000$. For any estimator \hat{q} for any of the methods we examine, we compute the percent bias, $100 \sum_{s=1}^{1,000} (\hat{q}^s - Q) / (1,000Q)$. We compute the proportion of the 1,000 95% confidence intervals based on \hat{q} and its corresponding variance estimate that cover Q . We also compute the ratio of the empirical variance of the 1,000 values of \hat{q} to the empirical variance of the 1,000 values of the *HT* point estimator. To investigate the accuracy of variance estimators, for each method we compute the ratio of the average of the 1,000 variance estimates over its corresponding empirical variance. Finally, to examine the stability of the variance estimator for each method, we compute the standard deviation of the 1,000 variance estimates. We present results for the first four quantities in the main text and for the last quantity in the Appendix.

3.2 Results

We first investigate the properties of *SynRep-R* and *SynRep-I* for the various settings of (M, R) . Figure 3.1, Figure 3.2, and Figure 3.3 display results for \bar{Y}_1 , \bar{Y}_2 , and β , respectively, for these two methods as well as for *Pseudo-Pop* and *Pseudo-SRS*. All four methods offer approximately unbiased point estimates of the three finite population quantities, with simulated percent biases generally around 1% or lower. These small biases originate primarily from the step of completing populations, as the biases in *Pseudo-Pop* are close to the biases in the other three methods. As expected, compared to the variance for *HT*, the simulated variances are increasingly inflated as M decreases. Holding $M = 10$ constant, decreasing R tends to

increase the simulated variances, although the effects are less pronounced than those from decreasing M . The variability in *SynRep-I* results with fixed M reflects Monte Carlo error. Taken together, these results suggest it is preferable to increase M rather than R when keeping MR constant. For example, when we compare *SynRep-R* with $(M = 10, R = 5)$ to *SynRep-I* with $M = 50$, the latter tends to result in smaller empirical variance with closer-to-nominal coverage rates. Similar benefits appear when comparing *SynRep-R* with $(M = 10, R = 25)$ to *SynRep-R* with $(M = 50, R = 5)$. This finding accords with results from Reiter (2008), who considered a similar trade-off for nested multiple imputation for partially synthetic and missing data. We note that using larger values of M also offers smaller variability in the estimated variances, as shown in the Appendix.

Figure 3.1 Repeated sampling properties of *SynRep-I* and *SynRep-R* for \bar{Y}_1 under different numbers of synthetic samples (M) and replicates (R) under a probability proportional to size design.

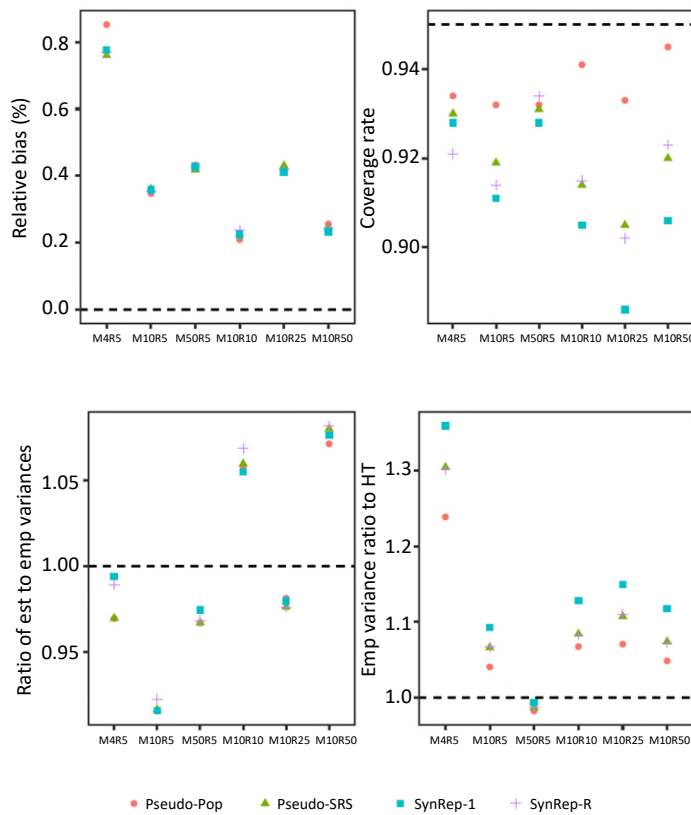
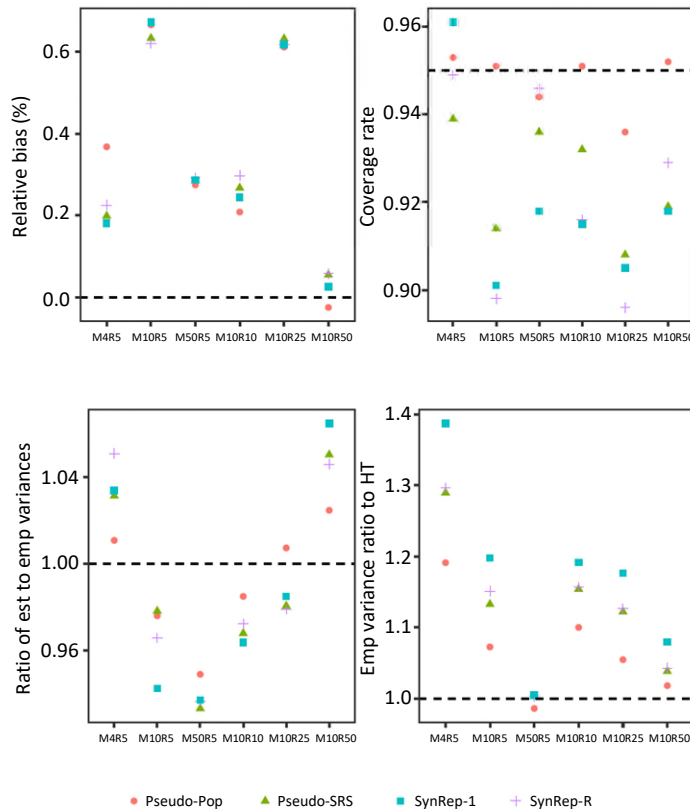


Figure 3.2 Repeated sampling properties of *SynRep-1* and *SynRep-R* for \bar{Y}_2 under different numbers of synthetic samples (M) and replicates (R) under a probability proportional to size design.



By comparing the ratios of the empirical variances to the variances for HT , we can see the effect on efficiency of the steps in the synthesis process. The variances generally increase as we go from *Pseudo-Pop* to *Pseudo-SRS* to *SynRep-R* or *SynRep-1*; that is, they increase as we add more steps that involve randomness. The variances for *SynRep-R* generally are slightly smaller than those for *SynRep-1*, reflecting the benefit for efficiency of the additional information from MR rather than M synthetic data sets. We note that the variance inflation from using synthetic data procedures versus HT largely disappears when $M = 50$.

Across all four synthetic data methods, the average variance estimates are reasonably similar to the empirical variances. Disparities from ratios of one apparently stem, once again, mainly from the step of completing the populations. The confidence interval coverage rates range from a low of 88% to a high of 96%, with most slightly below nominal. Coverage rates for *SynRep-R* and *SynRep-1* tend to be highest when $M = 50$, further reflecting the benefits of using a larger M . For $M \geq 10$, the coverage rates for *SynRep-R* tend to be higher than those for *SynRep-1*, although the difference is typically only a point or two.

The combining rules in (2.35) and (2.25) do result in negative variance estimates, as evident in Table 3.1. In the simulations, we use T_r^* and T_m^* to make confidence intervals when needed. As M increases, the number of negative variance estimates decreases. In fact, when $M = 50$, all of the variance estimates are

positive, offering additional support for making M large. The estimates of b_{syn} become less variable as M increases, which helps avoid the negative variances. Negative variance rates tend to be lower for *SynRep-R* than for *SynRep-I*, reflecting the benefits of increased datasets to estimate variance parameters. Although not shown in Table 3.1, the negative variance rates when $M = 10$ do not change much as we increase $R \geq 5$. We note that the negative variance rates for *SynReg-R* are similar to those for *Pseudo-SRS*. Evidently, when MR is large, the information available in \mathcal{D}_{syn} to estimate b_{syn} is on par with the information available in \mathcal{D}_{srs} .

Figure 3.3 Repeated sampling properties of *SynRep-I* and *SynRep-R* for β under different numbers of synthetic samples (M) and replicates (R) under a probability proportional to size design.

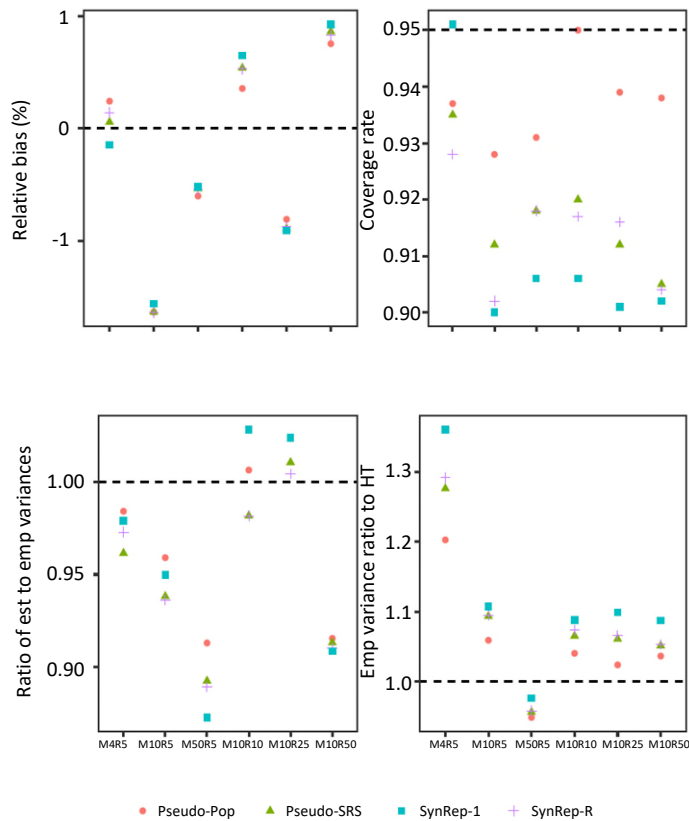
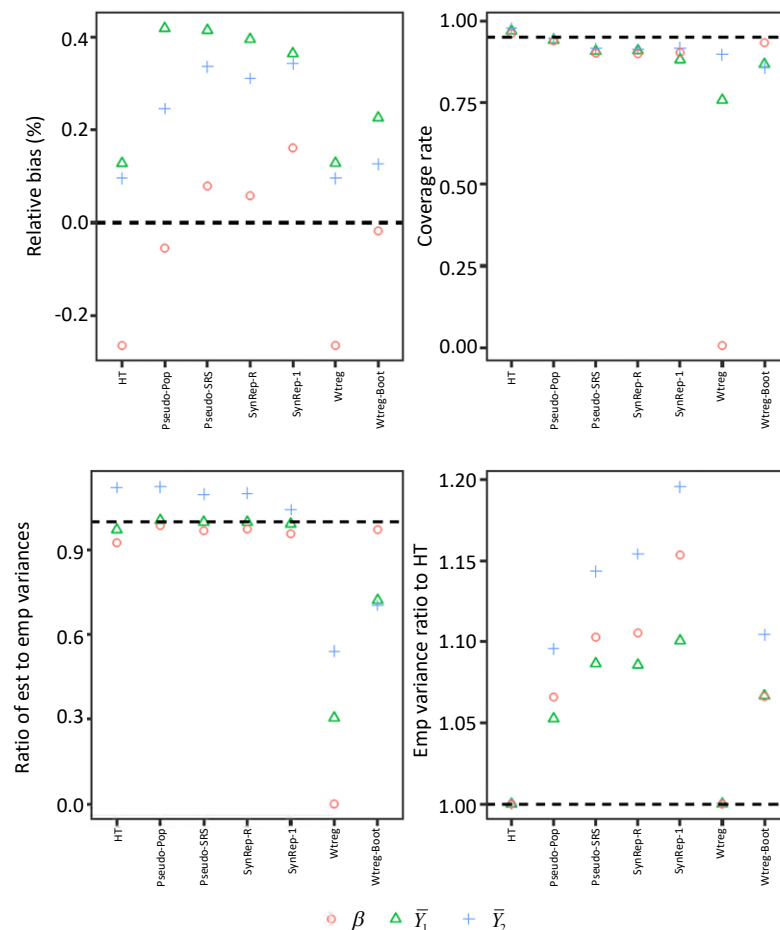


Table 3.1 Proportion of negative variance estimates in the PPS simulation studies. When $M = 50$, all variance estimates are positive.

(M, R)	Method	\bar{Y}_1	\bar{Y}_2	β
M4R5	<i>Pseudo-SRS</i>	0.09	0.15	0.13
M4R5	<i>SynRep-R</i>	0.11	0.17	0.13
M4R5	<i>SynRep-I</i>	0.17	0.26	0.22
M10R5	<i>Pseudo-SRS</i>	0.01	0.02	0.02
M10R5	<i>SynRep-R</i>	0.01	0.03	0.02
M10R5	<i>SynRep-I</i>	0.04	0.09	0.07

We next turn to compare *SynRep-R* and *SynRep-I* with other approaches, particularly *Wtreg*, *Wtreg-Boot*, and *SRSsyn*. Here, we set $M=10$ and, where relevant, $R=10$, and draw 500 repeated samples. Figure 3.4 summarizes the repeated sampling performances of the methods that account for survey weights. For all these methods, the point estimators have simulated percent biases that typically are negligible. For *SynRep-R* and *SynRep-I*, the average variance estimates are close to their corresponding empirical variances, and the coverage rates are close to nominal. For *Wtreg* and *Wtreg-Boot*, the variance estimators can underestimate the corresponding empirical variances severely, especially for \bar{Y}_1 and \bar{Y}_2 , resulting in confidence interval coverage rates that can be substantially lower than the nominal 95% level. The bootstrap step in *Wtreg-Boot* results in more reliable variance estimates compared to *Wtreg*, but *Wtreg-Boot* is not as well calibrated as *SynRep-R* and *SynRep-I*, which have closer to nominal coverage rates. As expected, *HT* results in accurate estimates with near nominal coverage rates. We note that Figure 3.4 does not display results for *Direct* and *SRSsyn* because they perform poorly for \bar{Y}_1 and \bar{Y}_2 . For these two methods, the simulated biases for \bar{Y}_1 and for \bar{Y}_2 are around 16% and 11%, respectively, with coverage rates near 0 and near 30%, respectively. These results emphasize the importance of accounting for informative designs when generating fully synthetic data that can be analyzed as simple random samples.

Figure 3.4 Repeated sampling properties of different quantities and procedures with $M=10$ synthetic samples and $R=10$ replicates under a probability proportional to size design.



Overall, the simulation studies suggest that *SynRep-R* and *SynRep-I* can provide approximately valid inferences, and they are superior inferentially to fully synthetic data that ignore the complex design. The Appendix also includes results of simulation studies where we sample \mathcal{D} via simple random samples. These confirm that the combining rules offer reasonable performance even without unequal probabilities of selection.

4. Illustration with ACS data

We illustrate *SynRep-R* and *SynRep-I* by letting \mathcal{D} be a subset of data from the 2021 ACS Public Use Microdata Sample for $n = 84,128$ individuals from the state of Michigan. The variables for our illustration include each participant's person-level weight, age, and total income. To mimic the variables in the simulations, we create a binary indicator Y_1 from age that equals one when someone is at least 65 years old; we refer to this indicator as senior status. For purposes of synthesis, we transform income by taking its cubic root. The synthesis models are then a Bernoulli distribution for Y_1 and a linear regression of the cubic root of total income on Y_1 . After synthesizing values of the cubic root of income, we raise them to the third power to get incomes on the original scale. We implement each method following the procedures from Section 3. For *SynRep-R* and *SynRep-I*, we set $M = 10$ and $R = 10$.

As population quantities, we estimate the population proportion of senior status individuals \bar{Y}_1 , the population mean of the income values, \bar{Y}_2 , and the coefficient β of Y_1 in the linear regression model of the cubic-root transformed income on senior status.

Figure 4.1 presents the point estimates and 95% confidence intervals for the three population quantities. Since *Direct* and *SRSsyn* ignore the sample design, they result in relatively inaccurate results, especially for \bar{Y}_1 . In contrast, the point estimates for the synthetic data methods that account for survey weights are closer to the *HT* point estimates. Additionally, the 95% confidence intervals for these methods largely overlap with the *HT* confidence intervals. We note, however, that *Wtreg* appears to suffer from underestimation of variance, particularly for β . Additionally, the confidence intervals for the pseudo-likelihood approaches can be narrower than those for *HT*, *SynRep-R*, and *SynRep-I*.

We also can examine potential disclosure risks for the synthetic data methods. Here, we mimic an attack scenario described by Kim et al. (2021), in which we consider an adversary who uses the synthetic data to estimate the largest income value in \mathcal{D} . Specifically, we examine differences between the maximum synthetic income in each synthetic dataset and the maximum income in \mathcal{D} . This evaluation is not intended to illustrate a rigorous and thorough process for assessing disclosure risks. Rather, we use this attack scenario mainly to compare the different synthesis procedures.

Table 4.1 presents the distributions of the differences for the synthesis methods that account for the survey design. Overall, the results are reasonably similar across the methods, suggesting they offer similar levels of protection in this scenario. All result in substantial differences between the largest synthetic and

observed incomes. The results suggest that an adversary taking this attack strategy is not likely to estimate the largest income accurately.

Figure 4.1 Point estimates and 95% confidence intervals for \bar{Y}_1 , \bar{Y}_2 , and β in the ACS data illustration. Results based on $M = 10$ synthetic samples and $R = 10$ replicates.

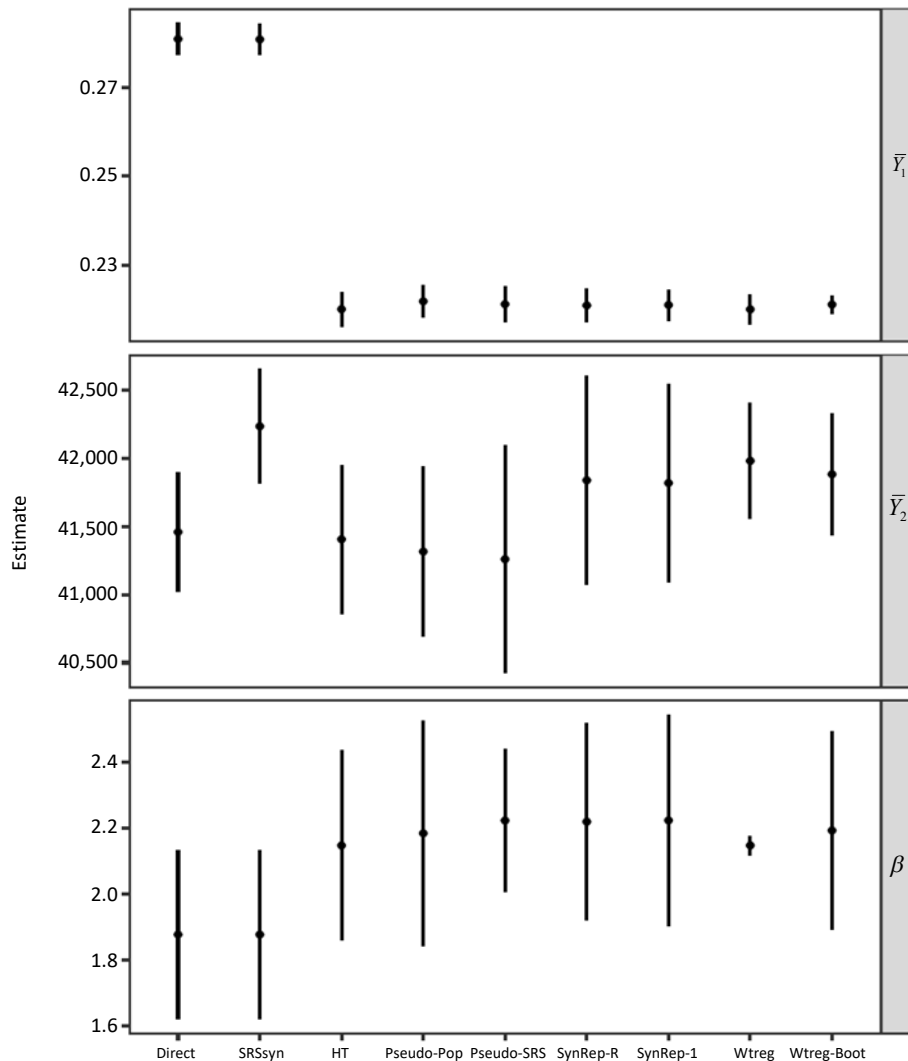


Table 4.1
Summaries of the differences (\$) in the largest income value in the synthetic and American Community Survey data. The actual largest value is \$1,029,000.

Method	Min.	1st Quartile	Median	Mean	3rd Quartile	Max.
<i>SynRep-R</i>	-424,323	-298,230	-252,984	-214,476	-139,874	465,380
<i>SynRep-1</i>	-371,466	-297,180	-287,199	-268,711	-267,428	-40,405
<i>Wtreg</i>	-440,253	-297,689	-242,095	-218,810	-159,766	707,411
<i>Wtreg-Boot</i>	-410,354	-275,398	-209,513	-174,444	-139,109	133,759

5. Discussion

SynRep-R and *SynRep-I* represent a general strategy for constructing fully synthetic data that account for complex sample designs: use the WFPBB to “undo” the design, then replace the confidential values with simulated values. Releasing multiple synthetic data sets, i.e., setting $MR > 1$, can increase statistical efficiency and facilitate variance estimation. However, agencies also can use the WFPBB approach with $MR = 1$. Although releasing a single synthetic data set may not enable approximately valid variance estimation for complex surveys, it still can be useful in certain settings, e.g., when the synthetic data are intended for code training or exploratory analyses where variance estimation is not essential.

As noted by a reviewer, several agencies implementing synthetic data approaches also provide means for users to check the quality of their synthetic data inferences. For example, users can submit their code to the agency that released the synthetic data, which then can run the code and report back disclosure-protected outputs to the user. This is known as validation of results (Barrientos, Bolton, Balmat, Reiter, de Figueiredo, Machanavajhala, Chen, Kneifel and DeLong, 2018). Alternatively, users can submit queries to a server that computes an analysis of the confidential and synthetic data, and reports back measures of similarity of the two analysis results, e.g., the overlap in the confidence intervals (Karr, Kohnen, Oganian, Reiter and Sanil, 2006). This is known as verification of results (Barrientos et al., 2018). With validation or verification, users of *SynRep-R* and *SynRep-I* may face an additional burden. If the agency directly runs the users’ submitted analysis code, the user may need to specify a survey-weighted version of the code for validation, even though they have used a simple random sample analysis for synthetic data. Of course, for many analyses, e.g., regression modeling, some users forego weighted analyses, in which case the issue is moot. It is also possible for the agency to automate validation or verification, in which case it may be able to turn users’ submitted queries into survey-weighted versions automatically in the background; this is an area for future research.

We chose to develop methods that enable agencies to follow the idea in Rubin (1993): release data that can be analyzed as simple random samples. This can make analyses easier for users, as they do not have to figure out how to deal with any weights on the file, e.g., in variance estimation. Releasing simple random samples could also help mitigate disclosure risks that may arise from releasing survey weights. For example, if the weights released on the synthetic files are sampled directly from the weight values in \mathcal{D} without alteration, the weights may reveal information about data subjects that is considered an unacceptable disclosure risk (Fienberg, 2010). Finally, releasing simple random samples avoids the need to estimate relationships between the weights and the outcome variables, which could be complicated in practice. Nonetheless, it would be interesting to compare risk and utility profiles of these approaches with those developed here.

There are many other topics related to the general strategy worth further investigation. First, in practice, survey weights can be highly variable and may not be strongly related to the survey variables of interest; this can cause survey-weighted estimates to have inflated variances. This can be remedied somewhat, for example, by using model-based approaches to smooth the weights (Beaumont, 2008; Xia and Elliott, 2016; Si, Trangucci, Gabry and Gelman, 2020). Synthetic data generation based on the WFPBB (or any other

approach) is not immune to these weighting issues. Thus, it would be interesting to examine if and how the synthesis model can reduce the effects of variance inflation from extreme weights.

Second, we focus on developing the fully synthetic data framework and corresponding combining rules, using simple settings and synthesis models to illustrate the methods. Conceptually, agencies can apply *SynRep-R* and *SynRep-I* to multivariate data and for various estimands of interest, e.g., subdomain means and multiple regression coefficients. In such cases, it may be advantageous to use flexible modeling approaches, such as tree-based models or other machine learning algorithms. Future work could investigate the performance of these synthesizers in combination with the pseudo-population and pseudo-SRS generation steps.

Third, we derive the combining rules assuming the original survey data are complete. Agencies could impute missing survey data and generate synthetic replicates simultaneously, possibly accounting for the complex design in the imputation model and synthesis approach. This strategy may necessitate new combining rules akin to those in Reiter (2004).

Fourth, we present *ad hoc* adjustments to deal with negative values of the variance estimates. We may be able to improve on those adjustments. For example, we may be able to adapt the strategy in Si and Reiter (2011), who develop inferential methods for fully synthetic data based on sampling from the distributions used in the derivations of the combining rules. Additionally, as pointed out by a reviewer, it may be beneficial to use the insight of Raab et al. (2018) for the sampling and synthesis components of the derivation in *SynRep-R*. This results in an alternative variance estimator, $(1 + M^{-1})b_{\text{syn}} - (1 + R^{-1})\bar{v}_{\text{syn}}$. Future work can investigate the performance of these alternative inference methods.

Fifth, it would be informative to generalize the implementation of *SynRep-R* and *SynRep-I* to other complex designs, such as the stratified multi-stage cluster sampling designs that are common in practice. Zhou, Elliott and Raghunathan (2016) have extended the WFPBB to account for strata, clustering, and survey weights in synthetic population generation. We expect that one could take simple random samples from these pseudo-populations and generate synthetic replicates, possibly using synthesis models that capture design information as suggested in Reiter (2002), and extend the combining rules presented here. It would be a natural extension to comprehensively assess the repeated sampling performances of *SynRep-R* and *SynRep-I* in such multi-stage complex samples.

Lastly, it would be useful to develop principled approaches to measuring disclosure risks for these methods. For *SynRep-R* and *SynRep-I*, conceptually one could estimate an adversary's posterior distribution for confidential data values given the released synthetic values, e.g., as described for simple settings in Reiter, Wang and Zhang (2014) and Hu, Reiter and Wang (2015). However, this would be computationally challenging in practice. One would need to account for the entire synthetic data generation process – including the bootstrapping, sampling, and synthesis – when computing this posterior distribution. Indeed, as far as we are aware, agencies that release synthetic data use *ad hoc* approaches to assessing disclosure risks, such as comparing the similarity of outlier values in the confidential and synthetic data as we illustrated here (Kinney, Reiter and Miranda, 2014). Developing disclosure risk methods is a major area for future research for all approaches to generating fully synthetic data.

Acknowledgements

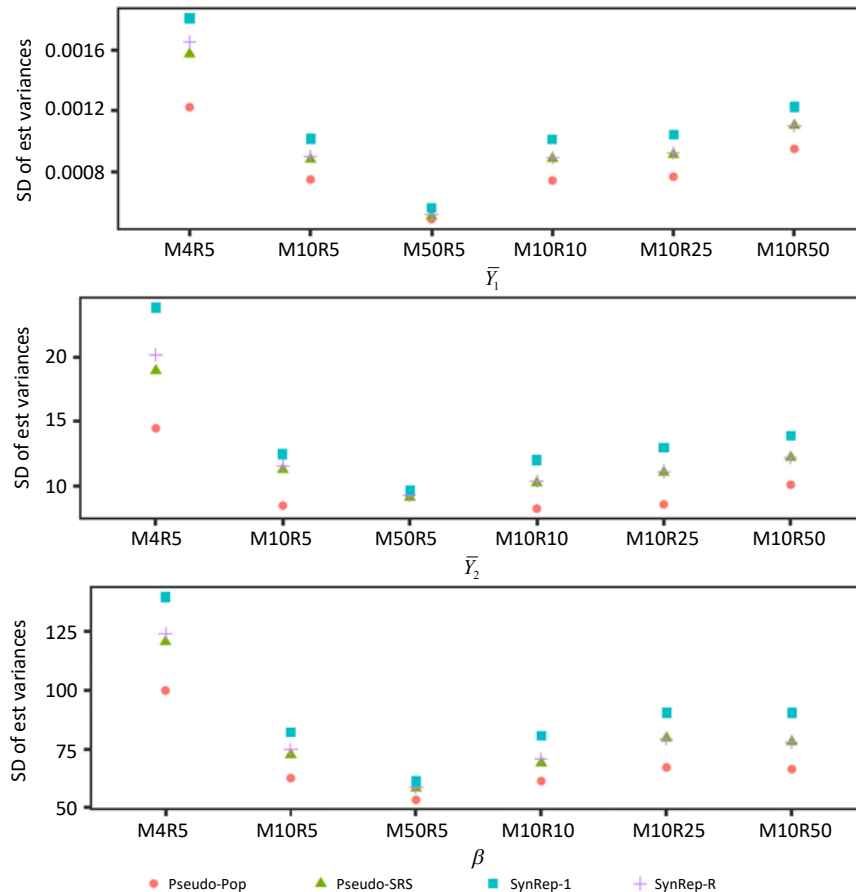
The work was funded by the U.S. National Science Foundation grant (SES 2217456) and a pilot project from the Michigan Center on the Demography of Aging with funding from the National Institute on Aging (P30 AG012846).

Appendix

A. Additional simulation results

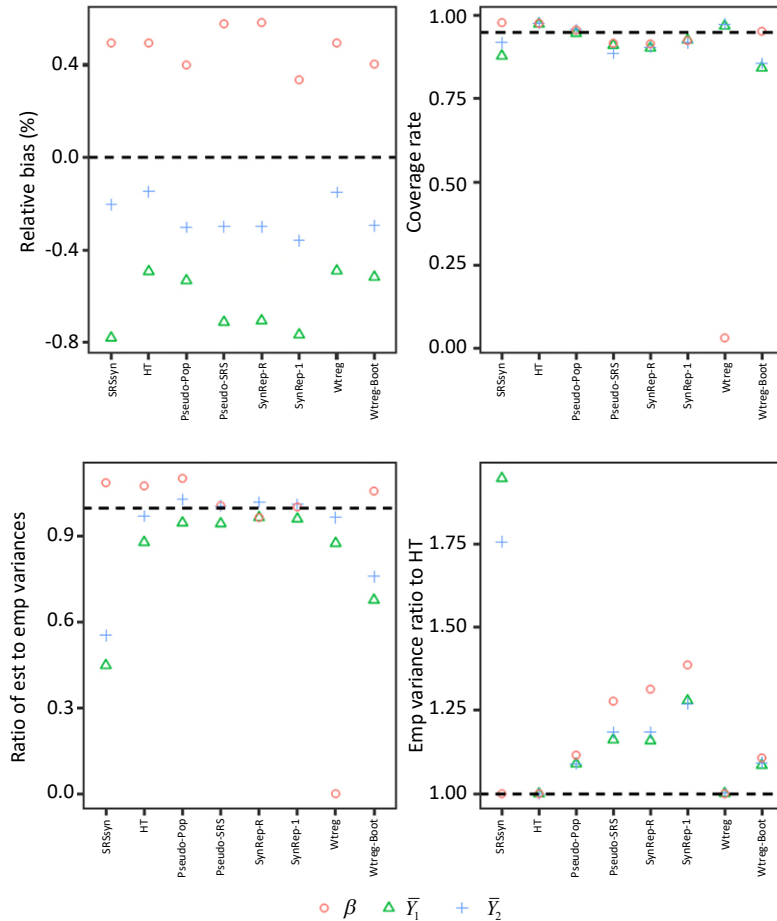
Figure A.1 displays the variability of the 1,000 values of estimated variances of the point estimators for β , \bar{Y}_1 , and \bar{Y}_2 for the simulation with the PPS design. The variability tends to decrease with M . Increasing R when M is held constant seems not to have much impact on the stability of the results. We see increased variability as the procedures introduce more steps that involve randomness; that is, as we go from *Pseudo-Pop* to *Pseudo-SRS* to *SynRep-R* and *SynRep-I*. The variability tends to be largest for *SynRep-I*.

Figure A.1 Standard deviation (SD) of estimated (est) variances of different population quantities with different procedures for different numbers of synthetic samples (M) and replicates (R) under a probability proportional to size design.



As another check of the validity of the combining rules, we repeat the simulations from Section 3 using a SRS in place of a PPS design. Specifically, we use the population described in Section 3.1, but we use a SRS of $n = 500$ records for each \mathcal{D} . Figure A.2 displays the results. Overall, the performances of *SynRep-R* and *SynRep-I* mirror the patterns seen for the PPS design in Section 3.

Figure A.2 Repeated sampling properties of different quantities and procedures with $M = 10$ synthetic samples and $R = 10$ replicates under a SRS design.



References

Barrientos, A.F., Bolton, A., Balmat, T., Reiter, J.P., de Figueiredo, J.M., Machanavajjhala, A., Chen, Y., Kneifel, C. and DeLong, M. (2018). Providing access to confidential research data through synthesis and verification: An application to data on employees of the U.S. federal government. *Annals of Applied Statistics*, 12, 1124-1156.

- Beaumont, J.-F. (2008). A new approach to weighting and inference in sample surveys. *Biometrika*, 95, 3, 539-553.
- Dong, Q., Elliott, M.R. and Raghunathan, T.E. (2014). [A nonparametric method to generate synthetic populations to adjust for complex sampling design features](#). *Survey Methodology*, 40, 1, 29-46. Paper available at <https://www150.statcan.gc.ca/n1/en/pub/12-001-x/2014001/article/14003-eng.pdf>.
- Drechsler, J. (2011). *Synthetic Datasets for Statistical Disclosure Control*. New York: Springer.
- Drechsler, J., and Reiter, J.P. (2010). Sampling with synthesis: A new approach for releasing public use census microdata. *Journal of the American Statistical Association*, 105, 1347-1357.
- Fienberg, S.E. (2010). The relevance or irrelevance of weights for confidentiality and statistical analyses. *Journal of Privacy and Confidentiality*, 1, 183-195.
- Gambino, J.G. (2021). R package pps: PPS Sampling. <https://cran.r-project.org/web/packages/pps/index.html>.
- Horvitz, D.G., and Thompson, D.J. (1952). A generalization of sampling without replacement from a finite universe. *Journal of the American Statistical Association*, 47(260), 663-685.
- Hu, J., Reiter, J.P. and Wang, Q. (2015). Disclosure risk evaluation for fully synthetic data. In *Privacy in Statistical Databases*, (Ed., J. Domingo-Ferrer), 185-199. Heidelberg: Springer.
- Karr, A.F., Kohnen, C.N., Oganian, A., Reiter, J.P. and Sanil, A.P. (2006). A framework for evaluating the utility of data altered to protect confidentiality. *The American Statistician*, 60, 224-232.
- Kim, H.J., Drechsler, J. and Thompson, K.J. (2021). Synthetic microdata for establishment surveys under informative sampling. *Journal of Royal Statistical Society, Series A*, 184, 255-281.
- Kinney, S.K., Reiter, J.P. and Miranda, J. (2014). Synlbd 2.0: Improving the Synthetic Longitudinal Business Database. *Statistical Journal of the International Association for Official Statistics*, 30, 129-135.
- Kinney, S.K., Reiter, J.P., Reznick, A.P., Miranda, J., Jarmin, R.S. and Abowd, J.M. (2011). Towards unrestricted public use business microdata: The synthetic Longitudinal Business Database. *International Statistical Review*, 79, 363-384.
- Little, R.J. (2004). To model or not to model? Competing modes of inference for finite population sampling inference for finite population sampling. *Journal of the American Statistical Association*, 99, 546-556.
- Meeden, G., Lazar, R. and Geyer, C.J. (2020). R package polyapost: Simulating from the Polya posterior. <https://cran.r-project.org/web/packages/polyapost/index.html>.

- Mitra, R., and Reiter, J.P. (2006). Adjusting survey weights when altering identifying design variables via synthetic data. In *Privacy in Statistical Databases*, (Eds., J. Domingo-Ferrer and L. Franconi), 177-188. New York: Springer-Verlag.
- Pfeffermann, D. (1993). The role of sampling weights when modeling survey data. *International Statistical Review*, 61(2), 317-337.
- Pfeffermann, D. (2011). [Modelling of complex survey data: Why model? Why is it a problem? How can we approach it?](https://www150.statcan.gc.ca/n1/en/pub/12-001-x/2011002/article/11602-eng.pdf) *Survey Methodology*, 37, 2, 115-136. Paper available at <https://www150.statcan.gc.ca/n1/en/pub/12-001-x/2011002/article/11602-eng.pdf>.
- Raab, G.M., Nowok, B. and Dibben, C. (2018). Practical data synthesis for large samples. *Journal of Privacy and Confidentiality*, 7(3), 67-97.
- Raghunathan, T.E. (2021). Synthetic data. *Annual Review of Statistics and Its Application*, 8, 129-140.
- Raghunathan, T.E., Reiter, J.P. and Rubin, D.B. (2003). Multiple imputation for statistical disclosure limitation. *Journal of Official Statistics*, 19, 1-16.
- Reiter, J.P. (2002). Satisfying disclosure restrictions with synthetic data sets. *Journal of Official Statistics*, 18, 531-544.
- Reiter, J.P. (2003). [Inference for partially synthetic, public use microdata sets](https://www150.statcan.gc.ca/n1/en/pub/12-001-x/2003002/article/6785-eng.pdf). *Survey Methodology*, 29, 2, 181-188. Paper available at <https://www150.statcan.gc.ca/n1/en/pub/12-001-x/2003002/article/6785-eng.pdf>.
- Reiter, J.P. (2004). [Simultaneous use of multiple imputation for missing data and disclosure limitation](https://www150.statcan.gc.ca/n1/en/pub/12-001-x/2004002/article/7755-eng.pdf). *Survey Methodology*, 30, 2, 235-242. Paper available at <https://www150.statcan.gc.ca/n1/en/pub/12-001-x/2004002/article/7755-eng.pdf>.
- Reiter, J.P. (2005a). Releasing multiply-imputed, synthetic public use microdata: An illustration and empirical study. *Journal of the Royal Statistical Society, Series A*, 168, 185-205.
- Reiter, J.P. (2005b). Significance tests for multi-component estimands from multiply imputed, synthetic microdata. *Journal of Statistical Planning and Inference*, 131, 365-377.
- Reiter, J.P. (2005c). Using CART to generate partially synthetic, public use microdata. *Journal of Official Statistics*, 21, 441-462.
- Reiter, J.P. (2008). Selecting the number of imputed datasets when using multiple imputation for missing data and disclosure limitation. *Statistics and Probability Letters*, 78, 15-20.

- Reiter, J.P. (2009). Using multiple imputation to integrate and disseminate confidential microdata. *International Statistical Review*, 77, 179-195.
- Reiter, J.P., and Drechsler, J. (2010). Releasing multiply-imputed, synthetic data generated in two stages to protect confidentiality. *Statistica Sinica*, 20, 405-422.
- Reiter, J.P., and Kinney, S.K. (2012). Inferentially valid, partially synthetic data: Generating from posterior predictive distributions not necessary. *Journal of Official Statistics*, 28, 583-590.
- Reiter, J.P., and Raghunathan, T.E. (2007). The multiple adaptations of multiple imputation. *Journal of the American Statistical Association*, 102, 1462-1471.
- Reiter, J.P., Raghunathan, T.E. and Kinney, S.K. (2006). [The importance of modeling the sampling design in multiple imputation for missing data](#). *Survey Methodology*, 32, 2, 143-149. Paper available at <https://www150.statcan.gc.ca/n1/en/pub/12-001-x/2006002/article/9548-eng.pdf>.
- Reiter, J.P., Wang, Q. and Zhang, B. (2014). Bayesian estimation of disclosure risks in multiply imputed, synthetic data. *Journal of Privacy and Confidentiality* 6, Article 2.
- Rubin, D.B. (1981). The Bayesian bootstrap. *The Annals of Statistics*, 9, 130-134.
- Rubin, D.B. (1987). *Multiple Imputation for Nonresponse in Surveys*. New York: John Wiley & Sons, Inc.
- Rubin, D.B. (1993). Discussion: Statistical disclosure limitation. *Journal of Official Statistics*, 9, 462-468.
- Savitsky, T.D., and Toth, D. (2016). Bayesian estimation under informative sampling. *Electronic Journal of Statistics*, 10, 1677-1708.
- Si, Y., and Reiter, J.P. (2011). A comparison of posterior simulation and inference by combining rules for multiple imputation. *Journal of Statistical Theory and Practice*, 5, 335-347.
- Si, Y., Trangucci, R., Gabry, J.S. and Gelman, A. (2020). [Bayesian hierarchical weighting adjustment and survey inference](#). *Survey Methodology*, 46, 2, 181-214. Paper available at <https://www150.statcan.gc.ca/n1/en/pub/12-001-x/2020002/article/00003-eng.pdf>.
- Stan Development Team (2024). Stan: A C++ library for probability and sampling. <http://mc-stan.org>.
- United Nations Economic Commission for Europe (2022). Synthetic Data for National Statistical Organizations. <https://statswiki.unece.org/display/SDS/Synthetic+Data+Sets+public?preview=%2F282330193%2F330369384%2FHLG-MOS+Synthetic+Data+Guide.docx>. Accessed: 2022-01-12.
- United States Bureau of the Census (2021). Accessing American Community Survey PUMS data. <https://www.census.gov/programs-surveys/acs/microdata/access.html>.

Williams, M.R., and Savitsky, T.D. (2021). Uncertainty estimation for pseudo-Bayesian inference under complex sampling. *International Statistical Review*, 89, 72-107.

Xia, X., and Elliott, M.R. (2016). Weight smoothing for generalized linear models using a Laplace prior. *Journal of Official Statistics*, 32, 507-539.

Zhou, H., Elliott, M.R. and Raghunathan, T.E. (2016). Synthetic multiple-imputation procedure for multistage complex samples. *Journal of Official Statistics*, 32, 231-256.