# Combining multiple observational data sources to estimate causal effects

Shu Yang[*1] and Peng Ding[*2]
[*1]Department of Statistics, North Carolina State University
[*2]Department of Statistics, University of California, Berkeley

## Abstract

The era of big data has witnessed an increasing availability of multiple data sources for statistical analyses. We consider estimation of causal effects combining big main data with unmeasured confounders and smaller validation data with supplementary information on these confounders. Under the unconfoundedness assumption with completely observed confounders, the smaller validation data allow for constructing consistent estimators for causal effects, but the big main data can only give error-prone estimators in general. However, by leveraging the information in the big main data in a principled way, we can improve the estimation efficiencies yet preserve the consistencies of the initial estimators based solely on the validation data. Our framework applies to asymptotically normal estimators, including the commonly-used regression imputation, weighting, and matching estimators, and does not require a correct specification of the model relating the unmeasured confounders to the observed variables. We also propose appropriate bootstrap procedures, which makes our method straightforward to implement using software routines for existing estimators.

*Keywords:* Calibration; Causal inference; Inverse probability weighting; Missing confounder; Two-phase sampling.

1

# 1 Introduction

Unmeasured confounding is an important and common problem in observational studies. Many methods have been proposed to deal with unmeasured confounding in causal inference, such as sensitivity analyses (e.g. Rosenbaum and Rubin, 1983a), instrumental variable approaches (e.g. Angrist et al., 1996), etc. However, sensitivity analyses cannot provide point estimation, and valid instrumental variables are often difficult to find in practice. We consider the setting where external validation data provide additional information on unmeasured confounders. To be more precise, the study includes a large main dataset representing the population of interest with unmeasured confounders and a smaller validation dataset with additional information about these confounders.

Our framework covers two common types of studies. First, we have a large main dataset, and then collect more information on unmeasured confounders for a subset of units, e.g., using a two-phase sampling design (Neyman, 1938, Cochran, 2007, Wang et al., 2009). Second, we have a smaller but carefully designed validation dataset with rich covariates, and then link it to a larger main dataset with fewer covariates. The second type of data is now ubiquitous. In the era of big data, extremely large data have become available for research purposes, such as electronic health records, claims databases, disease data registries, census data, to name a few (e.g., Imbens and Lancaster, 1994, Schneeweiss et al., 2005, Chatterjee et al., 2016). Although these datasets might not contain full confounder information that guarantees consistent causal effect estimation, they can be useful to increase efficiencies of statistical analyses.

In causal inference, Stürmer et al. (2005) propose a propensity score calibration method when the main data contain the outcome and an error-prone propensity score based on partial confounders, and the validation data supplement a gold standard propensity score based on all confounders. Stürmer et al. (2005) then apply a regression calibration technique to correct for the measurement error from the error-prone propensity score. This approach does not require the validation data to contain the outcome variable. However, this approach relies on the *surrogacy property* entailing that the outcome variable is conditionally independent of the error-prone propensity score given the gold standard propensity score and treatment. This surrogacy property is difficult to justify in practice, and its violations can lead to substantial biases (Stürmer et al., 2007, Lunt et al., 2012). Under the Bayesian framework, McCandless et al. (2012) specify a full parametric model of the joint distribution for the main and validation data, and treat the gold standard propensity score as a missing variable in the main data. Antonelli et al. (2017) combine the ideas of Bayesian model averaging, confounder selection, and missing data imputation into a single framework in this context. Enders et al. (2018) use simulation to show that multiple imputation is more robust than two-phase logistic regression against misspecification of imputation models. Lin and Chen (2014) develop a two-stage calibration method, which summarizes the confounding information through propensity scores and combines the results from the main and validation data. Their two-stage calibration focuses on the regression context with a correctly specified outcome model. Unfortunately, regression parameters, especially in the logistic regression model used by Lin and Chen (2014), may not be the causal parameters of interest in general (Freedman, 2008).

In this article, we propose a general framework to estimate causal effects in the setting where the big main data have unmeasured confounders, but the smaller external validation data provide supplementary information on these confounders. Under the assumption of ignorable treatment assignment, causal effects can be identified and estimated from the validation data, using commonly-used estimators, such as regression imputation, (augmented) inverse probability weighting (Horvitz and Thompson 1952, Rosenbaum and Rubin 1983b; Robins et al., 1994, Bang and Robins, 2005, Cao et al., 2009), and matching (e.g., Rubin, 1973, Rosenbaum, 1989, Heckman et al., 1997, Hi-

rano et al., 2003, Hansen, 2004, Rubin, 2006, Abadie and Imbens, 2006, Stuart, 2010, Abadie and Imbens, 2016). However, these estimators based solely on the validation data may not be efficient. We leverage the correlation between the initial estimator from the validation data and the error-prone estimator from the main data to improve the efficiency over the initial estimator. This idea is similar to the two-stage calibration in Lin and Chen (2014); however, their method focuses only on regression parameters and requires the validation data to be a simple random sample from the main data. Alternatively, the empirical likelihood is also an attractive approach to combine multiple data sources (Chen and Sitter, 1999, Qin, 2000, Chen et al., 2002 and Chen et al., 2003). However, the empirical likelihood approach needs sophisticated programming, and its computation can be heavy when data become large. Our method is practically simple, because we only need to compute commonly-used estimators that can be easily implemented by existing software routines. More-over, Lin and Chen (2014) and the empirical likelihood approach can only deal with regular and asymptotically linear (RAL) estimators often formulated by moment conditions, but our framework can also deal with non-RAL estimators, such as matching estimators. We also propose a unified bootstrap procedure based on resampling the linear expansions of the estimators, which is simple to implement and works for both RAL and matching estimators.

Furthermore, we relax the assumption that the validation data is a random sample from the study population of interest. We also link the proposed method to existing methods for missing data, viewing the additional confounders as missing values for units outside of the validation data. In contrast to most existing methods in the missing data literature, the proposed method does not need to specify the missing data model relating the unmeasured confounders with the observed variables.

For simplicity of exposition, we use "IID" for "identically and independently distributed", $1(\cdot)$ for the indicator function, $\xi^{\otimes 2} = \xi\xi^{\mathrm{T}}$ for a vector or matrix $\xi$, "plim" for the probability limit of a random sequence, and $A_n \cong B_n$ for two random sequences satisfying $A_n = B_n + o_P(n^{-1/2})$ with $n$ being the generic sample size. We relegate all regularity conditions for asymptotic analyses to the on-line supplementary material.

## 2 Basic setup

### 2.1 Notation: causal effect and two data sources

Following Neyman (1923) and Rubin (1974), we use the potential outcomes framework to define causal effects. Suppose that the treatment is a binary variable $A \in \{0, 1\}$, with 0 and 1 being the labels for control and active treatments, respectively. For each level of treatment $a \in \{0, 1\}$, we assume that there exists a potential outcome $Y(a)$, representing the outcome had the subject, possibly contrary to the fact, been given treatment $a$. The observed outcome is $Y = Y(A) = AY(1) + (1 - A)Y(0)$. Let a vector of pretreatment covariates be $(X, U)$, where $X$ is observed for all units, but $U$ may not be observed for some units.

Although we can extend our discussion to multiple data sources, for simplicity of exposition, we first consider a study with two data sources. The validation data have observations $\mathcal{O}_2 = \{(A_j, X_j, U_j, Y_j) : j \in \mathcal{S}_2\}$ with sample size $n_2 = |\mathcal{S}_2|$. The main data have observations $\mathcal{O}_1 = \{(A_i, X_i, Y_i) : i \in \mathcal{S}_1 \backslash \mathcal{S}_2\} \cup \mathcal{O}_2$ with sample size $n_1 = |\mathcal{S}_1|$. In our formulation, we consider the case with $\mathcal{S}_2 \subset \mathcal{S}_1$, and let $\rho = \lim_{n_2 \to \infty} n_2/n_1 \in [0, 1]$. If one has two separate main and validation datasets, the main dataset in our context combines these two datasets. Although the main dataset is larger, i.e., $n_1 > n_2$, it does not contain full information on important covariates $U$. Under

3

a superpopulation model, we assume that $\{A_i, X_i, U_i, Y_i(0), Y_i(1) : i \in \mathcal{S}_1\}$ are IID for all $i \in \mathcal{S}_1$, and therefore the observations in $\mathcal{O}_1$ are also IID. The following assumption links the main and validation data.

**Assumption 1** *The index set $\mathcal{S}_2$ for the validation data of size $n_2$ is a simple random sample from $\mathcal{S}_1$.*

Under Assumption 1, $\{A_j, X_j, U_j, Y_j(0), Y_j(1) : j \in \mathcal{S}_2\}$ and the observations in $\mathcal{O}_2$ of the validation data are also IID, respectively. We shall relax Assumption 1 to allow $\mathcal{S}_2$ to be a general probability sample from $\mathcal{S}_1$ in Section 7. But Assumption 1 makes the presentation simpler.

**Example 1** *Two-phase sampling design is an example that results in the observed data structure. In a study, some variables (e.g. $A$, $X$, and $Y$) may be relatively cheaper, while some variables (e.g. $U$) are more expensive to obtain. A two-phase sampling design (Neyman, 1938, Cochran, 2007, Wang et al., 2009) can reduce the cost of the study: in the first phase, the easy-to-obtain variables are measured for all units, and in the second phase, additional expensive variables are measured for a selected validation sample.*

**Example 2** *Another example is highly relevant in the era of big data, where one links small data with full information on $(A, X, U, Y)$ to external big data with only $(A, X, Y)$. Chatterjee et al. (2016) recently consider this scenario for parametric regression analyses.*

Without loss of generality, we first consider the average causal effect (ACE)

$$\tau = E\{Y(1) - Y(0)\}, \tag{1}$$

and will discuss extensions to other causal estimands in Section 4.1. Because of the IID assumption, we drop the indices $i$ and $j$ in the expectations in (1) and later equations.

In what follows, we define the conditional means of the outcome as

$$\mu_a(X, U) = E(Y \mid A = a, X, U), \quad \mu_a(X) = E(Y \mid A = a, X),$$

the conditional variances of the outcome as

$$\sigma_a^2(X, U) = \mathrm{var}(Y \mid A = a, X, U), \quad \sigma_a^2(X) = \mathrm{var}(Y \mid A = a, X),$$

the conditional probabilities of the treatment as

$$e(X, U) = P(A = 1 \mid X, U), \quad e(X) = P(A = 1 \mid X).$$

## 2.2 Identification and model assumptions

A fundamental problem in causal inference is that we can observe at most one potential outcome for a unit. Following Rosenbaum and Rubin (1983b), we make the following assumptions to identify causal effects.

**Assumption 2 (Ignorability)** $Y(a) \perp\!\!\!\perp A \mid (X, U)$ *for $a = 0$ and $1$.*

4

Under Assumption 2, the treatment assignment is ignorable in $\mathcal{O}_2$ given $(X, U)$. However, the treatment assignment is only "latent" ignorable in $\mathcal{O}_1 \backslash \mathcal{O}_2$ given $X$ and the latent variable $U$ (Frangakis and Rubin, 1999, Jin and Rubin, 2008).

Moreover, we require adequate overlap between the treatment and control covariate distributions, quantified by the following assumption on the *propensity score* $e(X, U)$.

**Assumption 3 (Overlap)** *There exist constants $c_1$ and $c_2$ such that with probability 1, $0 < c_1 \leq e(X, U) \leq c_2 < 1$.*

Under Assumptions 2 and 3, $P\{A = 1 \mid X, U, Y(1)\} = P\{A = 1 \mid X, U, Y(0)\} = e(X, U)$, and $E\{Y(a) \mid X, U\} = E\{Y(a) \mid A = a, X, U\} = \mu_a(X, U)$. The ACE $\tau$ can then be estimated through regression imputation, inverse probability weighting (IPW), augmented inverse probability weighting (AIPW), or matching. See Rosenbaum (2002), Imbens (2004) and Rubin (2006) for surveys of these estimators.

In practice, the outcome distribution and the propensity score are often unknown and therefore need to be modeled and estimated.

**Assumption 4 (Outcome model)** *The parametric model $\mu_a(X, U; \beta_a)$ is a correct specification for $\mu_a(X, U)$, for $a = 0, 1$; i.e., $\mu_a(X, U) = \mu_a(X, U; \beta_a^*)$, where $\beta_a^*$ is the true model parameter, for $a = 0, 1$.*

**Assumption 5 (Propensity score model)** *The parametric model $e(X, U; \alpha)$ is a correct specification for $e(X, U)$; i.e., $e(X, U) = e(X, U; \alpha^*)$, where $\alpha^*$ is the true model parameter.*

The consistency of different estimators requires different model assumptions.

# 3   Methodology and important estimators

## 3.1   Review of commonly-used estimators based on validation data

The validation data $\{(A_j, X_j, U_j, Y_j) : j \in \mathcal{S}_2\}$ contain observations of all confounders $(X, U)$. Therefore, under Assumptions 2 and 3, $\tau$ is identifiable and can be estimated by some commonly-used estimator solely from the validation data, denoted by $\hat{\tau}_2$. Although the main data do not contain the full confounding information, we leverage the information on the common variables $(A, X, Y)$ as in the main data to improve the efficiency of $\hat{\tau}_2$. Before presenting the general theory, we first review important estimators that are widely-used in practice.

Let $\mu_a(X, U; \beta_a)$ be a working model for $\mu_a(X, U)$, for $a = 0, 1$, and $e(X, U; \alpha)$ be a working model for $e(X, U)$. We construct consistent estimators $\hat{\beta}_a$ ($a = 0, 1$) and $\hat{\alpha}$ based on $\mathcal{O}_2$, with probability limits $\beta_a^*$ ($a = 0, 1$) and $\alpha^*$, respectively. Under Assumption 4, $\mu_a(X, U; \beta_a^*) = \mu_a(X, U)$, and under Assumption 5, $e(X, U; \alpha^*) = e(X, U)$.

**Example 3 (Regression imputation)** *The regression imputation estimator is $\hat{\tau}_{\mathrm{reg},2} = n_2^{-1} \sum_{j \in \mathcal{S}_2} \hat{\tau}_{\mathrm{reg},2,j}$, where*

$$\hat{\tau}_{\mathrm{reg},2,j} = \mu_1(X_j, U_j; \hat{\beta}_1) - \mu_0(X_j, U_j; \hat{\beta}_0).$$

$\hat{\tau}_{\mathrm{reg},2}$ *is consistent for $\tau$ under Assumption 4.*

**Example 4 (Inverse probability weighting)** *The IPW estimator is $\hat{\tau}_{\mathrm{IPW},2} = n_2^{-1} \sum_{j \in \mathcal{S}_2} \hat{\tau}_{\mathrm{IPW},2,j}$, where*

$$\hat{\tau}_{\mathrm{IPW},2,j} = \frac{A_j Y_j}{e(X_j, U_j; \hat{\alpha})} - \frac{(1 - A_j) Y_j}{1 - e(X_j, U_j; \hat{\alpha})}.$$

$\hat{\tau}_{\mathrm{IPW},2}$ *is consistent for $\tau$ under Assumption 5.*

The Horvitz–Thompson-type estimator $\hat{\tau}_{\mathrm{IPW},2}$ has large variability, and is often inferior to the Hajek-type estimator (Hájek, 1971). We do not present the Hajek-type estimator because we can improve it by the AIPW estimator below. The AIPW estimator employs both the propensity score and the outcome models.

**Example 5 (Augmented inverse probability weighting)** *Define the residual outcome as $R_j = Y_j - \mu_1(X_j, U_j; \hat{\beta}_1)$ for treated units and $R_j = Y_j - \mu_0(X_j, U_j; \hat{\beta}_0)$ for control units. The AIPW estimator is $\hat{\tau}_{\mathrm{AIPW},2} = n_2^{-1} \sum_{j \in \mathcal{S}_2} \hat{\tau}_{\mathrm{AIPW},2,j}$, where*

$$\hat{\tau}_{\mathrm{AIPW},2,j} = \frac{A_j R_j}{e(X_j, U_j; \hat{\alpha})} + \mu_1(X_j, U_j; \hat{\beta}_1) - \frac{(1 - A_j) R_j}{1 - e(X_j, U_j; \hat{\alpha})} - \mu_0(X_j, U_j; \hat{\beta}_0). \quad (2)$$

$\hat{\tau}_{\mathrm{AIPW},2}$ *is doubly robust in the sense that it is consistent if either Assumption 4 or 5 holds. Moreover, it is locally efficient if both Assumptions 4 and 5 hold (Bang and Robins, 2005, Tsiatis, 2006, Cao et al., 2009).*

Matching estimators are also widely used in practice. To fix ideas, we consider matching with replacement with the number of matches fixed at $M$. Matching estimators hinge on imputing the missing potential outcome for each unit. To be precise, for unit $j$, the potential outcome under $A_j$ is the observed outcome $Y_j$; the (counterfactual) potential outcome under $1 - A_j$ is not observed but can be imputed by the average of the observed outcomes of the nearest $M$ units with $1 - A_j$. Let these matched units for unit $j$ be indexed by $\mathcal{J}_{d,V,j}$, where the subscripts $d$ and $V$ denote the dataset $\mathcal{O}_d$ and the matching variable $V$ (e.g. $V = (X, U)$), respectively. Without loss of generality, we use the Euclidean distance to determine neighbors; the discussion applies to other distances (Abadie and Imbens, 2006). Let $K_{d,V,j} = \sum_{l \in \mathcal{S}_d} 1(j \in \mathcal{J}_{d,V,l})$ be the number of times that unit $j$ is used as a match based on the matching variable $V$ in $\mathcal{O}_d$.

**Example 6 (Matching)** *Define the imputed potential outcomes as*

$$\hat{Y}_j(1) = \begin{cases} M^{-1} \sum_{l \in \mathcal{J}_{2,(X,U),j}} Y_l & \text{if } A_j = 0, \\ Y_j & \text{if } A_j = 1, \end{cases} \quad \hat{Y}_j(0) = \begin{cases} Y_j & \text{if } A_j = 0, \\ M^{-1} \sum_{l \in \mathcal{J}_{2,(X,U),j}} Y_l & \text{if } A_j = 1. \end{cases}$$

*Then the matching estimator of $\tau$ is*

$$\hat{\tau}_{\mathrm{mat},2}^{(0)} = n_2^{-1} \sum_{j \in \mathcal{S}_2} \{\hat{Y}_j(1) - \hat{Y}_j(0)\} = n_2^{-1} \sum_{j \in \mathcal{S}_2} (2A_j - 1) \left( Y_j - M^{-1} \sum_{l \in \mathcal{J}_{2,(X,U),j}} Y_l \right).$$

*Abadie and Imbens (2006) obtain the decomposition:*

$$n_2^{1/2} (\hat{\tau}_{\mathrm{mat},2}^{(0)} - \tau) = B_2 + D_2,$$

*where*

$$B_2 = n_2^{-1/2} \sum_{j \in \mathcal{S}_2} (2A_j - 1) \left[ M^{-1} \sum_{l \in \mathcal{J}_{2,(X,U),j}} \left\{ \mu_{1-A_j}(X_j, U_j) - \mu_{1-A_j}(X_l, U_l) \right\} \right], \quad (3)$$

$$D_2 = n_2^{-1/2} \sum_{j \in \mathcal{S}_2} [\mu_1(X_j, U_j) - \mu_0(X_j, U_j) - \tau$$
$$+ (2A_j - 1) \left\{ 1 + M^{-1} K_{2,(X,U),j} \right\} \left\{ Y_j - \mu_{A_j}(X_j, U_j) \right\} ].$$

*The difference $\mu_{1-A_j}(X_j, U_j) - \mu_{1-A_j}(X_l, U_l)$ in (3) accounts for the matching discrepancy, and therefore $B_2$ contributes to the asymptotic bias of the matching estimator. Abadie and Imbens (2006) show that the matching estimators have non-negligible biases with $p \geq 2$. Let $\hat{\mu}_a(X, U)$ be an estimator for $\mu_a(X, U)$, obtained either parametrically, e.g., by a linear regression estimator, or nonparametrically, for $a = 0, 1$. Abadie and Imbens (2006) propose a bias-corrected matching estimator*

$$\hat{\tau}_{\mathrm{mat},2} = \hat{\tau}_{\mathrm{mat},2}^{(0)} - n_2^{-1/2} \hat{B}_2,$$

*where $\hat{B}_2$ is an estimator for $B_2$ by replacing $\mu_a(X, U)$ with $\hat{\mu}_a(X, U)$.*

## 3.2 A general strategy

We give a general strategy for efficient estimation of the ACE by utilizing both the main and validation data. In Sections 3.3 and 3.4, we will provide examples to elucidate the proposed strategy with specific estimators.

Although the estimators based on the validation data $\mathcal{O}_2$ are consistent for $\tau$ under certain regularity conditions, they are inefficient without using the main data $\mathcal{O}_1$. However, the main data $\mathcal{O}_1$ do not contain important confounders $U$; if we naively use the estimators in Examples 3–6 with $U$ being empty, then the corresponding estimators can be inconsistent for $\tau$ and thus are error-prone in general. Moreover, for robustness consideration, we do not want to impose additional modeling assumptions linking $U$ and $(A, X, Y)$.

Our strategy is straightforward: we apply the same error-prone procedure to both the main and validation data. The key insight is that the difference of the two error-prone estimates is consistent for 0 and can be used to improve efficiency of the initial estimator due to its association with $\hat{\tau}_2$. Let an error-prone estimator of $\tau$ from the main data be $\hat{\tau}_{1,\mathrm{ep}}$, which converges to some constant $\tau_{\mathrm{ep}}$, not necessarily the same as $\tau$. Applying the same method to the validation data $\{(A_j, X_j, Y_j) : j \in \mathcal{S}_2\}$, we can obtain another error-prone estimator $\hat{\tau}_{2,\mathrm{ep}}$. More generally, we can consider $\tau_{\mathrm{ep}}$ to be an $L$-dimensional vector of parameters identifiable based on the joint distribution of $(A, X, Y)$, and $\hat{\tau}_{1,\mathrm{ep}}$ and $\hat{\tau}_{2,\mathrm{ep}}$ to be the corresponding estimators from the main and validation data, respectively. For example, $\hat{\tau}_{d,\mathrm{ep}}$ can contain estimators of $\tau$ using different methods based on $\mathcal{O}_d$.

We consider a class of estimators satisfying

$$n_2^{1/2} \begin{pmatrix} \hat{\tau}_2 - \tau \\ \hat{\tau}_{2,\mathrm{ep}} - \hat{\tau}_{1,\mathrm{ep}} \end{pmatrix} \to \mathcal{N} \left\{ 0_{L+1}, \begin{pmatrix} v_2 & \Gamma^{\mathrm{T}} \\ \Gamma & V \end{pmatrix} \right\}, \quad (4)$$

in distribution, as $n_2 \to \infty$, which is general enough to include all the estimators reviewed in Examples 3–6. Heuristically, if (4) holds exactly rather than asymptotically, by the multivariate normal theory, we have the following the conditional distribution

$$n_2^{1/2}(\hat{\tau}_2 - \tau) \mid n_2^{1/2}(\hat{\tau}_{2,\mathrm{ep}} - \hat{\tau}_{1,\mathrm{ep}}) \sim \mathcal{N} \left\{ n_2^{1/2} \Gamma^{\mathrm{T}} V^{-1}(\hat{\tau}_{2,\mathrm{ep}} - \hat{\tau}_{1,\mathrm{ep}}), v_2 - \Gamma^{\mathrm{T}} V^{-1} \Gamma \right\}.$$

Let $\hat{v}_2$, $\hat{\Gamma}$ and $\hat{V}$ be consistent estimators for $v_2$, $\Gamma$ and $V$. We set $n_2^{1/2}(\hat{\tau}_2 - \tau)$ to equal its estimated conditional mean $n_2^{1/2}\hat{\Gamma}^{\mathrm{T}}\hat{V}^{-1}(\hat{\tau}_{2,\mathrm{ep}} - \hat{\tau}_{1,\mathrm{ep}})$, leading to an estimating equation for $\tau$:

$$n_2^{1/2}(\hat{\tau}_2 - \tau) = n_2^{1/2}\hat{\Gamma}^{\mathrm{T}}\hat{V}^{-1}(\hat{\tau}_{2,\mathrm{ep}} - \hat{\tau}_{1,\mathrm{ep}}).$$

Solving this equation for $\tau$, we obtain the estimator

$$\hat{\tau} = \hat{\tau}_2 - \hat{\Gamma}^{\mathrm{T}}\hat{V}^{-1}(\hat{\tau}_{2,\mathrm{ep}} - \hat{\tau}_{1,\mathrm{ep}}). \tag{5}$$

**Proposition 1** *Under Assumption 1 and certain regularity conditions, if (4) holds, then $\hat{\tau}$ is consistent for $\tau$, and*

$$n_2^{1/2}(\hat{\tau} - \tau) \to \mathcal{N}(0, v_2 - \Gamma^{\mathrm{T}}V^{-1}\Gamma), \tag{6}$$

*in distribution, as $n_2 \to \infty$. Given a nonzero $\Gamma$, the asymptotic variance, $v_2 - \Gamma^{\mathrm{T}}V^{-1}\Gamma$, is smaller than the asymptotic variance of $\hat{\tau}_2$, $v_2$.*

The consistency of $\hat{\tau}$ does not require any component in $\hat{\tau}_{1,\mathrm{ep}}$ and $\hat{\tau}_{2,\mathrm{ep}}$ to correctly estimate $\tau$. That is, these estimators can be error prone. The requirement for the error-prone estimators is minimal, as long as they are consistent for the same (finite) parameters. Under Assumption 1, $\hat{\tau}_{1,\mathrm{ep}} - \hat{\tau}_{2,\mathrm{ep}}$ is consistent for a vector of zeros, as $n_2 \to \infty$.

We can estimate the asymptotic variance of $\hat{\tau}$ by

$$\hat{v} = (\hat{v}_2 - \hat{\Gamma}^{\mathrm{T}}\hat{V}^{-1}\hat{\Gamma})/n_2. \tag{7}$$

**Remark 1** *We construct the error prone estimators $\hat{\tau}_{1,\mathrm{ep}}$ and $\hat{\tau}_{2,\mathrm{ep}}$ based on $\mathcal{O}_1$ and $\mathcal{O}_2$, respectively. Another intuitive way is to construct $\hat{\tau}_{1,\mathrm{ep}}$ and $\hat{\tau}_{2,\mathrm{ep}}$ based on $\mathcal{O}_1 \backslash \mathcal{O}_2$ and $\mathcal{O}_2$, respectively. In general, we can construct the error prone estimators based on different subsets of $\mathcal{O}_1$ and $\mathcal{O}_2$ as long as their difference converges in probability to zero. We show in the supplementary material that our construction maximizes the variance reduction for $\hat{\tau}_2$, $\Gamma^{\mathrm{T}}V^{-1}\Gamma$, given the procedure of the error prone estimators.*

**Remark 2** *We can view (5) as the best consistent estimator of $\tau$ among all linear combinations $\{\hat{\tau}_2 + \lambda^{\mathrm{T}}(\hat{\tau}_{2,\mathrm{ep}} - \hat{\tau}_{1,\mathrm{ep}}) : \lambda \in \mathbb{R}^L\}$, in the sense that (5) achieves the minimal asymptotic variance among this class of consistent estimators. Similar ideas appeared in design-optimal regression estimation in survey sampling (Deville and Särndal, 1992, Fuller, 2009), regression analyses (Chen and Chen, 2000; Chen, 2002 and Wang and Wang, 2015), improved prediction in high dimensional datasets (Boonstra et al., 2012), and meta-analysis (Collaboration, 2009). In the supplementary material, we show that the proposed estimator in (5) is the best estimator of $\tau$ among the class of estimators $\{\hat{\tau} = f(\hat{\tau}_2, \hat{\tau}_{1,\mathrm{ep}}, \hat{\tau}_{2,\mathrm{ep}}): f(x, y, z)$ is a smooth function of $(x, y, z)$, and $\hat{\tau}$ is consistent for $\tau\}$, in the sense that (5) achieves the minimal asymptotic variance among this class.*

**Remark 3** *The choice of the error-prone estimators will affect the efficiency of $\hat{\tau}$. From (6), for a given $\hat{\tau}_2$, to improve the efficiency of $\hat{\tau}$ with a 1-dimensional error-prone estimator, we would like this estimator to have a small variance $V$ and a large correlation with $\hat{\tau}_2$, $\Gamma$. In principle, increasing the dimension of the error-prone estimator would not decrease the asymptotic efficiency gain as shown in the supplementary material. However, it would also increase the complexity of implementation and harm the finite sample properties. To "optimize" the trade-off, we suggest choosing the error-prone estimator to be the same type as the initial estimator $\hat{\tau}_2$. For example, if $\hat{\tau}_2$ is an AIPW estimator, we can choose $\hat{\tau}_{d,\mathrm{ep}}$ to be an AIPW estimator without using U in a possibly misspecified propensity score model. The simulation in Section 5 confirms that this choice is reasonable.*

To close this subsection, we comment on the existing literature and the advantages of our strategy. The proposed estimator $\hat{\tau}$ in (5) utilizes both the main and validation data and improves the efficiency of the estimator based solely on the validation data. In economics, Imbens and Lancaster (1994) propose to use the generalized method of moments (Hansen, 1982) for utilizing the main data which provide moments of the marginal distribution of some economic variables. In survey sampling, calibration is a standard technique to integrate auxiliary information in estimation or handle nonresponse; see, e.g., Chen and Chen (2000), Wu and Sitter (2001), Kott (2006), Chang and Kott (2008) and Kim et al. (2016). An important issue is how to specify optimal calibration equations; see, for example, Deville and Särndal (1992), Robins et al. (1994), Wu and Sitter (2001), and Lumley et al. (2011). Other researchers developed constrained empirical likelihood methods to calibrate auxiliary information from the main data; see, e.g., Chen and Sitter (1999), Qin (2000), Chen et al. (2002) and Chen et al. (2003).

Compared to these methods, the proposed framework is attractive because it is simple to implement which requires only standard software routines for existing methods, and it can deal with estimators that cannot be derived from moment conditions, e.g., matching estimators. Moreover, our framework does not require a correct model specification of the relationship between unmeasured covariates $U$ and measured variables $(A, X, Y)$.

## 3.3 Regular asymptotically linear (RAL) estimators

We first elucidate the proposed method with RAL estimators.

From the validation data, we consider the case when $\hat{\tau}_2 - \tau$ is RAL; i.e., it can be asymptotically approximated by a sum of IID random vectors with mean 0:

$$\hat{\tau}_2 - \tau \cong n_2^{-1} \sum_{j \in \mathcal{S}_2} \psi(A_j, X_j, U_j, Y_j), \tag{8}$$

where $\{\psi(A_j, X_j, U_j, Y_j) : j \in \mathcal{S}_2\}$ are IID with mean 0. The random vector $\psi(A, X, U, Y)$ is called the influence function of $\hat{\tau}_2$ with $E(\psi) = 0$ and $E(\psi^2) < \infty$ (e.g. Bickel et al., 1993). Regarding regularity conditions, see, for example, Newey (1990).

Let $e(X; \gamma)$ be an error-prone propensity score model for $e(X)$, and $\mu_a(X; \eta_a)$ be an error-prone outcome regression model for $\mu_a(X)$, for $a = 0, 1$. The corresponding error-prone estimators of the ACE can be obtained from the main data $\mathcal{O}_1$ and the validation data $\mathcal{O}_2$. We consider the case when $\hat{\tau}_{d,\text{ep}}$ is RAL:

$$\hat{\tau}_{d,\text{ep}} - \tau_{\text{ep}} \cong n_d^{-1} \sum_{j \in \mathcal{S}_d} \phi(A_j, X_j, Y_j), \qquad (d = 1, 2), \tag{9}$$

where $\{\phi(A_j, X_j, Y_j) : j \in \mathcal{S}_d\}$ are IID with mean 0.

**Theorem 1** *Under certain regularity conditions, (4) holds for the RAL estimators (8) and (9), where $v_2 = \text{var}\{\psi(A, X, U, Y)\}$, $\Gamma = (1 - \rho)\text{cov}\{\psi(A, X, U, Y), \phi(A, X, Y)\}$, and $V = (1 - \rho) \times \text{var}\{\phi(A, X, Y)\}$.*

To derive $\hat{\Gamma}$ and $\hat{V}$ for RAL estimators, let $\hat{\phi}_d(A, X, Y)$ and $\hat{\psi}(A, X, U, Y)$ be estimators of $\phi(A, X, Y)$ and $\psi(A, X, U, Y)$ by replacing $E(\cdot)$ with the empirical measure and unknown parameters with their corresponding estimators. Note that the subscript $d$ in $\hat{\phi}_d(A, X, Y)$ indicates that it

is obtained based on $\mathcal{O}_d$. Then, we can estimate $\Gamma$ and $V$ by

$$
\hat{\Gamma} = \widehat{\text{cov}}(\hat{\tau}_2, \hat{\tau}_{2,\text{ep}} - \hat{\tau}_{1,\text{ep}}) = \left(1 - \frac{n_2}{n_1}\right) \frac{1}{n_2} \sum_{j \in S_2} \hat{\psi}(A_j, X_j, U_j, Y_j)\hat{\phi}_2(A_j, X_j, Y_j),
$$

$$
\hat{V} = \widehat{\text{var}}(\hat{\tau}_{1,\text{ep}} - \hat{\tau}_{2,\text{ep}}) = \left(1 - \frac{n_2}{n_1}\right) \frac{1}{n_1} \sum_{i \in S_1} \left\{\hat{\phi}_1(A_i, X_i, Y_i)\right\}^{\otimes 2}.
$$

Finally, we can obtain the estimator and its variance estimator by (5) and (7), respectively.

The commonly-used RAL estimators include the regression imputation and (augmented) inverse probability weighting estimators. Because the influence functions for $\hat{\tau}_{\text{reg},2}$ and $\hat{\tau}_{\text{IPW},2}$ are standard, we present the details in the supplementary material. Below, we state only the influence function for $\hat{\tau}_{\text{AIPW},2}$.

For the outcome model, let $S_a(A, X, U, Y; \beta_a)$ be the estimating function for $\beta_a$, e.g.,

$$
S_a(A, X, U, Y; \beta_a) = \frac{\partial \mu_a(X, U; \beta_a)}{\partial \beta_a}\{Y - \mu_a(X, U; \beta_a)\},
$$

for $a = 0, 1$, which is a standard choice for the conditional mean model. For the propensity score model, let $S(A, X, U; \alpha)$ be the estimating function for $\alpha$, e.g.,

$$
S(A, X, U; \alpha) = \frac{A - e(X, U; \alpha)}{e(X, U; \alpha)\{1 - e(X, U; \alpha)\}} \frac{\partial e(X, U; \alpha)}{\partial \alpha},
$$

which is the score function from the likelihood of a binary response model. Moreover, let

$$
\Sigma_{\alpha\alpha} = E\left\{S^{\otimes 2}(A, X, U; \alpha)\right\} = E\left[\frac{1}{e(X, U; \alpha^*)\{1 - e(X, U; \alpha^*)\}}\left\{\frac{\partial e(X, U; \alpha^*)}{\partial \alpha}\right\}^{\otimes 2}\right]
$$

be the Fisher information matrix for $\alpha$ in the propensity score model. In addition, let $\hat{\beta}_a$ $(a = 0, 1)$ and $\hat{\alpha}$ be the estimators solving the corresponding empirical estimating equations based on $\mathcal{O}_2$, with probability limits $\beta_a^*$ $(a = 0, 1)$ and $\alpha^*$, respectively.

**Lemma 1 (Augmented inverse probability weighting)** *For simplicity, denote* $e_j^* = e(X_j, U_j; \alpha^*)$, $\dot{e}_j^* = \partial e(X_j, U_j; \alpha^*)/\partial\alpha^\mathrm{T}$, $S_j^* = S(A_j, X_j, U_j; \alpha^*)$, $\mu_{aj}^* = \mu_a(X_j, U_j; \beta_a^*)$, $\dot{\mu}_{aj}^* = \partial\mu_a(X_j, U_j; \beta_a^*)/\partial\beta_a^\mathrm{T}$, $S_{aj}^* = S_a(A_j, X_j, U_j, Y_j; \beta_a^*)$, *and* $\dot{S}_{aj}^* = \partial S_a(A_j, X_j, U_j, Y_j; \beta_a^*)/\partial\beta_a^\mathrm{T}$ *for* $a = 0, 1$. *Under Assumption 4 or 5,* $\hat{\tau}_{\text{AIPW},2}$ *has the influence function*

$$
\psi_{\text{AIPW}}(A_j, X_j, U_j, Y_j) = \frac{A_j Y_j}{e_j^*} + \left(1 - \frac{A_j}{e_j^*}\right)\mu_{1j}^*
$$

$$
- \frac{(1 - A_j)Y_j}{1 - e_j^*} - \left(1 - \frac{1 - A_j}{1 - e_j^*}\right)\mu_{0j}^* - \tau + H_{\text{AIPW}}\Sigma_{\alpha\alpha}^{-1}S_j^*
$$

$$
+ E\left\{\left(1 - \frac{1 - A}{1 - e^*}\right)\dot{\mu}_0^*\right\}\left\{E(\dot{S}_0^*)\right\}^{-1}S_{0j}^* \tag{10}
$$

$$
- E\left\{\left(1 - \frac{A}{e^*}\right)\dot{\mu}_1^*\right\}\left\{E(\dot{S}_1^*)\right\}^{-1}S_{1j}^*, \tag{11}
$$

*where*

$$
H_{\text{AIPW}} = E\left[\left\{\frac{A(Y - \mu_1^*)}{(e^*)^2} - \frac{(1 - A)(Y - \mu_0^*)}{(1 - e^*)^2}\right\}\dot{e}^*\right].
$$

10

Lemma 1 follows from standard asymptotic theory, but as far as we know it has not appeared in the literature. Lunceford and Davidian (2004) suggest using the terms without (10) and (11) for $\psi_{\mathrm{AIPW}}$, which, however, works only when both Assumptions 4 and 5 hold. Otherwise, the resulting variance estimator is not consistent if either Assumption 4 or 5 does not hold, as shown by simulation in Funk et al. (2011). The correction terms in (10) and (11) also make the variance estimator doubly robust in the sense that the variance estimator for $\hat{\tau}_{\mathrm{AIPW},2}$ is consistent if either Assumption 4 or 5 holds, not necessarily both.

For error-prone estimators, we can obtain the influence functions similarly. The subtlety is that both the propensity score and outcome models can be misspecified. For simplicity of the presentation, we defer the exact formulas to the online supplementary material.

### 3.4 Matching estimators

We then elucidate the proposed method with non-RAL estimators. An important class of non-RAL estimators for the ACE are the matching estimators. The matching estimators are not regular estimators because the functional forms are not smooth due to the fixed numbers of matches (Abadie and Imbens, 2008). Continuing with Example 6, Abadie and Imbens (2006) express the bias-corrected matching estimator $\hat{\tau}_{\mathrm{mat},2}$ in a linear form as

$$\hat{\tau}_{\mathrm{mat},2} - \tau \cong n_2^{-1} \sum_{j \in \mathcal{S}_2} \psi_{\mathrm{mat},j}, \tag{12}$$

where

$$\psi_{\mathrm{mat},j} = \mu_1(X_j, U_j) - \mu_0(X_j, U_j) - \tau + (2A_j - 1)\left\{1 + M^{-1}K_{2,(X,U),j}\right\}\left\{Y_j - \mu_{A_j}(X_j, U_j)\right\}. \tag{13}$$

Similarly, $\hat{\tau}_{\mathrm{mat},d,\mathrm{ep}}$ has a linear form

$$\hat{\tau}_{\mathrm{mat},d,\mathrm{ep}} - \tau_{\mathrm{ep}} \cong n_d^{-1} \sum_{j \in \mathcal{S}_d} \phi_{\mathrm{mat},d,j}, \tag{14}$$

where

$$\phi_{\mathrm{mat},d,j} = \mu_1(X_j) - \mu_0(X_j) - \tau_{\mathrm{ep}} + (2A_j - 1)\left(1 + M^{-1}K_{d,X,j}\right)\left\{Y_j - \mu_{A_j}(X_j)\right\}. \tag{15}$$

**Theorem 2** *Under certain regularity conditions, (4) holds for the matching estimators (12) and (14), where*

$$
\begin{aligned}
v_2 &= \operatorname{var}\left\{\tau(X,U)\right\} + \operatorname{plim}\left[n_2^{-1}\sum_{j \in \mathcal{S}_2}\left\{1 + M^{-1}K_{2,(X,U),j}\right\}^2 \sigma_{A_j}^2(X_j, U_j)\right], \\
\Gamma &= (1-\rho)\left(\operatorname{cov}\left\{\mu_1(X,U) - \mu_0(X,U), \mu_1(X) - \mu_0(X)\right\}\right. \\
&\qquad \left. + \operatorname{plim}\left[n_2^{-1}\sum_{j \in \mathcal{S}_2}\left\{1 + M^{-1}K_{2,(X,U),j}\right\}\left(1 + M^{-1}K_{2,X,j}\right)\sigma_{A_j}^2(X_j, U_j)\right]\right), \\
V &= (1-\rho)\left[\operatorname{var}\left\{\mu_1(X) - \mu_0(X)\right\} + \operatorname{plim}\left\{n_2^{-1}\sum_{j \in S_2}\left(1 + M^{-1}K_{2,X,j}\right)^2 \sigma_{A_j}^2(X_j)\right\}\right].
\end{aligned}
$$

11

The existence of the probability limits in Theorem 2 are guaranteed by the regularity conditions specified in the supplementary material (c.f. Abadie and Imbens, 2006).

To estimate $(v_2, \Gamma, V)$ in Theorem 2, we need to estimate the conditional mean and variance functions of the outcome given covariates. Following Abadie and Imbens (2006), we can estimate these functions via matching units with the same treatment level. We will discuss an alternative bootstrap strategy in the next subsection.

### 3.5 Bootstrap variance estimation

The asymptotic results in Theorems 1 and 2 allow for variance estimation of $\hat{\tau}$. In addition, we also consider the bootstrap for variance estimation, which is simpler to implement and often has better finite sample performances (Otsu and Rai, 2016). This is particularly important for matching estimators because the analytic variance formulas involve nonparametric estimation of the conditional variances $\sigma_a^2(x, u)$ and $\sigma_a^2(x)$.

There are two approaches for obtaining bootstrap observations: (a) the original observations; and (b) the asymptotic linear terms of the proposed estimator. For RAL estimators, bootstrapping the original observations will yield valid variance estimators (Efron and Tibshirani, 1986, Shao and Tu, 2012). However, for matching estimators, Abadie and Imbens (2008) show that due to lack of smoothness in their functional form, the bootstrap based on approach (a) does not apply for variance estimation. This is mainly because the bootstrap based on approach (a) cannot preserve the distribution of the numbers of times that the units are used as matches. As a remedy, Otsu and Rai (2016) propose to construct the bootstrap counterparts by resampling based on approach (b) for the matching estimator.

To unify the notation, let $\psi_j$ indicate $\psi(A_j, X_j, U_j, Y_j)$ for RAL $\hat{\tau}_2$ and $\psi_{\mathrm{mat},j}$ for $\hat{\tau}_{\mathrm{mat},2}$; and similar definitions apply to $\phi_{d,j}$ $(d = 1, 2)$. Let $\hat{\psi}_j$ and $\hat{\phi}_{d,j}$ be their estimated version by replacing the population quantities by the estimated quantities $(d = 1, 2)$. Following Otsu and Rai (2016), for $b = 1, \ldots, B$, we construct the bootstrap replicates for the proposed estimators as follows:

**Step** 1. Sample $n_1$ units from $\mathcal{S}_1$ with replacement as $\mathcal{S}_1^{*(b)}$, treating the units with observed $U$ as the bootstrap validation data $\mathcal{S}_2^{*(b)}$.

**Step** 2. Compute the bootstrap replicates of $\hat{\tau}_2 - \tau$ and $\hat{\tau}_{d,\mathrm{ep}} - \tau_{\mathrm{ep}}$ as

$$
\hat{\tau}_2^{(b)} - \hat{\tau}_2 = n_2^{-1} \sum_{j \in \mathcal{S}_2^{*(b)}} \hat{\psi}_j,
$$

$$
\hat{\tau}_{d,\mathrm{ep}}^{(b)} - \hat{\tau}_{d,\mathrm{ep}} = n_d^{-1} \sum_{j \in \mathcal{S}_d^{*(b)}} \hat{\phi}_{d,j}, \qquad (d = 1, 2).
$$

Based on the bootstrap replicates, we estimate $\Gamma, V$ and $v_2$ by

$$
\hat{\Gamma} = (B-1)^{-1} \sum_{b=1}^{B} (\hat{\tau}_2^{(b)} - \hat{\tau}_2)(\hat{\tau}_{2,\mathrm{ep}}^{(b)} - \hat{\tau}_{1,\mathrm{ep}}^{(b)} - \hat{\tau}_{2,\mathrm{ep}} + \hat{\tau}_{1,\mathrm{ep}}), \tag{16}
$$

$$
\hat{V} = (B-1)^{-1} \sum_{b=1}^{B} (\hat{\tau}_{2,\mathrm{ep}}^{(b)} - \hat{\tau}_{1,\mathrm{ep}}^{(b)} - \hat{\tau}_{2,\mathrm{ep}} + \hat{\tau}_{1,\mathrm{ep}})^{\otimes 2}, \tag{17}
$$

$$
\hat{v}_2 = (B-1)^{-1} \sum_{b=1}^{B} (\hat{\tau}_2^{(b)} - \hat{\tau}_2)^2. \tag{18}
$$

Finally, we estimate the asymptotic variance of $\hat{\tau}$ by (7), i.e., $\hat{v} = (\hat{v}_2 - \hat{\Gamma}^{\mathrm{T}} \hat{V}^{-1} \hat{\Gamma})/n_2$.

**Theorem 3** *Under certain regularity conditions, $(\hat{\Gamma}, \hat{V}, \hat{v}_2, \hat{v})$ are consistent for $\{\Gamma, V, \mathrm{var}(\hat{\tau}_2), \mathrm{var}(\hat{\tau})\}$.*

**Remark 4** *If the ratio of $n_2$ and $n_1$ is small, the above bootstrap approach may be unstable, because it is likely that some bootstrap validation data contain only a few or even zero observations. In this case, we use an alternative bootstrap approach, where we sample $n_2$ units from $\mathcal{S}_2$ with replacement as $\mathcal{S}_2^*$, sample $n_1 - n_2$ units from $\mathcal{S}_1 \backslash \mathcal{S}_2$ with replacement, combined with $\mathcal{S}_2^*$, as $\mathcal{S}_1^*$, and obtain the proposed estimators based on $\mathcal{S}_1^*$ and $\mathcal{S}_2^*$. This approach guarantees that the bootstrap validation data contain $n_2$ observations.*

**Remark 5** *It is worthwhile to comment on a computational issue. When the main data have a substantially large size, the computation for the bootstrap can be demanding if we follow Steps 1 and 2 above. In this case, we can use subsampling (Politis et al., 1999) or the Bag of Little Bootstrap (Kleiner et al., 2014) to reduce the computational burden. More interestingly, when $n_1 \to \infty$ and $\rho = 0$, i.e., the validation data contain a small fraction of the main data, $\Gamma$ and $V$ reduce to $\mathrm{cov}(\hat{\tau}_2, \hat{\tau}_{2,\mathrm{ep}})$ and $\mathrm{var}(\hat{\tau}_{2,\mathrm{ep}})$, respectively. That is, when the size of the main data is substantially large, we can ignore the uncertainty of $\hat{\tau}_{1,\mathrm{ep}}$ and treat it as a constant, which is a regime recently considered by Chatterjee et al. (2016). In this case, we need only to bootstrap the validation data, which is computationally simpler.*

## 4 Extensions

### 4.1 Other causal estimands

Our strategy extends to a wide class of causal estimands, as long as (4) holds. For example, we can consider the average causal effects over a subset of population (Crump et al., 2006, Li et al., 2016), including the average causal effect on the treated.

We can also consider nonlinear causal estimands. For example, for a binary outcome, the log of the causal risk ratio is

$$\log \mathrm{CRR} = \log \frac{P\{Y(1) = 1\}}{P\{Y(0) = 1\}} = \log \frac{E\{Y(1)\}}{E\{Y(0)\}},$$

and the log of the causal odds ratio is

$$\log \mathrm{COR} = \log \frac{P\{Y(1) = 1\}/P\{Y(1) = 0\}}{P\{Y(0) = 1\}/P\{Y(0) = 0\}} = \log \frac{E\{Y(1)\}/[1 - E\{Y(1)\}]}{E\{Y(0)\}/[1 - E\{Y(0)\}]}.$$

We give a brief discussion for the $\log \mathrm{CRR}$ as an illustration. The key insight is that under Assumptions 2 and 3, we can estimate $E\{Y(a)\}$ with commonly-used estimators from $\mathcal{O}_2$, denoted by $\hat{E}\{Y(a)\}$, for $a = 0, 1$. We can then obtain an estimator for the $\log \mathrm{CRR}$ as $\log[\hat{E}\{Y(1)\}/\hat{E}\{Y(0)\}]$. Similarly, we can obtain error-prone estimators for the $\log \mathrm{CRR}$ from both $\mathcal{O}_1$ and $\mathcal{O}_2$ using only covariates $X$. By the Taylor expansion, we can linearize these estimators and establish a similar result as (4), which serves as the basis to construct an improved estimator for the $\log \mathrm{CRR}$.

## 4.2 Design issue: optimal sample size allocation

As a design issue, we consider planning a study to obtain the data structure in Section 2 subject to a cost constraint. The goal is to find the optimal design, specifically the sample allocation, that minimizes the variance of the proposed estimator subject to a cost constraint, as in the classical two-phase sampling (Cochran, 2007).

Suppose that it costs $C_1$ to collect $(A, X, Y)$ for each unit, and $C_2$ to collect $U$ for each unit. Thus, the total cost of the study is

$$C = n_1 C_1 + n_2 C_2. \tag{19}$$

The variance of the proposed estimator $\hat{\tau}$ is of the form

$$n_2^{-1} v_2 - (n_2^{-1} - n_1^{-1})\gamma, \tag{20}$$

e.g., for RAL estimators,

$$\gamma = \text{cov}\{\psi(A, X, U, Y), \phi(A, X, Y)\}^{\text{T}} [\text{var}\{\phi(A, X, Y)\}]^{-1} \text{cov}\{\psi(A, X, U, Y), \phi(A, X, Y)\}$$

is the variance of the projection of $\psi(A, X, U, Y)$ onto the linear space spanned by $\phi(A, X, Y)$. Minimizing (20) with respect to $n_1$ and $n_2$ subject to the constraint (19) yields the optimal $n_1^*$ and $n_2^*$, which satisfy

$$\rho^* = \frac{n_2^*}{n_1^*} = \left\{ (1 - R_{\psi|\phi}^2) \times \frac{C_1}{C_2} \right\}^{1/2}, \tag{21}$$

where $R_{\psi|\phi}^2 = \gamma/v_2$ is the squared multiple correlation coefficient of $\psi(A, X, U, Y)$ on $\phi(A, X, Y)$, which measures the association between the initial estimator and the error-prone estimator. We derive (21) by the Lagrange multipliers, and relegate the details to the supplementary material. Not surprisingly, (21) shows that the sizes of the validation data and the main data should be inversely proportional to the square-root of the costs. In addition, from (21), a large size $n_2$ for the validation data is more desirable when the association between the initial estimator and the error-prone estimator is small.

## 4.3 Multiple data sources

We have considered the setting with two data sources, and we can easily extend the theory to the setting with multiple data sources $\mathcal{O}_1, ..., \mathcal{O}_K$, where $\mathcal{O}_1, ..., \mathcal{O}_{K-1}$ contain partial covariate information, and the validation data, $\mathcal{O}_K$, contain full information for $(A, X, U, Y)$. For example, for $d = 1, \ldots, K - 1$, $\mathcal{O}_d$ contains variables $(A, V_d, Y)$ where $V_d \subsetneq (X, U)$. Each dataset $\mathcal{O}_d$, indexed by $\mathcal{S}_d$, has size $n_d$ for $d = 1, \ldots, K$. This type of data structure arises from a multi-phase sampling as an extension of Example 1 or multiple sources of "big data" as an extension of Example 2.

Let $\hat{\tau}_K$ be the initial estimator for $\tau$ from the validation data $\mathcal{O}_K$, and $\hat{\tau}_{d,\text{ep}}$ be the error-prone estimator for $\tau$ from $\mathcal{O}_d$ ($d = 1, \ldots, K - 1$). Let $\hat{\tau}_{d,K,\text{ep}}$ be the estimator obtained by applying the same error-prone estimator for $O_d$ to $\mathcal{O}_K$, so that $\hat{\tau}_{d,\text{ep}} - \hat{\tau}_{d,K,\text{ep}}$ is consistent for 0, for $d = 1, \ldots, K - 1$. Assume that

$$n_K^{1/2} \begin{pmatrix} \hat{\tau}_K - \tau \\ \hat{\tau}_{1,\text{ep}} - \hat{\tau}_{1,K,\text{ep}} \\ \vdots \\ \hat{\tau}_{K-1,\text{ep}} - \hat{\tau}_{K-1,K,\text{ep}} \end{pmatrix} \to \mathcal{N} \left\{ \begin{pmatrix} 0 \\ 0_L \end{pmatrix}, \begin{pmatrix} v_K & \Gamma^{\text{T}} \\ \Gamma & V \end{pmatrix} \right\},$$

14

in distribution, as $n_K \to \infty$, where $L = \sum_{d=1}^{K-1} \dim(\hat{\tau}_{d,\mathrm{ep}})$. If $\Gamma$ and $V$ have consistent estimators $\hat{\Gamma}$ and $\hat{V}$, respectively, then, extending the proposed method in Section 3, we can use

$$\hat{\tau} = \hat{\tau}_K - \hat{\Gamma}^{\mathrm{T}}\hat{V}^{-1} \left( \begin{array}{ccc} \hat{\tau}_{1,\mathrm{ep}}^{\mathrm{T}} - \hat{\tau}_{1,K,\mathrm{ep}}^{\mathrm{T}}, & \cdots, & \hat{\tau}_{K-1,\mathrm{ep}}^{\mathrm{T}} - \hat{\tau}_{K-1,K,\mathrm{ep}}^{\mathrm{T}} \end{array} \right)^{\mathrm{T}}$$

to estimate $\tau$. The estimator $\hat{\tau}$ is consistent for $\tau$ with the asymptotic variance $v_K - \Gamma^{\mathrm{T}}V^{-1}\Gamma$, which is smaller than the asymptotic variance of $\hat{\tau}_K$, $v_K$, if $\Gamma$ is non-zero. Similar to the reasoning in Remark 3, using more data sources will improve the asymptotic estimation efficiency of $\tau$.

# 5 Simulation

In this section, we conduct a simulation study to evaluate the finite sample performance of the proposed estimators. In our data generating model, the covariates are $X_i \sim \mathrm{Unif}(0,2)$ and $U_i = 0.5 + 0.5X_i - 2\sin(X_i) + 2\mathrm{sign}\{\sin(5X_i)\} + \epsilon_i$, where $\epsilon_i \sim \mathrm{Unif}(-0.5, 0.5)$. The potential outcomes are $Y_i(0) = -X_i - U_i + \epsilon_i(0)$ and $Y_i(1) = -X_i + 4U_i + \epsilon_i(1)$, where $\epsilon_i(0) \sim \mathcal{N}(0,1)$, $\epsilon_i(1) \sim \mathcal{N}(0,1)$, and $\epsilon_i(0)$ and $\epsilon_i(1)$ are independent. Therefore, the true value of the ACE is $\tau = E(5U_i)$. The treatment indicator $A_i$ follows $\mathrm{Bernoulli}(\pi_i)$ with $\mathrm{logit}(\pi_i) = 1 - 0.5X_i - 0.5U_i$. The main data $\mathcal{O}_1$ consist of $n_1$ units, and the validation data $\mathcal{O}_2$ consist of $n_2$ units randomly selected from the main data.

The initial estimators are the regression imputation, (A)IPW and matching estimators applied solely to the validation data, denoted by $\hat{\tau}_{\mathrm{reg},2}$, $\hat{\tau}_{\mathrm{IPW},2}$, $\hat{\tau}_{\mathrm{AIPW},2}$ and $\hat{\tau}_{\mathrm{mat},2}$, respectively. To distinguish the estimators constructed based on different error-prone methods, we assign each proposed estimator a name with the form $\hat{\tau}_{\mathrm{method},2\&\mathrm{methods}}$, where "method,2" indicates the initial estimator applied to the validation data $\mathcal{O}_2$, and "methods" indicates the error-prone estimator(s) used to improve the efficiency of the initial estimator. For example, $\hat{\tau}_{\mathrm{reg},2\&\mathrm{IPW}}$ indicates the initial estimator is the regression imputation estimator and the error-prone estimator is the IPW estimator. We compare the proposed estimators with the initial estimators in terms of percentages of reduction of mean squared errors, defined as $\{1 - \mathrm{MSE}(\hat{\tau}_{\mathrm{method},2\&\mathrm{methods}})/\mathrm{MSE}(\hat{\tau}_{\mathrm{method},2})\} \times 100\%$. To demonstrate the robustness of the proposed estimator against misspecification of the imputation model, we consider the multiple imputation (MI, Rubin, 1987) estimator, denoted by $\hat{\tau}_{\mathrm{mi}}$, which uses a regression model of $U$ given $(A, X, Y)$ for imputation. We implement MI using the "mice" package in R with $m = 10$.

Based on a point estimate $\hat{\tau}$ and a variance estimate $\hat{v}$ obtained by the asymptotic variance formula or the bootstrap method described in Section 3.5, we construct a Wald-type $95\%$ confidence interval $(\hat{\tau} - z_{0.975}\hat{v}^{1/2}, \hat{\tau} + z_{0.975}\hat{v}^{1/2})$, where $z_{0.975}$ is the $97.5\%$ quantile of the standard normal distribution. We further compare the variance estimators in terms of empirical coverage rates.

Figure 1 shows the simulation results over $2,000$ Monte Carlo samples for $(n_1, n_2) = (1000, 200)$ and $(n_1, n_2) = (1000, 500)$. The multiple imputation estimator is biased due to the misspecification of the imputation model. In all scenarios, the proposed estimators are unbiased and improve the initial estimators. Using the error-prone estimator of the same type of the initial estimator achieves a substantial efficiency gain, and the efficiency gain from incorporating additional error-prone estimator is not significantly important. Because of the practical simplicity, we recommend using the same type of error-prone estimator to improve the efficiency of the initial estimator. Confidence intervals constructed from the asymptotic variance formula and the bootstrap method work well, in the sense that the empirical coverage rate of the confidence intervals is close to the nominal coverage rate. In our settings, the matching estimator has the smallest efficiency gain among all types of estimators.
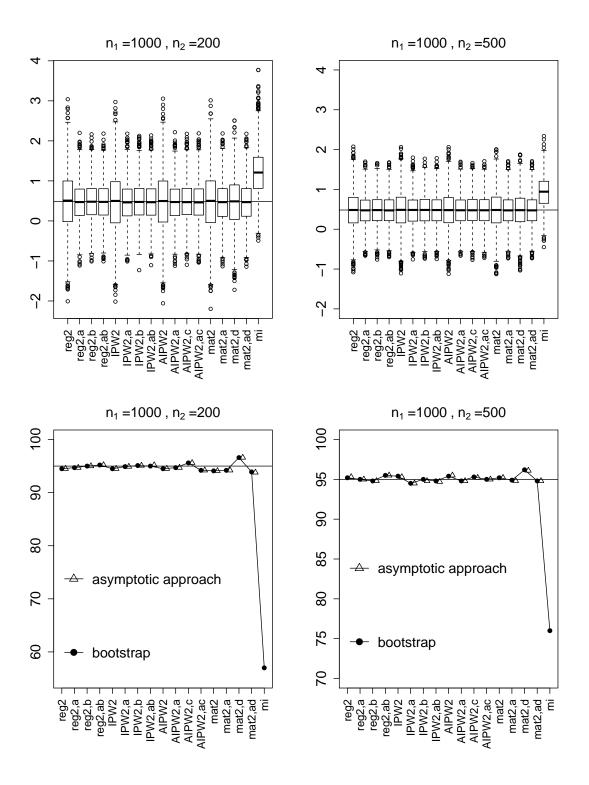
Figure 1: Simulation results of point estimates (top panels) and coverage rates (bottom panels): the subscripts "a", "b", "c", "d" stand for methods "reg", "IPW", "AIPW", "mat", respectively, "reg2" is $\hat{\tau}_{\text{reg},2}$, "reg2,method" is $\hat{\tau}_{\text{reg},2\&\text{method}}$, other notation is defined similarly, and "mi" is $\hat{\tau}_{\text{mi}}$.

16

# 6 Application

We present an analysis to evaluate the effect of chronic obstructive pulmonary disease (COPD) on the development of herpes zoster (HZ). COPD is a chronic inflammatory lung disease that causes obstructed airflow from the lungs, which can cause systematic inflammation and dysregulate a patient's immune function. The hypothesis is that people with COPD are at increased risk of developing HZ. Yang et al. (2011) find a positive association between COPD and development of HZ; however, they do not control for important counfounders between COPD and HZ, for example, cigarette smoking and alcohol consumption.

We analyze the main data from the 2005 Longitudinal Health Insurance Database (LHID, Yang et al., 2011) and the validation data from the 2005 National Health Interview Survey conducted by the National Health Research Institute and the Bureau of Health Promotion in Taiwan (Lin and Chen, 2014). The 2005 LHID consist of $42,430$ subjects followed from the date of cohort entry on January 1, 2004 until the development of HZ or December 31, 2006, whichever came first. Among those, there are $8,486$ subjects with COPD, denoted by $A = 1$, and $33,944$ subjects without COPD, denoted by $A = 0$. The outcome $Y$ was the development of HZ during follow up (1, having HZ and 0, not having HZ). The observed prevalence of HZ among COPD and non-COPD subjects are $3.7\%$ and $2.2\%$ in the main data and $2.5\%$ and $0.8\%$ in the validation data.

The confounders $X$ available from the main data were age, sex, diabetes mellitus, hypertension, coronary artery disease, chronic liver disease, autoimmune disease, and cancer. However, important confounders $U$, including cigarette smoking and alcohol consumption, were not available. The validation data $\mathcal{O}_2$ use the same inclusion criteria as in the main study and consist of $1,148$ subjects who were comparable to the subjects in the main data. Among those, 244 subjects were diagnosed of COPD, and 904 subjects were not. In addition to all variables available from the main data, cigarette smoking and alcohol consumption were measured. In our formulation, the main data $\mathcal{O}_1$ combine the LHID data and the validation data. Table 4 in Lin and Chen (2014) shows summary statistics on demographic characteristics and comorbid disorders for COPD and Non-COPD subjects in the main and validation data. Because the common covariates in the main and validation data are comparable, it is reasonable to assume that the validation sample is a simple random sample from the main data. Moreover, the difference in distributions of alcohol consumption between COPD and non-COPD subjects is not statistical significant in the validation data. But, the COPD subjects tended to have higher cumulative smoking rates than the non-COPD subjects in the validation data.

We obtain the initial estimators applied solely to the validation data and the proposed estimators applied to both data. As suggested by the simulation in Section 5, we use the same type of the error-prone estimator as the initial estimator. Following Stürmer et al. (2005) and Lin and Chen (2014), we use the propensity score to accommodate the high-dimensional confounders. Specifically, we fit logistic regression models for the propensity score $e(X, U; \alpha)$ and the error-prone propensity score $e(X; \gamma)$ based on $\{(A_j, X_j, U_j) : j \in \mathcal{S}_2\}$ and $\{(A_i, X_i) : j \in \mathcal{S}_1\}$, respectively. We fit logistic regression models for the outcome mean function $\mu_a(X, U)$ based on a linear predictor $\{1, e(X, U; \hat{\alpha})\}^{\mathrm{T}} \beta_a$, and for $\mu_a(X)$ based on a linear predictor $\{1, e(X; \hat{\gamma})\}^{\mathrm{T}} \eta_a$, for $a = 0, 1$.

We first estimate the ACE $\tau$. Table 1 shows the results for the average COPD effect on the development of HZ. We find no big differences in the point estimates between our proposed estimators and the corresponding initial estimators, but large reductions in the estimated standard errors of the proposed estimators. As a result, all $95\%$ confidence intervals based on the initial estimators include 0, but the $95\%$ confidence intervals based on the proposed estimators do not include 0, except for $\hat{\tau}_{\mathrm{mat2\&mat}}$. As demonstrated by the simulation in Section 5, the variance reduction by utilizing the

main data is the smallest for the matching estimator. From the results, on average, COPD increases the percentage of developing HZ by $1.55\%$.

We also estimate the log of the causal risk ratio of HZ with COPD. The initial IPW estimate from the validation data is $\log \widehat{\mathrm{CRR}}_{\mathrm{IPW},2} = 1.10$ ($95\%$ confidence interval: $0.02,\ 2.18$). In contrast, the proposed estimate by using the error-prone IPW estimators is $\log \widehat{\mathrm{CRR}}_{\mathrm{IPW},2\&\mathrm{IPW}} = 0.57$ ($95\%$ confidence interval: $0.41,\ 0.72$), which is much more accurate than the initial IPW estimate.

Table 1: Point estimate, bootstrapped standard error and $95\%$ confidence interval

| | Est | SE | 95% CI | | Est | SE | 95% CI |
|---|---|---|---|---|---|---|---|
| $\hat{\tau}_{\mathrm{reg},2}$ | 0.0178 | 0.0112 | $(-0.0047, 0.0402)$ | $\hat{\tau}_{\mathrm{reg},2\&\mathrm{reg}}$ | 0.0155 | 0.0023 | $(0.0109, 0.0200)$ |
| $\hat{\tau}_{\mathrm{IPW},2}$ | 0.0175 | 0.0111 | $(-0.0048, 0.0398)$ | $\hat{\tau}_{\mathrm{IPW},2\&\mathrm{IPW}}$ | 0.0155 | 0.0024 | $(0.0108, 0.0202)$ |
| $\hat{\tau}_{\mathrm{AIPW},2}$ | 0.0179 | 0.0111 | $(-0.0044, 0.0402)$ | $\hat{\tau}_{\mathrm{AIPW},2\&\mathrm{AIPW}}$ | 0.0156 | 0.0024 | $(0.0109, 0.0203)$ |
| $\hat{\tau}_{\mathrm{mat},2}$ | 0.0077 | 0.0092 | $(-0.0106, 0.021)$ | $\hat{\tau}_{\mathrm{mat},2\&\mathrm{mat}}$ | 0.0079 | 0.0053 | $(-0.0027, 0.0184)$ |

# 7 Relaxing Assumption 1

In previous sections, we invoked Assumption 1 that $\mathcal{S}_2$ is a random sample from $\mathcal{S}_1$. We now relax this assumption and link our framework to existing methods for missing data. Let $I_i$ be the indicator of selecting unit $i$ into the validation data, i.e., $I_i = 1$ if $i \in \mathcal{S}_2$ and $I_i = 0$ if $i \notin \mathcal{S}_2$. Alternatively, $I_i$ can be viewed as the missingness indicator of $U_i$. Under Assumption 1, $I \perp\!\!\!\perp (A, X, U, Y)$; i.e., $U$ is missing completely at random. We now relax it to $I \perp\!\!\!\perp U \mid (A, X, Y)$, i.e., $U$ is missing at random. In this case, the selection of $\mathcal{S}_2$ from $\mathcal{S}_1$ can depend on a probability design, which is common in observational studies, e.g., an outcome-dependent two-phase sampling (Breslow et al., 2003, Wang et al., 2009).

We assume that each unit in the main data is subjected to an independent Bernoulli trial which determines whether the unit is selected into the validation data. For simplicity, we further assume that the inclusion probability $P(I = 1 \mid A, X, U, Y) = P(I = 1 \mid A, X, Y) \equiv \pi(A, X, Y)$ is known as in two-phase sampling. Otherwise, we need to fit a model for the missing data indicator $I$ given $(A, X, Y)$. We summarize the above in the following assumption.

**Assumption 6** $\{(I_i, A_i, X_i, U_i, Y_i) : i \in \mathcal{S}_1\}$ *are IID with* $I \perp\!\!\!\perp U \mid (A, X, Y)$. $\mathcal{S}_2$ *is selected from* $\mathcal{S}_1$ *with a known inclusion probability* $\pi(A, X, Y) > 0$.

In what follows, we use $\pi$ for $\pi(A, X, Y)$ and $\pi_j$ for $\pi(A_j, X_j, Y_j)$ for shorthand. Because of Assumption 6, we drop the indices $i$ and $j$ in the expectations, covariances, and variances, which are taken with respect to both the sampling and superpopulation models.

## 7.1 RAL estimators

For the illustration of RAL estimators, we focus on the AIPW estimator of the ACE $\tau$, because the regression imputation and inverse probability weighting estimators are its special cases. Let $\hat{\alpha}$ and $\hat{\beta}_a$ solve the weighted estimating equations $\sum_{j \in \mathcal{S}_2} \pi_j^{-1} S(A_j, X_j, U_j; \alpha) = 0$ and $\sum_{j \in \mathcal{S}_2} \pi_j^{-1} S_a(A_j, X_j, U_j, Y_j; \beta_a = 0$, and let $\alpha^*$ and $\beta_a^*$ satisfy $E\{S(A, X, U; \alpha^*)\} = 0$ and $E\{S_a(A, X, U, Y; \beta_a^*)\} = 0$. Under suitable regularity condition, $\hat{\alpha} \to \alpha^*$ and $\hat{\beta}_a \to \beta_a^*$ in probability, for $a = 0, 1$. Let the

initial estimator for $\tau$ be the Hajek-type estimator (Hájek, 1971):

$$\hat{\tau}_2 = \frac{\sum_{j \in \mathcal{S}_2} \pi_j^{-1} \hat{\tau}_{\text{AIPW},2,j}}{\sum_{j \in \mathcal{S}_2} \pi_j^{-1}}, \tag{22}$$

where $\hat{\tau}_{\text{AIPW},2,j}$ has the same form as (2). Under regularity conditions, Assumption 4 or 5, and Assumption 6, we show in the supplementary material that

$$\hat{\tau}_2 - \tau \cong n_1^{-1} \sum_{j \in \mathcal{S}_1} \pi_j^{-1} I_j \psi(A_j, X_j, U_j, Y_j), \tag{23}$$

where $\psi(A, X, U, Y)$ is given by (11). Because $\{\pi_j^{-1} I_j \psi(A_j, X_j, U_j, Y_j) : j \in \mathcal{S}_1\}$ are IID with mean 0, $\hat{\tau}_2$ is consistent for $\tau$.

Similarly, let $\hat{\gamma}_d$ and $\hat{\eta}_{d,a}$ solve the weighted estimating equation $\sum_{j \in \mathcal{S}_d} \pi_j^{-1} S(A_j, X_j; \gamma) = 0$ and $\sum_{j \in \mathcal{S}_d} \pi_j^{-1} S_a(A_j, X_j, Y_j; \eta_a) = 0$, and let $\gamma^*$ and $\eta_a^*$ satisfy $E\{S(A_j, X_j; \gamma^*)\} = 0$ and $E\{S_a(A_j, X_j, Y_j; \eta_a^*)\} = 0$ . Under suitable regularity condition, $\hat{\gamma}_d \to \gamma^*$ and $\hat{\eta}_{d,a} \to \eta_a^*$ in probability, for $a = 0, 1$ and $d = 1, 2$. Let the error-prone estimators be

$$\hat{\tau}_{1,\text{ep}} = n_1^{-1} \sum_{i \in \mathcal{S}_1} \hat{\tau}_{\text{AIPW},1,\text{ep},i}, \qquad \hat{\tau}_{2,\text{ep}} = \frac{\sum_{j \in \mathcal{S}_2} \pi_j^{-1} \hat{\tau}_{\text{AIPW},2,\text{ep},j}}{\sum_{j \in \mathcal{S}_2} \pi_j^{-1}}, \tag{24}$$

where $\hat{\tau}_{\text{AIPW},d,\text{ep},j}$ has the same form as (S8) in the supplementary material. Following a similar derivation for (23), we have

$$\hat{\tau}_{1,\text{ep}} - \tau_{\text{ep}} \cong n_1^{-1} \sum_{i \in \mathcal{S}_1} \phi(A_i, X_i, Y_i), \quad \hat{\tau}_{2,\text{ep}} - \tau_{\text{ep}} \cong n_1^{-1} \sum_{j \in \mathcal{S}_1} \pi_j^{-1} I_j \phi(A_j, X_j, Y_j), \tag{25}$$

where $\phi(A, X, Y)$ is given by (S9) in the supplementary material. Because both $\{\phi(A_i, X_i, Y_i) : i \in \mathcal{S}_1\}$ and $\{\pi_j^{-1} I_j \phi(A_j, X_j, Y_j) : j \in \mathcal{S}_1\}$ are IID with mean 0, $\hat{\tau}_{1,\text{ep}}$ and $\hat{\tau}_{2,\text{ep}}$ are consistent for $\tau_{\text{ep}}$.

**Theorem 4** *Under certain regularity conditions, (4) holds for the Hajek-type estimators (22) and (24), where $\rho = \text{plim}_{n_2 \to \infty} (n_2/n_1)$, $v_2 = \rho \times \text{var}\{\pi^{-1} I \psi(A, X, U, Y)\}$, $\Gamma = \rho \times \text{cov}\{\pi^{-1} I \psi(A, X, U, Y), (\pi^{-1} I - 1) \phi(A, X, Y)\}$ and $V = \rho \times \text{var}\{(\pi^{-1} I - 1) \phi(A, X, Y)\}$.*

Similar to Section 3.3, we can construct a consistent variance estimator for $\hat{\tau}$ by replacing the variances and covariance in Theorem 4 with their sample analogs.

## 7.2 Matching estimators

Recall that $\mathcal{J}_{d,V,l}$ is the index set of matches for unit $l$ based on data $\mathcal{O}_d$ and the matching variable $V$, which can be $(X, U)$ or $X$. Define $\delta_{d,V,(j,l)} = 1$ if $j \in \mathcal{J}_{d,V,l}$ and $\delta_{d,V,(j,l)} = 0$ otherwise. Now, we denote $K_{d,V,j} = \pi_j \sum_{l \in \mathcal{S}_d} \pi_l^{-1} 1\{A_l = 1 - A_j\} \delta_{d,V,(j,l)}$ as the weighted number of times that unit $j$ is used as a match. If $\pi_j$ is a constant for all $j \in \mathcal{S}_d$, then $K_{d,V,j}$ reduces to the number of times that unit $j$ is used as a match defined in Section 3.1, which justifies using the same notation as before.

Let the initial matching estimator for $\tau$ be the Hajek-type estimator:

19

$$\hat{\tau}_{\text{mat},2}^{(0)} = \left(\sum_{j\in\mathcal{S}_2} \pi_j^{-1}\right)^{-1} \sum_{j\in\mathcal{S}_2} \pi_j^{-1}(2A_j - 1)\left(Y_j - M^{-1}\sum_{l\in\mathcal{J}_{2,(X,U),j}} Y_l\right)$$

$$= \left(\sum_{j\in\mathcal{S}_2} \pi_j^{-1}\right)^{-1} \sum_{j\in\mathcal{S}_2} \pi_j^{-1}(2A_j - 1)\left\{1 + M^{-1}K_{2,(X,U),j}\right\}Y_j.$$

Let a bias-corrected matching estimator be

$$\hat{\tau}_{\text{mat},2} = \hat{\tau}_{\text{mat},2}^{(0)} - n_1^{-1/2}\hat{B}_2, \tag{26}$$

where

$$\hat{B}_2 = n_1^{-1/2}\sum_{j\in\mathcal{S}_2}\pi_j^{-1}(2A_j - 1)\left[M^{-1}\sum_{l\in\mathcal{J}_{2,(X,U),j}}\left\{\hat{\mu}_{1-A_j}(X_j, U_j) - \hat{\mu}_{1-A_j}(X_l, U_l)\right\}\right],$$

We show in the supplementary material that

$$\hat{\tau}_{\text{mat},2} - \tau \cong n_1^{-1}\sum_{j\in\mathcal{S}_2}\pi_j^{-1}\psi_{\text{mat},j}, \tag{27}$$

where $\psi_{\text{mat},j}$ is defined in (13) with the new definition of $K_{2,(X,U),j}$.

Similarly, we obtain error-prone matching estimators and express them as

$$\hat{\tau}_{\text{mat},1,\text{ep}} - \tau_{\text{ep}} \cong n_1^{-1}\sum_{j\in\mathcal{S}_1}\phi_{\text{mat},1,j}, \quad \hat{\tau}_{\text{mat},2,\text{ep}} - \tau_{\text{ep}} \cong n_1^{-1}\sum_{j\in\mathcal{S}_2}\pi_j^{-1}\phi_{\text{mat},2,j}, \tag{28}$$

where $\phi_{\text{mat},d,j}$ is defined in (15) with the new definition of $K_{d,X,j}$.

From the above decompositions, $\hat{\tau}_{\text{mat},2}$ is consistent for $\tau$, and $\hat{\tau}_{\text{mat},1,\text{ep}} - \hat{\tau}_{\text{mat},2,\text{ep}}$ is consistent for 0.

**Theorem 5** *Under certain regularity conditions, (4) holds for the estimators (26) and $\hat{\tau}_{\text{mat},d,\text{ep}}$ (d = 1, 2), where $\rho = \text{plim}_{n_2\to\infty}(n_2/n_1)$,*

$$v_2 = \rho \times \left(E\left[\frac{1-\pi}{\pi}\left\{\tau(X,U) - \tau\right\}^2\right]\right.$$

$$\left. + \text{plim}\left[n_1^{-1}\sum_{j\in\mathcal{S}_1}\frac{1-\pi_j}{\pi_j}\left\{1 + M^{-1}K_{2,(X,U),j}\right\}^2\sigma_{A_j}^2(X_j, U_j)\right]\right),$$

$$\Gamma = \rho \times E\left[\frac{1-\pi}{\pi}\left\{\mu_1(X,U) - \mu_0(X,U) - \tau\right\}\left\{\mu_1(X) - \mu_0(X) - \tau_{\text{ep}}\right\}\right]$$

$$+\rho \times \text{plim}\left[n_1^{-1}\sum_{j\in\mathcal{S}_1}\frac{1-\pi_j}{\pi_j}\left\{1 + M^{-1}K_{2,(X,U),j}\right\}\left(1 + M^{-1}K_{2,X,j}\right)\sigma_{A_j}^2(X_j, U_j)\right],$$

$$V = \rho \times E\left[\frac{1-\pi}{\pi}\left\{\mu_1(X) - \mu_0(X) - \tau_{\text{ep}}\right\}^2\right]$$

$$+\rho \times \text{plim}\left[n_1^{-1}\sum_{j\in\mathcal{S}_1}\frac{1-\pi_j}{\pi_j}\left(1 + M^{-1}K_{1,X,j}\right)^2\sigma_{A_j}^2(X_j)\right].$$

We can construct variance estimators based on the formulas in Theorem 5. However, this will again involve estimating the conditional variances $\sigma_0^2(x)$ and $\sigma_1^2(x)$. We recommend using a bootstrap variance estimator proposed in the next subsection.

## 7.3 A bootstrap variance estimation procedure

The asymptotic linear forms (23), (25), (27), and (28) are useful for the bootstrap variance estimation. For $b = 1, \ldots, B$, we construct the bootstrap replicates fas follows:

**Step 1.** Sample $n_1$ units from $\mathcal{S}_1$ with replacement as $\mathcal{S}_1^{*(b)}$.

**Step 2.** Compute the bootstrap replicates of $\hat{\tau}_2 - \tau$ and $\hat{\tau}_{d,\mathrm{ep}} - \tau_{\mathrm{ep}}$ as

$$\hat{\tau}_2^{(b)} - \hat{\tau}_2 = n_1^{-1} \sum_{i \in \mathcal{S}_1^{*(b)}} \pi_i^{-1} I_i \hat{\psi}_i,$$

$$\hat{\tau}_{1,\mathrm{ep}}^{(b)} - \hat{\tau}_{1,\mathrm{ep}} = n_1^{-1} \sum_{i \in \mathcal{S}_1^{*(b)}} \hat{\phi}_{1,i},$$

$$\hat{\tau}_{2,\mathrm{ep}}^{(b)} - \hat{\tau}_{2,\mathrm{ep}} = n_1^{-1} \sum_{i \in \mathcal{S}_1^{*(b)}} \pi_i^{-1} I_i \hat{\phi}_{2,i},$$

where $(\hat{\psi}_i, \hat{\phi}_{d,i})$ are the estimated versions of $(\psi_i, \phi_i)$ from $\mathcal{O}_d$ $(d = 1, 2)$.

We estimate $\Gamma, V$ and $v_2$ by (16)–(18) based on the above bootstrap replicates, and $\mathrm{var}(\hat{\tau})$ by (7), i.e., $\hat{v} = \hat{v}_2 - \hat{\Gamma}^{\mathrm{T}} \hat{V}^{-1} \hat{\Gamma}$.

**Theorem 6** *Under certain regularity conditions, $(\hat{\Gamma}, \hat{V}, \hat{v}_2, \hat{v})$ are consistent for $\{\Gamma, V, \mathrm{var}(\hat{\tau}_2), \mathrm{var}(\hat{\tau})\}$.*

For RAL estimators, we can also use the classical nonparametric bootstrap based on resampling the IID observations $\{(I_i, A_i, X_i, U_i, Y_i) : i \in \mathcal{S}_1\}$ and repeating the analysis as for the original data. The above bootstrap procedure based on resampling the linear forms are particularly useful for the matching estimator.

## 7.4 Connection with missing data

As a final remark, we express the proposed estimator in a linear form that has appeared in the missing data literature.

**Proposition 2** *Under certain regularity conditions and Assumption 6, $\hat{\tau}$ has an asymptotic linear form*

$$n_1^{1/2}(\hat{\tau} - \tau) = n_1^{-1/2} \sum_{i \in \mathcal{S}_1} \left\{ \frac{I_i}{\pi_i} \psi_i - \left( \frac{I_i}{\pi_i} - 1 \right) \Gamma V^{-1} \phi_i \right\} + o_P(1), \tag{29}$$

*where $\psi_i$ is $\psi(A_i, X_i, U_i, Y_i)$ for RAL estimators and $\psi_{\mathrm{mat},i}$ for the matching estimator, and a similar definition applies to $\phi_i$. Under Assumption 1, $\pi_i \equiv \rho$.*

Expression (29) is within a class of estimators in the missing data literature with the form

$$n_1^{1/2}(\hat{\tau} - \tau) = n_1^{-1/2} \sum_{i \in \mathcal{S}_1} \left\{ \frac{I_i}{\pi_i} s(A_i, X_i, U_i, Y_i) - \left( \frac{I_i}{\pi_i} - 1 \right) \kappa(A_i, X_i, Y_i) \right\} + o_P(1), \tag{30}$$

21

where $\pi = E(I \mid A, X, U, Y)$, $s(A, X, U, Y)$ satisfies $E\{s(A, X, U, Y)\} = 0$, and $s(A, X, U, Y)$ and $\kappa(A, X, Y)$ are square integrable. Given $s(A, X, U, Y)$, the optimal choice of $\kappa(A, X, Y)$ is $\kappa_{\mathrm{opt}}(A, X, Y) = E\{s(A, X, U, Y) \mid A, X, Y\}$, which minimizes the asymptotic variance of (30) (Robins et al., 1994, Wang et al., 2009). However, $\kappa_{\mathrm{opt}}(A, X, Y)$ requires a correct specification of the missing data model $f(U \mid A, X, Y)$. In our approach, instead of specifying the missing data model, we specify the error-prone estimators and utilize an estimator that is consistent for zero to improve the efficiency of the initial estimator. This is more attractive and closer to empirical practice than calculating $\kappa_{\mathrm{opt}}(A, X, Y)$, because practitioners only need to apply their favorite estimators to the main and validation data using existing software. See also Chen and Chen (2000) for a similar discussion in the regression context under Assumption 1.

# 8 Discussion

Depending on the roles in statistical inference, there are two types of *big data*: one with large-sample sizes and the other with richer covariates. In our discussion, the main observational data have a larger sample size, and the validation observational data have more covariates. Although some counterexamples exist (Pearl, 2009, 2010, Ding and Miratrix, 2015, Ding et al., 2017), it is more reliable to draw causal inference from the validation data. The proposed strategy is applicable even the number of covariates is high or the sample size is small in the validation data. In this case, we can consider $\hat{\tau}_2$ to be the double machine learning estimators (Chernozhukov et al., 2018) that use flexible machine learning methods for estimating regression and propensity score functions while retain the property in (4). Our framework allows for more efficient estimators of the causal effects by further combining information in the main data, without imposing any parametric models for the partially observed covariates. Coupled with the bootstrap, our estimators require only software implementations of standard estimators, and thus are attractive for practitioners who want to combine multiple observational data sources.

The key insight is to leverage an estimator of zero to improve the efficiency of the initial estimator. If a certain feature is transportable across datasets, we can construct a consistent estimator of zero. We have shown that if the validation data are simple random samples from the main data, the distribution of $(A, X, Y)$ is transportable from the validation data to the main data. We then construct a consistent estimator of zero by taking the difference of the estimators based on $(A, X, Y)$ from the two datasets. In the presence of heterogeneity between two data sources, the transportability of the whole distribution of $(A, X, Y)$ can be stringent. However, if we are willing to assume the conditional distribution of $Y$ given $(A, X)$ is transportable, we can then take the error prone estimators to be the regression coefficients of $Y$ on $(A, X)$ from the two datasets. As suggested by one of the reviewers, if the subgroups of two samples are comparable, we can construct the error prone estimators based on the subgroups. Similarly, this construction of error prone estimators can adopt to different transportability assumptions based on the subject matter knowledge.

In the worst case, the heterogeneity is intrinsic between the two samples, and we cannot construct two error prone estimators with the same limit. We can still conduct a sensitivity analysis combining two data. Instead of (4), we assume

$$n_2^{1/2} \begin{pmatrix} \hat{\tau}_2 - \tau \\ \hat{\tau}_{2,\mathrm{ep}} - \hat{\tau}_{1,\mathrm{ep}} - \delta \end{pmatrix} \to \mathcal{N} \left\{ 0_{L+1}, \begin{pmatrix} v_2 & \Gamma^{\mathrm{T}} \\ \Gamma & V \end{pmatrix} \right\}, \tag{31}$$

where $\delta$ is the sensitivity parameter, quantifying the systematic difference between $\hat{\tau}_{2,\mathrm{ep}}$ and $\hat{\tau}_{1,\mathrm{ep}}$. The adjusted estimator becomes $\hat{\tau}_{\mathrm{adj}}(\delta) = \hat{\tau}_2 - \hat{\Gamma}^{\mathrm{T}} \hat{V}^{-1} (\hat{\tau}_{2,\mathrm{ep}} - \hat{\tau}_{1,\mathrm{ep}} - \delta)$. With different values

of $\delta$, the estimator $\hat{\tau}_{\text{adj}}(\delta)$ can provide valuable insight on the impact of the heterogeneity of the two data, allowing an investigator to assess the extent to which the heterogeneity may alter causal inferences.

## Acknowledgments

## Supplemenary materials

The online supplemenatry material contains technical details and proofs. The R pacakge "IntegrativeCI" is available at `https://github.com/shuyang1987/IntegrativeCI` to perform the proposed estimators.

## References

Abadie, A. and Imbens, G. W. (2006). Large sample properties of matching estimators for average treatment effects, *Econometrica* **74**: 235–267.

Abadie, A. and Imbens, G. W. (2008). On the failure of the bootstrap for matching estimators, *Econometrica* **76**: 1537–1557.

Abadie, A. and Imbens, G. W. (2016). Matching on the estimated propensity score, *Econometrica* **84**: 781–807.

Angrist, J. D., Imbens, G. W. and Rubin, D. B. (1996). Identification of causal effects using instrumental variables, *J. Am. Stat. Assoc.* **91**: 444–455.

Antonelli, J., Zigler, C. and Dominici, F. (2017). Guided bayesian imputation to adjust for confounding when combining heterogeneous data sources in comparative effectiveness research, *Biostatistics* **18**: 553–568.

Bang, H. and Robins, J. M. (2005). Doubly robust estimation in missing data and causal inference models, *Biometrics* **61**: 962–973.

Bickel, P. J., Klaassen, C., Ritov, Y. and Wellner, J. (1993). *Efficient and Adaptive Inference in Semiparametric Models*, Johns Hopkins University Press, Baltimore.

Boonstra, P. S., Taylor, J. M. and Mukherjee, B. (2012). Incorporating auxiliary information for improved prediction in high-dimensional datasets: an ensemble of shrinkage approaches, *Biostatistics* **14**: 259–272.

Breslow, N., McNeney, B. and Wellner, J. A. (2003). Large sample theory for semiparametric regression models with two-phase, outcome dependent sampling, *Ann. Statist.* **31**: 1110–1139.

Cao, W., Tsiatis, A. A. and Davidian, M. (2009). Improving efficiency and robustness of the doubly robust estimator for a population mean with incomplete data, *Biometrika* **96**: 723–734.

Chang, T. and Kott, P. S. (2008). Using calibration weighting to adjust for nonresponse under a plausible model, *Biometrika* **95**: 555–571.

Chatterjee, N., Chen, Y. H., Maas, P. and Carroll, R. J. (2016). Constrained maximum likelihood estimation for model calibration using summary-level information from external big data sources, *J. Am. Stat. Assoc.* **111**: 107–117.

Chen, J. and Sitter, R. (1999). A pseudo empirical likelihood approach to the effective use of auxiliary information in complex surveys, *Statist. Sinica* **9**: 385–406.

Chen, J., Sitter, R. and Wu, C. (2002). Using empirical likelihood methods to obtain range restricted weights in regression estimators for surveys, *Biometrika* **89**: 230–237.

Chen, S. X., Leung, D. H. Y. and Qin, J. (2003). Information recovery in a study with surrogate endpoints, *J. Am. Stat. Assoc.* **98**: 1052–1062.

Chen, Y. H. (2002). Cox regression in cohort studies with validation sampling, *J. R. Stat. Soc. Ser. B.* **64**: 51–62.

Chen, Y. H. and Chen, H. (2000). A unified approach to regression analysis under double-sampling designs, *J. R. Stat. Soc. Ser. B.* **62**: 449–460.

Chernozhukov, V., Chetverikov, D., Demirer, M., Duflo, E., Hansen, C., Newey, W. and Robins, J. (2018). Double/debiased machine learning for treatment and structural parameters, *The Econometrics Journal* **21**: C1–C68.

Cochran, W. G. (2007). *Sampling Techniques*, 3 edn, New York: John Wiley & Sons, Inc.

Collaboration, F. S. (2009). Systematically missing confounders in individual participant data meta-analysis of observational cohort studies, *Statistics in Medicine* **28**: 1218–1237.

Crump, R., Hotz, V. J., Imbens, G. and Mitnik, O. (2006). Moving the goalposts: Addressing limited overlap in the estimation of average treatment effects by changing the estimand, *Technical report*, 330, National Bureau of Economic Research, Cambridge, MA.
**URL:** *http://www.nber.org/papers/T0330*

Deville, J.-C. and Särndal, C.-E. (1992). Calibration estimators in survey sampling, *J. Am. Stat. Assoc.* **87**: 376–382.

Ding, P. and Miratrix, L. W. (2015). To adjust or not to adjust? sensitivity analysis of M-bias and butterfly-bias, *Journal of Causal Inference* **3**: 41–57.

Ding, P., Vanderweele, T. and Robins, J. (2017). Instrumental variables as bias amplifiers with general outcome and confounding, *Biometrika* **104**: 291–302.

Efron, B. and Tibshirani, R. (1986). Bootstrap methods for standard errors, confidence intervals, and other measures of statistical accuracy, *Statistical Science* **1**: 54–75.

Enders, D., Kollhorst, B., Engel, S., Linder, R. and Pigeot, I. (2018). Comparison of multiple imputation and two-phase logistic regression to analyse two-phase case–control studies with rich phase 1: a simulation study, *Journal of Statistical Computation and Simulation* **88**: 2201–2214.

Fan, J. and Gijbels, I. (1996). *Local Polynomial Modelling and Its Applications*, CRC Press, Chapman & Hall, London.

Frangakis, C. E. and Rubin, D. B. (1999). Addressing complications of intention-to-treat analysis in the combined presence of all-or-none treatment-noncompliance and subsequent missing outcomes, *Biometrika* **86**: 365–379.

Freedman, D. A. (2008). Randomization does not justify logistic regression, *Statist. Sci.* **23**: 237–249.

Fuller, W. A. (2009). *Sampling Statistics*, Wiley, Hoboken, NJ.

Funk, M. J., Westreich, D., Wiesen, C., Stürmer, T., Brookhart, M. A. and Davidian, M. (2011). Doubly robust estimation of causal effects, *Am. J. Epidemiol.* **173**: 761–767.

Hájek, J. (1971). Comment on "an essay on the logical foundations of survey sampling, part one," by D. Basu, *in* V. P. V. P. Godambe and D. A. Sprott (eds), *Foundations of Survey Sampling*, Toronto: Holt, Rinehart, and Winston, p. 236.

Hansen, B. B. (2004). Full matching in an observational study of coaching for the SAT, *J. Am. Stat. Assoc.* **99**: 609–618.

Hansen, L. P. (1982). Large sample properties of generalized method of moments estimators, *Econometrica* **50**: 1029–1054.

Heckman, J. J., Ichimura, H. and Todd, P. E. (1997). Matching as an econometric evaluation estimator: Evidence from evaluating a job training programme, *Rev. Econ. Stud.* **64**: 605–654.

Hirano, K., Imbens, G. W. and Ridder, G. (2003). Efficient estimation of average treatment effects using the estimated propensity score, *Econometrica* **71**: 1161–1189.

Horvitz, D. G. and Thompson, D. J. (1952). A generalization of sampling without replacement from a finite universe, *J. Am. Stat. Assoc.* **47**: 663–685.

Imbens, G. W. (2004). Nonparametric estimation of average treatment effects under exogeneity: A review, *Rev. Econ. Stat.* **86**: 4–29.

Imbens, G. W. and Lancaster, T. (1994). Combining micro and macro data in microeconometric models, *The Review of Economic Studies* **61**: 655–680.

Jin, H. and Rubin, D. B. (2008). Principal stratification for causal inference with extended partial compliance, *J. Am. Stat. Assoc.* **103**: 101–111.

Kim, J. K., Kwon, Y. and Paik, M. C. (2016). Calibrated propensity score method for survey nonresponse in cluster sampling, *Biometrika* **103**: 461–473.

Kleiner, A., Talwalkar, A., Sarkar, P. and Jordan, M. I. (2014). A scalable bootstrap for massive data, *J. R. Stat. Soc. Ser. B.* **76**: 795–816.

Kott, P. S. (2006). Using calibration weighting to adjust for nonresponse and coverage errors, *Survey Methodology* **32**: 133–142.

Li, F., Morgan, K. L. and Zaslavsky, A. M. (2016). Balancing covariates via propensity score weighting, *J. Am. Stat. Assoc.* (DOI:10.1080/01621459.2016.1260466).

Lin, H. W. and Chen, Y. H. (2014). Adjustment for missing confounders in studies based on observational databases: 2-stage calibration combining propensity scores from primary and validation data, *Am. J. Epidemiol.* **180**: 308–317.

Lumley, T., Shaw, P. A. and Dai, J. Y. (2011). Connections between survey calibration estimators and semiparametric models for incomplete data, *International Statistical Review* **79**: 200–220.

Lunceford, J. K. and Davidian, M. (2004). Stratification and weighting via the propensity score in estimation of causal treatment effects: a comparative study, *Stat. Med.* **23**: 2937–2960.

Lunt, M., Glynn, R. J., Rothman, K. J., Avorn, J. and Stürmer, T. (2012). Propensity score calibration in the absence of surrogacy, *American Journal of Epidemiology* **175**: 1294–1302.

Mason, D. M. and Newton, M. A. (1992). A rank statistics approach to the consistency of a general bootstrap, *Ann. Statist.* **20**: 1611–1624.

McCandless, L. C., Richardson, S. and Best, N. (2012). Adjustment for missing confounders using external validation data and propensity scores, *J. Am. Stat. Assoc.* **107**: 40–51.

Newey, W. K. (1990). Semiparametric efficiency bounds, *Journal of Applied Econometrics* **5**: 99–135.

Newey, W. K. (1997). Convergence rates and asymptotic normality for series estimators, *J. Econometrics* **79**: 147–168.

Neyman, J. (1923). Sur les applications de la thar des probabilities aux experiences Agaricales: Essay de principle. English translation of excerpts by Dabrowska, D. and Speed, T., *Statist. Sci.* **5**: 465–472.

Neyman, J. (1938). Contribution to the theory of sampling human populations, *J. Am. Stat. Assoc.* **33**: 101–116.

Otsu, T. and Rai, Y. (2016). Bootstrap inference of matching estimators for average treatment effects, *J. Am. Stat. Assoc.* **112**: 1720–1732.

Pearl, J. (2009). Letter to the editor: Remarks on the method of propensity score, *Statistics in Medicine* **28**: 1420–1423.

Pearl, J. (2010). On a class of bias-amplifying variables that endanger effect estimates, *in* P. Grunwald and P. Spirtes (eds), *the Twenty-Sixth Conference on Uncertainty in Artificial Intelligence*, Association for Uncertainty in Artificial Intelligence, Corvallis, OR, pp. 425–432.

Politis, D. N., Romano, J. P. and Wolf, M. (1999). *Subsampling*, Springer-Verlag: New York.

Qin, J. (2000). Combining parametric and empirical likelihoods, *Biometrika* **87**: 484–490.

Robins, J. M., Rotnitzky, A. and Zhao, L. P. (1994). Estimation of regression coefficients when some regressors are not always observed, *J. Am. Stat. Assoc.* **89**: 846–866.

Rosenbaum, P. R. (1989). Optimal matching for observational studies, *J. Am. Stat. Assoc.* **84**: 1024–1032.

Rosenbaum, P. R. (2002). *Observational Studies*, 2 edn, Springer, New York.

Rosenbaum, P. R. and Rubin, D. B. (1983a). Assessing sensitivity to an unobserved binary covariate in an observational study with binary outcome, *J. R. Stat. Soc. Ser. B.* **45**: 212–218.

Rosenbaum, P. R. and Rubin, D. B. (1983b). The central role of the propensity score in observational studies for causal effects, *Biometrika* **70**: 41–55.

Rubin, D. B. (1973). Matching to remove bias in observational studies, *Biometrics* **29**: 159–183.

Rubin, D. B. (1974). Estimating causal effects of treatments in randomized and nonrandomized studies., *J. Educ. Psychol.* **66**: 688–701.

Rubin, D. B. (1987). *Multiple Imputation for Nonresponse in Surveys*, New York: Wiley.

Rubin, D. B. (2006). *Matched Sampling for Causal Effects*, Cambridge University Press, New York, NY.

Schneeweiss, S., Glynn, R. J., Tsai, E. H., Avorn, J. and Solomon, D. H. (2005). Adjusting for unmeasured confounders in pharmacoepidemiologic claims data using external information: the example of cox2 inhibitors and myocardial infarction, *Epidemiology* **16**: 17–24.

Shao, J. and Tu, D. (2012). *The Jackknife and Bootstrap*, Springer, New York.

Stuart, E. A. (2010). Matching methods for causal inference: A review and a look forward, *Statist. Sci.* **25**: 1–21.

Stürmer, T., Schneeweiss, S., Avorn, J. and Glynn, R. J. (2005). Adjusting effect estimates for unmeasured confounding with validation data using propensity score calibration, *Am. J. Epidemiol.* **162**: 279–289.

Stürmer, T., Schneeweiss, S., Rothman, K. J., Avorn, J. and Glynn, R. J. (2007). Performance of propensity score calibratio–a simulation study, *Am. J. Epidemiol.* **165**: 1110–1118.

Tsiatis, A. (2006). *Semiparametric Theory and Missing Data*, Springer, New York.

Wang, W., Scharfstein, D., Tan, Z. and MacKenzie, E. J. (2009). Causal inference in outcome-dependent two-phase sampling designs, *J. R. Stat. Soc. Ser. B.* **71**: 947–969.

Wang, X. and Wang, Q. (2015). Semiparametric linear transformation model with differential measurement error and validation sampling, *Journal of Multivariate Analysis* **141**: 67–80.

Wu, C. and Sitter, R. R. (2001). A model-calibration approach to using complete auxiliary information from survey data, *J. Am. Stat. Assoc.* **96**: 185–193.

Yang, Y. W., Chen, Y. H., Wang, K. H., Wang, C. Y. and Lin, H. W. (2011). Risk of herpes zoster among patients with chronic obstructive pulmonary disease: a population-based study, *Can. Med. Assoc. J.* **183**: 275–280.

27

# Supplementary materials for "Combining multiple observational data sources to estimate causal effects"

## S1   Proof of Remark 1

For simplicity, we show the result in Remark 1 for RAL estimators in Section 3.3. Suppose that we construct $\hat{\tau}_{1,\text{ep}}$ and $\hat{\tau}_{2,\text{ep}}$ based on different subsets subsets of $\mathcal{O}_1$ and $\mathcal{O}_2$, respectively, and that $\hat{\tau}_{2,\text{ep}} - \hat{\tau}_{1,\text{ep}}$ is consistent for zero. We express

$$n_2^{1/2}(\hat{\tau}_{2,\text{ep}} - \hat{\tau}_{1,\text{ep}}) \cong n_2^{-1/2} \left\{ \sum_{j \in \mathcal{S}_2} \gamma_j \phi(A_j, X_j, Y_j) + \sum_{i \in \mathcal{S}_1 \backslash \mathcal{S}_2} \gamma_i \phi(A_i, X_i, Y_i) \right\}, \quad \text{(S1)}$$

for a vector of constants $\gamma = \{\gamma_i \colon i \in \mathcal{S}_1\}$. Because we require that $\hat{\tau}_{1,\text{ep}} - \hat{\tau}_{2,\text{ep}}$ is consistent for zero, we have that $\sum_{i \in \mathcal{S}_1} \gamma_i = 0$. For example, if we compute error-prone estimators $\hat{\tau}_{1,\text{ep}}$ and $\hat{\tau}_{2,\text{ep}}$ based on $\mathcal{O}_1$ and $\mathcal{O}_2$, respectively, then this choice corresponds to $\gamma_j = n_2/n_1$ for $j \in \mathcal{S}_2$ and $\gamma_i = n_2/n_1 - 1$ for $i \in \mathcal{S}_1 \backslash \mathcal{S}_2$; if we compute error-prone estimators $\hat{\tau}_{1,\text{ep}}$ and $\hat{\tau}_{2,\text{ep}}$ based on $\mathcal{O}_1 \backslash \mathcal{O}_2$ and $\mathcal{O}_2$, respectively, then this choice corresponds to $\gamma_j = 1$ for $j \in \mathcal{S}_2$ and $\gamma_i = n_2/(n_1 - n_2)$ for $i \in \mathcal{S}_1 \backslash \mathcal{S}_2$.

We now show that our choice achieves the largest variance reduction $\Gamma^{\text{T}} V^{-1} \Gamma$ for $\hat{\tau}_2$ over all $\gamma$ satisfying $\sum_{i \in \mathcal{S}_1} \gamma_i = 0$.

Considering (8) and (S1), we calculate $\Gamma^{\text{T}} V^{-1} \Gamma$ as

$$\left\{ \text{plim} \ \frac{\left( n_2^{-1} \sum_{j \in \mathcal{S}_2} \gamma_j \right)^2}{n_2^{-1} \sum_{i \in \mathcal{S}_1} \gamma_i^2} \right\} \times V_{\psi|\phi},$$

where $V_{\psi|\phi} = \text{cov}\{\psi(A, X, U, Y), \phi(A, X, Y)\}^{\text{T}} \text{var}\{\phi(A, X, Y)\}^{-1} \text{cov}\{\psi(A, X, U, Y), \phi(A, X, Y)\}$.

Therefore, to maximize $\Gamma^{\text{T}} V^{-1} \Gamma$, it suffices to maximize $\left( \sum_{j \in \mathcal{S}_2} \gamma_j \right)^2 / (\sum_{i \in \mathcal{S}_1} \gamma_i^2)$, subject to the constraint $\sum_{i \in \mathcal{S}_1} \gamma_i = 0$. Using the Lagrange multiplier technique, we maximize the objective function

$$f(\gamma, \lambda) = \frac{\left( \sum_{j \in \mathcal{S}_2} \gamma_j \right)^2}{\sum_{i \in \mathcal{S}_1} \gamma_i^2} - 2\lambda \sum_{i \in \mathcal{S}_1} \gamma_i.$$

We take partial derivatives of $f(\gamma, \lambda)$ with respect to $\gamma$ and $\lambda$, and consider the system of equations

$$\frac{\partial f(\gamma, \lambda)}{\partial \gamma_j} = \frac{2 \left( \sum_{j \in \mathcal{S}_2} \gamma_j \right) \left( \sum_{i \in \mathcal{S}_1} \gamma_i^2 \right) - 2\gamma_j \left( \sum_{j \in \mathcal{S}_2} \gamma_j \right)^2}{\left( \sum_{i \in \mathcal{S}_1} \gamma_i^2 \right)^2} - 2\lambda = 0 \ \ (j \in \mathcal{S}_2), \quad \text{(S2)}$$

$$\frac{\partial f(\gamma, \lambda)}{\partial \gamma_i} = \frac{-2\gamma_i \left( \sum_{j \in \mathcal{S}_2} \gamma_j \right)^2}{\left( \sum_{i \in \mathcal{S}_1} \gamma_i^2 \right)^2} - 2\lambda = 0 \ \ (i \in \mathcal{S}_1 \backslash \mathcal{S}_2), \quad \text{(S3)}$$

and $\sum_{i \in \mathcal{S}_1} \gamma_i = 0$. The above system of equations suggests that the optimal choice $\gamma^*$ satisfies that $\gamma_i^* = c_1^*$ for $i \in \mathcal{S}_1 \backslash \mathcal{S}_2$, $\gamma_j^* = c_2^*$ for $i \in \mathcal{S}_2$, and $(n_1 - n_2)c_1^* + n_2 c_2^* = 0$ or $c_1^*/c_2^* = -n_2/(n_1 - n_2)$. This implies that the largest variance reduction is

$$
\left\{ \operatorname{plim} \frac{\left(n_2^{-1} \sum_{j \in \mathcal{S}_2} \gamma_j^*\right)^2}{n_2^{-1} \sum_{i \in \mathcal{S}_1} \gamma_i^{*2}} \right\} \times V_{\psi|\phi} = \left\{ \operatorname{plim} \frac{c_2^{*2}}{n_2^{-1} \left\{ (n_1 - n_2)c_1^{*2} + n_2 c_2^{*2} \right\}} \right\} \times V_{\psi|\phi} = (1 - \rho) \times V_{\psi|\phi}.
$$

As shown in Theorem 1, our choice achieves the largest variance reduction.

## S2 Proof of Remark 2

For a smooth function $f(x, y, z)$, let $\partial f(x, y, z)/\partial x$, $\partial f(x, y, z)/\partial y$ and $\partial f(x, y, z)/\partial z$ be its partial derivatives. By the Taylor expansion, we have

$$
\begin{aligned}
\hat{\tau} &= f(\hat{\tau}_2, \hat{\tau}_{1,\text{ep}}, \hat{\tau}_{2,\text{ep}}) \\
&\cong f(\tau, \tau_{\text{ep}}, \tau_{\text{ep}}) + \frac{\partial f}{\partial x}(\tau, \tau_{\text{ep}}, \tau_{\text{ep}})(\hat{\tau}_2 - \tau) \\
&\quad + \frac{\partial f}{\partial y^{\mathrm{T}}}(\tau, \tau_{\text{ep}}, \tau_{\text{ep}})(\hat{\tau}_{1,\text{ep}} - \tau_{\text{ep}}) + \frac{\partial f}{\partial z^{\mathrm{T}}}(\tau, \tau_{\text{ep}}, \tau_{\text{ep}})(\hat{\tau}_{2,\text{ep}} - \tau_{\text{ep}}) \\
&\equiv l_0 + l_1 \hat{\tau}_2 + l_2^{\mathrm{T}} \hat{\tau}_{1,\text{ep}} + l_3^{\mathrm{T}} \hat{\tau}_{2,\text{ep}}, \quad\quad\quad\quad (S4)
\end{aligned}
$$

where

$$
\begin{aligned}
l_0 &= f(\tau, \tau_{\text{ep}}, \tau_{\text{ep}}) - \frac{\partial f(\tau, \tau_{\text{ep}}, \tau_{\text{ep}})}{\partial x} \tau - \left\{ \frac{\partial f(\tau, \tau_{\text{ep}}, \tau_{\text{ep}})}{\partial y^{\mathrm{T}}} + \frac{\partial f(\tau, \tau_{\text{ep}}, \tau_{\text{ep}})}{\partial z^{\mathrm{T}}} \right\} \tau_{\text{ep}}, \\
l_1 &= \frac{\partial f(\tau, \tau_{\text{ep}}, \tau_{\text{ep}})}{\partial x}, \quad l_2 = \frac{\partial f(\tau, \tau_{\text{ep}}, \tau_{\text{ep}})}{\partial y}, \quad l_3 = \frac{\partial f(\tau, \tau_{\text{ep}}, \tau_{\text{ep}})}{\partial z}.
\end{aligned}
$$

Because $\hat{\tau}$ is consistent for $\tau$, letting $n_2$ go to infinity in (S4), we have $\tau = l_0 + l_1 \tau + l_2^{\mathrm{T}} \tau_{\text{ep}} + l_3^{\mathrm{T}} \tau_{\text{ep}}$ for all $\tau$ and $\tau_{\text{ep}}$. Then, it follows $l_0 = 0$, $l_1 = 1$, and $l_2 = -l_3$. Therefore, $\hat{\tau} \cong \hat{\tau}_2 - l_3^{\mathrm{T}}(\hat{\tau}_{2,\text{ep}} - \hat{\tau}_{1,\text{ep}})$. By the joint normality of $(\hat{\tau}, \hat{\tau}_{1,\text{ep}}, \hat{\tau}_{2,\text{ep}})$, the optimal $l_3$ minimizing $\operatorname{var}\{\hat{\tau}_2 - l_3^{\mathrm{T}}(\hat{\tau}_{2,\text{ep}} - \hat{\tau}_{1,\text{ep}})\}$ is $\Gamma^{\mathrm{T}} V^{-1}$. Therefore, the proposed estimator in (5) is optimal among the class of consistent estimators that is a smooth function of $(\hat{\tau}, \hat{\tau}_{1,\text{ep}}, \hat{\tau}_{2,\text{ep}})$.

## S3 Proof of Remark 3

We can view $v_2 - \Gamma^{\mathrm{T}} V^{-1} \Gamma$ and $\Gamma_1^{\mathrm{T}} V_{11}^{-1} \Gamma_1$ as the conditional variances in a multivariate normal vector. Then the conclusion holds immediately because conditioning on more variables will decrease the variance for a multivariate normal vector. For algebraic completeness, we give a formal proof below.

Decompose the $L$-dimensional error-prone estimator $\hat{\tau}_{d,\text{ep}}$ into two components with $L_1$ and $L_2$ dimensions respectively. Then, $\Gamma$ and $V$ have the corresponding partitions:

$$
\Gamma = \begin{pmatrix} \Gamma_1 \\ \Gamma_2 \end{pmatrix}, \quad V = \begin{pmatrix} V_{11} & V_{12} \\ V_{12}^{\mathrm{T}} & V_{22} \end{pmatrix}.
$$

We assume $V$ is invertable, and therefore $V_{11}$ and $V_{22}$ are also invertable. To show that increasing the dimension of the error-prone estimator would not decrease the efficiency gain, according to (6), it suffices to show that $\Gamma^{\mathrm{T}}V^{-1}\Gamma \geq \Gamma_1^{\mathrm{T}}V_{11}^{-1}\Gamma_1$. Toward this end, note that

$$V^{-1} = \begin{pmatrix} V_{11}^{-1} + V_{11}^{-1}V_{12}(V_{22} - V_{12}^{\mathrm{T}}V_{11}^{-1}V_{12})^{-1}V_{12}^{\mathrm{T}}V_{11}^{-1} & -V_{11}^{-1}V_{12}(V_{22} - V_{12}^{\mathrm{T}}V_{11}^{-1}V_{12})^{-1} \\ -(V_{22} - V_{12}^{\mathrm{T}}V_{11}^{-1}V_{12})^{-1}V_{12}^{\mathrm{T}}V_{11}^{-1} & (V_{22} - V_{12}^{\mathrm{T}}V_{11}^{-1}V_{12})^{-1} \end{pmatrix},$$

so that

$$\begin{aligned}
\Gamma^{\mathrm{T}}V^{-1}\Gamma &= \Gamma_1^{\mathrm{T}}V_{11}^{-1}\Gamma_1 + \Gamma_1^{\mathrm{T}}\{V_{11}^{-1}V_{12}(V_{22} - V_{12}^{\mathrm{T}}V_{11}^{-1}V_{12})^{-1}V_{12}^{\mathrm{T}}V_{11}^{-1}\}\Gamma_1 \\
&\quad -\Gamma_1^{\mathrm{T}}V_{11}^{-1}V_{12}(V_{22} - V_{12}^{\mathrm{T}}V_{11}^{-1}V_{12})^{-1}\Gamma_2 \\
&\quad -\Gamma_2^{\mathrm{T}}(V_{22} - V_{12}^{\mathrm{T}}V_{11}^{-1}V_{12})^{-1}V_{12}^{\mathrm{T}}V_{11}^{-1}\Gamma_1 \\
&\quad +\Gamma_2^{\mathrm{T}}(V_{22} - V_{12}^{\mathrm{T}}V_{11}^{-1}V_{12})^{-1}\Gamma_2. \tag{S5}
\end{aligned}$$

Let $\theta = (V_{22} - V_{12}^{\mathrm{T}}V_{11}^{-1}V_{12})^{-1/2}V_{12}^{\mathrm{T}}V_{11}^{-1}\Gamma_1$ and $\delta = (V_{22} - V_{12}^{\mathrm{T}}V_{11}^{-1}V_{12})^{-1/2}\Gamma_2$. Then, (S5) becomes

$$\Gamma^{\mathrm{T}}V^{-1}\Gamma = \Gamma_1^{\mathrm{T}}V_{11}^{-1}\Gamma_1 + \theta^{\mathrm{T}}\theta - \theta^{\mathrm{T}}\delta - \delta^{\mathrm{T}}\theta + \delta^{\mathrm{T}}\delta = \Gamma_1^{\mathrm{T}}V_{11}^{-1}\Gamma_1 + (\theta - \delta)^{\mathrm{T}}(\theta - \delta) \geq \Gamma_1^{\mathrm{T}}V_{11}^{-1}\Gamma_1.$$

## S4   Proof for Theorem 1

The asymptotic normality holds for $n_2^{1/2}(\hat{\tau}_2 - \tau, \hat{\tau}_{2,\mathrm{ep}} - \hat{\tau}_{1,\mathrm{ep}})$ by the moment conditions for the RAL estimators and the central limit theorem. We then show the asymptotic variance formula in (4).

Based on the linear form in (8), $v_2 = \mathrm{var}\{\psi(A, X, U, Y)\}$. Based on the linear form in (9),

$$n_2^{1/2}(\hat{\tau}_{2,\mathrm{ep}} - \hat{\tau}_{1,\mathrm{ep}}) \cong n_2^{-1/2}\sum_{j \in \mathcal{S}_2}\phi(A_j, X_j, Y_j) - \rho^{1/2}n_1^{-1/2}\sum_{i \in \mathcal{S}_1}\phi(A_i, X_i, Y_i). \tag{S6}$$

The first term in (S6) contributes $\mathrm{var}\{\phi(A, X, Y)\}$, the second term contributes $\rho\,\mathrm{var}\{\phi(A, X, Y)\}$, and their correlation contributes $-2\rho^{1/2}\rho^{1/2}\mathrm{var}\{\phi(A, X, Y)\}$. Therefore,

$$\begin{aligned}
V &= \mathrm{var}\{\phi(A, X, Y)\} + \rho\,\mathrm{var}\{\phi(A, X, Y)\} - 2\rho^{1/2}\rho^{1/2}\mathrm{var}\{\phi(A, X, Y)\} \\
&= (1 - \rho) \times \mathrm{var}\{\phi(A, X, Y)\}.
\end{aligned}$$

To obtain the expression for $\Gamma$, we re-write

$$n_2^{1/2}(\hat{\tau}_{2,\mathrm{ep}} - \hat{\tau}_{1,\mathrm{ep}}) \cong n_2^{-1/2}\sum_{j \in \mathcal{S}_2}(1 - \rho)\phi(A_j, X_j, Y_j) - \rho n_2^{-1/2}\sum_{i \in \mathcal{S}_1 \backslash \mathcal{S}_2}\phi(A_i, X_i, Y_i).$$

Because observations in $\mathcal{S}_2$ and $\mathcal{S}_1\backslash\mathcal{S}_2$ are independent, simple calculations give

$$\Gamma = (1 - \rho)\mathrm{cov}\{\psi(A, X, U, Y), \phi(A, X, Y)\}.$$

## S5   Influence functions for RAL estimators based on the validation data

We review the influence functions for $\hat{\tau}_{\mathrm{reg},2}$ and $\hat{\tau}_{\mathrm{IPW},2}$ based on the validation data.

**Lemma S1 (Regression imputation)** *Under Assumption 4, $\hat{\tau}_{\text{reg},2}$ has the influence function*

$$\psi_{\text{reg}}(A_j, X_j, U_j, Y_j) = \mu_1(X_j, U_j; \beta_1^*) - \mu_0(X_j, U_j; \beta_0^*) - \tau$$

$$-E\left\{\frac{\partial\mu_1(X, U; \beta_1^*)}{\partial\beta_1^{\text{T}}}\right\} E\left\{\frac{\partial S_1(A, X, U, Y; \beta_1^*)}{\partial\beta_1^{\text{T}}}\right\}^{-1} S_1(A_j, X_j, U_j, Y_j; \beta_1^*)$$

$$+E\left\{\frac{\partial\mu_0(X, U; \beta_0^*)}{\partial\beta_0^{\text{T}}}\right\} E\left\{\frac{\partial S_0(A, X, U, Y; \beta_0^*)}{\partial\beta_0^{\text{T}}}\right\}^{-1} S_0(A_j, X_j, U_j, Y_j; \beta_0^*).$$

**Lemma S2 (Inverse probability weighting)** *Under Assumption 5, $\hat{\tau}_{\text{IPW},2}$ has the influence function*

$$\psi_{\text{IPW}}(A_j, X_j, U_j, Y_j) = \frac{A_j Y_j}{e(X_j, U_j; \alpha^*)} - \frac{(1 - A_j)Y_j}{1 - e(X_j, U_j; \alpha^*)} - \tau - H_{\text{IPW}}\Sigma_{\alpha\alpha}^{-1} S(A_j, X_j, U_j; \alpha^*),$$

*where*

$$H_{\text{IPW}} = E\left(\left[\frac{AY}{e(X, U; \alpha^*)^2} - \frac{(1 - A)Y}{\{1 - e(X, U; \alpha^*)\}^2}\right]\frac{\partial e(X, U; \alpha^*)}{\partial\alpha^{\text{T}}}\right).$$

Lunceford and Davidian (2004) derive Lemmas S1 and S2.

## S6  Proof of Lemma 1

We write $\hat{\tau}_{\text{AIPW},2} = \hat{\tau}_{\text{AIPW},2}(\hat{\alpha}, \hat{\beta}_0, \hat{\beta}_1)$ to emphasize its dependence on the parameter estimates $(\hat{\alpha}, \hat{\beta}_0, \hat{\beta}_1)$. By the Taylor expansion,

$$\hat{\tau}_{\text{AIPW},2}(\hat{\alpha}, \hat{\beta}_0, \hat{\beta}_1)$$

$$\cong \hat{\tau}_{\text{AIPW},2}(\alpha^*, \beta_0^*, \beta_1^*) + E\left\{\frac{\partial\hat{\tau}_{\text{AIPW},2}(\alpha^*, \beta_0^*, \beta_1^*)}{\partial\alpha^{\text{T}}}\right\}(\hat{\alpha} - \alpha^*)$$

$$+E\left\{\frac{\partial\hat{\tau}_{\text{AIPW},2}(\alpha^*, \beta_0^*, \beta_1^*)}{\partial\beta_0^{\text{T}}}\right\}(\hat{\beta}_0 - \beta_0^*) + E\left\{\frac{\partial\hat{\tau}_{\text{AIPW},2}(\alpha^*, \beta_0^*, \beta_1^*)}{\partial\beta_1^{\text{T}}}\right\}(\hat{\beta}_1 - \beta_1^*)$$

$$\cong \hat{\tau}_{\text{AIPW},2}(\alpha^*, \beta_0^*, \beta_1^*)$$

$$+n_2^{-1}\sum_{j\in\mathcal{S}_2} E\left\{\frac{\partial\hat{\tau}_{\text{AIPW},2}(\alpha^*, \beta_0^*, \beta_1^*)}{\partial\alpha^{\text{T}}}\right\} E\left\{\frac{\partial S(A, X, U; \alpha^*)}{\partial\alpha^{\text{T}}}\right\}^{-1} S(A_j, X_j, U_j; \alpha^*) \quad\text{(S7)}$$

$$-n_2^{-1}\sum_{j\in\mathcal{S}_2} E\left\{\frac{\partial\hat{\tau}_{\text{AIPW},2}(\alpha^*, \beta_0^*, \beta_1^*)}{\partial\beta_0^{\text{T}}}\right\} E\left\{\frac{\partial S_0(A, X, U, Y; \beta_0^*)}{\partial\beta_0^{\text{T}}}\right\}^{-1} S_0(A_j, X_j, U_j, Y_j; \beta_0^*)$$

$$-n_2^{-1}\sum_{j\in\mathcal{S}_2} E\left\{\frac{\partial\hat{\tau}_{\text{AIPW},2}(\alpha^*, \beta_0^*, \beta_1^*)}{\partial\beta_1^{\text{T}}}\right\} E\left\{\frac{\partial S_0(A, X, U, Y; \beta_0^*)}{\partial\beta_1^{\text{T}}}\right\}^{-1} S_1(A_j, X_j, U_j, Y_j; \beta_1^*).$$

We have the following calculations:

$$E\left\{\frac{\partial\hat{\tau}_{\text{AIPW},2}(\alpha^*, \beta_0^*, \beta_1^*)}{\partial\alpha^{\text{T}}}\right\} = H_{\text{AIPW}},$$

$$E\left\{\frac{\partial\hat{\tau}_{\text{AIPW},2}(\alpha^*, \beta_0^*, \beta_1^*)}{\partial\beta_0^{\text{T}}}\right\} = -E\left[\left\{1 - \frac{1 - A}{1 - e(X, U; \alpha^*)}\right\}\frac{\partial\mu_0(X; \beta_0^*)}{\partial\beta_0^{\text{T}}}\right],$$

$$E\left\{\frac{\partial\hat{\tau}_{\text{AIPW},2}(\alpha^*, \beta_0^*, \beta_1^*)}{\partial\beta_1^{\text{T}}}\right\} = E\left[\left\{1 - \frac{A}{e(X, U; \alpha^*)}\right\}\frac{\partial\mu_1(X, U; \beta_1^*)}{\partial\beta_1^{\text{T}}}\right].$$

Under Assumption 4, $H_{\mathrm{AIPW}} = 0$. Under Assumption 5,

$$E\left\{\frac{\partial S(A, X, U; \alpha^*)}{\partial \alpha^{\mathrm{T}}}\right\} = E\left\{S(A, X, U; \alpha^*)^{\otimes 2}\right\} = \Sigma_{\alpha\alpha}.$$

Therefore, we can always replace $E\left\{\partial S(A, X, U; \alpha^*)/\partial \alpha^{\mathrm{T}}\right\}$ by $\Sigma_{\alpha\alpha}$ in expression (S7) if either Assumption 4 or 5 holds. Thus, we can derive the influence function for the AIPW estimator.

## S7 Lemmas for error-prone estimators

The error-prone estimators can be viewed as the initial estimators in Examples 3–5 with $U$ being null. The following results are similar to Lemmas S1–1 with a subtle difference that neither the propensity score or the outcome model is correctly specified. For completeness, we establish the results for the error-prone estimators.

Let $\mu_a(X; \eta_a)$ be a working model for $\mu_a(X)$, for $a = 0, 1$, and $e(X; \gamma)$ be a working model for $e(X)$. Let $S_a(A, X, Y; \eta_a)$ be the estimating function for $\eta_a$, e.g.,

$$S_a(A, X, Y; \eta_a) = \frac{\partial \mu_a(X; \eta_a)}{\partial \eta_a}\{Y - \mu_a(X; \eta_a)\},$$

for $a = 0, 1$. Let $S(A, X; \gamma)$ be the estimating function for $\gamma$, e.g.,

$$S(A, X; \gamma) = \frac{A - e(X; \gamma)}{e(X; \gamma)\{1 - e(X; \gamma)\}}\frac{\partial e(X; \gamma)}{\partial \gamma}.$$

We further define

$$\Sigma_{\gamma\gamma} = E\left\{\frac{\partial S(A, X; \gamma^*)}{\partial \gamma^{\mathrm{T}}}\right\},$$

which does not necessarily equal to the Fisher information matrix for a misspecified propensity score model. Let $\hat{\eta}_{d,a}$ $(a = 0, 1)$ and $\hat{\gamma}_d$ be the estimators solving the corresponding estimating equations based on $\mathcal{S}_d$, and let $\eta_a^*$ $(a = 0, 1)$ and $\gamma^*$ satisfy $E\{S_a(A, X, Y; \eta_a^*)\} = 0$ $(a = 0, 1)$ and $E\{S(A, X; \gamma^*)\} = 0$. Under suitable regularity conditions, $\hat{\eta}_{d,a}$ $(a = 0, 1)$ and $\hat{\gamma}_d$ have probability limits $\eta_a^*$ $(a = 0, 1)$ and $\gamma^*$.

**Lemma S3 (Regression imputation)** *The error-prone regression imputation estimator for $\tau$ is $\hat{\tau}_{\mathrm{reg},d,\mathrm{ep}} = n_d^{-1}\sum_{j\in\mathcal{S}_d}\hat{\tau}_{\mathrm{reg},d,\mathrm{ep},j}$, where*

$$\hat{\tau}_{\mathrm{reg},d,\mathrm{ep},j} = \mu_1(X_j; \hat{\eta}_{d,1}) - \mu_0(X_j; \hat{\eta}_{d,0}).$$

*It has probability limit $\tau_{\mathrm{ep}} = E\left\{\mu_1(X; \eta_1^*) - \mu_0(X; \eta_0^*)\right\}$ and influence function*

$$
\begin{aligned}
\phi_{\mathrm{reg}}(A_j, X_j, Y_j) &= \mu_1(X_j; \eta_1^*) - \mu_0(X_j; \eta_0^*) - \tau_{\mathrm{ep}} \\
&\quad - E\left\{\frac{\partial \mu_1(X; \eta_1^*)}{\partial \eta_1^{\mathrm{T}}}\right\} E\left\{\frac{\partial S_1(A, X, Y; \eta_1^*)}{\partial \eta_1^{\mathrm{T}}}\right\}^{-1} S_1(A_j, X_j, Y_j; \eta_1^*) \\
&\quad + E\left\{\frac{\partial \mu_0(X; \eta_0^*)}{\partial \eta_0^{\mathrm{T}}}\right\} E\left\{\frac{\partial S_0(A, X, Y; \eta_0^*)}{\partial \eta_0^{\mathrm{T}}}\right\}^{-1} S_0(A_j, X_j, Y_j; \eta_0^*).
\end{aligned}
$$

**Lemma S4 (Inverse probability weighting)** *The error-prone IPW estimator for $\tau$ is $\hat{\tau}_{\text{IPW},d,\text{ep}} = n_d^{-1} \sum_{j \in \mathcal{S}_d} \hat{\tau}_{\text{IPW},d,\text{ep},j}$, where*

$$\hat{\tau}_{\text{IPW},d,\text{ep},j} = \frac{A_j Y_j}{e(X_j; \hat{\gamma}_d)} - \frac{(1 - A_j)Y_j}{1 - e(X_j; \hat{\gamma}_d)}.$$

*It has probability limit*

$$\tau_{\text{ep}} = E\left\{ \frac{AY}{e(X; \gamma^*)} - \frac{(1 - A)Y}{1 - e(X; \gamma^*)} \right\}$$

*and influence function*

$$\phi_{\text{IPW}}(A_j, X_j, Y_j) = \left\{ \frac{A_j Y_j}{e(X_j; \gamma^*)} - \frac{(1 - A_j)Y_j}{1 - e(X_j; \gamma^*)} - \tau_{\text{ep}} \right\} - \tilde{H}_{\text{IPW}} \Sigma_{\gamma\gamma}^{-1} S(A_j, X_j; \gamma^*),$$

*where*

$$\tilde{H}_{\text{IPW}} = E\left( \left[ \frac{AY}{e(X; \gamma^*)^2} - \frac{(1 - A)Y}{\{1 - e(X; \gamma^*)\}^2} \right] \frac{\partial e(X; \gamma^*)}{\partial \gamma^{\text{T}}} \right).$$

**Lemma S5 (Augmented inverse probability weighting)** *Define the residual outcome as $\tilde{R}_{d,j} = Y_j - \mu_1(X_j; \hat{\eta}_{d,1})$ for treated units and $\tilde{R}_{d,j} = Y_j - \mu_0(X_j; \hat{\eta}_{d,0})$ for control units, for $d = 1, 2$. The error-prone AIPW estimator for $\tau$ is $\hat{\tau}_{\text{AIPW},d,\text{ep}} = n_d^{-1} \sum_{j \in \mathcal{S}_d} \hat{\tau}_{\text{AIPW},d,\text{ep},j}$, where*

$$\hat{\tau}_{\text{AIPW},d,\text{ep},j} = \frac{A_j \tilde{R}_{d,j}}{e(X_j; \hat{\gamma}_d)} + \mu_1(X_j; \hat{\eta}_{d,1}) - \frac{(1 - A_j)\tilde{R}_{d,j}}{1 - e(X_j; \hat{\gamma}_d)} - \mu_0(X_j; \hat{\eta}_{d,0}). \tag{S8}$$

*It has probability limit*

$$\tau_{\text{ep}} = E\left[ \frac{AY}{e(X; \gamma^*)} + \left\{ 1 - \frac{A}{e(X; \gamma^*)} \right\} \mu_1(X; \eta_1^*) - \frac{(1 - A)Y}{1 - e(X; \gamma^*)} - \left\{ 1 - \frac{1 - A}{1 - e(X; \gamma^*)} \right\} \mu_0(X; \eta_0^*) \right]$$

*and influence function*

$$
\begin{aligned}
&\phi_{\text{AIPW}}(A_j, X_j, Y_j) \\
&= \frac{A_j Y_j}{e(X_j; \gamma^*)} + \left\{ 1 - \frac{A_j}{e(X_j; \gamma^*)} \right\} \mu_1(X_j; \eta_1^*) \\
&\quad - \frac{(1 - A_j)Y_j}{1 - e(X_j; \gamma^*)} - \left\{ 1 - \frac{1 - A_j}{1 - e(X_j; \gamma^*)} \right\} \mu_0(X_j; \eta_0^*) - \tau_{\text{ep}} + \tilde{H}_{\text{AIPW}} \Sigma_{\gamma\gamma}^{-1} S(A_j, X_j; \gamma^*) \\
&\quad + E\left[ \left\{ 1 - \frac{1 - A}{1 - e(X; \gamma^*)} \right\} \frac{\partial \mu_0(X; \eta_0^*)}{\partial \eta_0^{\text{T}}} \right] E\left\{ \frac{\partial S_0(A, X, Y; \eta_0^*)}{\partial \eta_0^{\text{T}}} \right\}^{-1} S_0(A_j, X_j, Y_j; \eta_0^*) \\
&\quad - E\left[ \left\{ 1 - \frac{A}{e(X; \gamma^*)} \right\} \frac{\partial \mu_1(X; \eta_1^*)}{\partial \eta_1^{\text{T}}} \right] E\left\{ \frac{\partial S_1(A, X, Y; \eta_1^*)}{\partial \eta_1^{\text{T}}} \right\}^{-1} S_1(A_j, X_j, Y_j; \eta_1^*), \tag{S9}
\end{aligned}
$$

*where*

$$\tilde{H}_{\text{AIPW}} = E\left( \left[ \frac{A\{Y - \mu_1(X; \eta_1^*)\}}{e(X; \gamma^*)^2} - \frac{(1 - A)\{Y - \mu_0(X; \eta_0^*)\}}{\{1 - e(X; \gamma^*)\}^2} \right] \frac{\partial e(X; \gamma^*)}{\partial \gamma^{\text{T}}} \right).$$

The results in Lemmas S3 and S4 can be obtained by the Taylor expansion. The proof for Lemma S5 is similar to that for Lemma 1 and therefore omitted.

## S8 Assumptions for the matching estimator

We review the assumptions for the matching estimators, which can also be found in Abadie and Imbens (2006).

**Assumption S1 (Population distributions)** *(i) $(X, U)$ is continuously distributed on a compact and convex support. The density of $(X, U)$ is bounded and bounded away from zero on its support.*

*(ii) For $a = 0, 1$, $\mu_a(x, u)$ and $\sigma_a^2(x, u)$ are Lipschitz, $\sigma_a^2(x, u)$ is bounded away from zero, and $E(Y^4 \mid A = a, X = x, U = u)$ is bounded uniformly over its support.*

*(iii) for $a = 0, 1$, $\mu_a(x)$ and $\sigma_a^2(x)$ are Lipschitz, $\sigma_a^2(x)$ is bounded away from zero, and $E(Y^4 \mid A = a, X = x)$ is bounded uniformly over its support.*

Assumption S1 (i) can be relaxed by allowing $(X, U)$ to have discrete components. We only need to obtain results on each level of discrete covariates and derive the same result. Assumption S1 (ii) requires the conditional mean and variance functions to be bounded and satisfy certain smoothness conditions, which are rather mild. In fact, Assumption S1 (ii) implies Assumption S1 (iii). To be more transparent, we state Assumption S1 (iii) explicitly.

**Assumption S2 (Estimators of mean functions)** *For $a = 0, 1$, the estimators $\hat{\mu}_a(x, u)$ and $\hat{\mu}_a(x)$ satisfy the following asymptotic conditions: (i) $|\hat{\mu}_a(x, u) - \mu_a(x, u)| = o_P\left\{n_2^{-1/2 + 1/\dim(x,u)}\right\}$; and (ii) $|\hat{\mu}_a(x) - \mu_a(x)| = o_P\left\{n_2^{-1/2 + 1/\dim(x)}\right\}$.*

If $\hat{\mu}_a(x, u)$ and $\hat{\mu}_a(x)$ are obtained under correctly specified parametric models, then Assumption S2 holds. If $\hat{\mu}_a(x, u)$ and $\hat{\mu}_a(x)$ are obtained using nonparametric methods, such as power series regression (Newey, 1997) or kernel regression (Fan and Gijbels, 1996) estimators, we need to select their tuning parameters properly to ensure Assumption S2. Assumption S2 is needed so that the bias correction terms achieve fast convergence; e.g., $n_2^{1/2}(\hat{B}_2 - B_2) \to 0$ in probability, as $n_2 \to \infty$.

## S9 Proof of Theorem 2

First, we express

$$
\begin{aligned}
&n_2^{1/2}(\hat{\tau}_{\mathrm{mat},2} - \tau) \\
={}& n_2^{-1/2} \sum_{j \in \mathcal{S}_2} \hat{\psi}_{\mathrm{mat},j} + o_P(1) \\
={}& n_2^{-1/2} \sum_{j \in \mathcal{S}_2} \{\mu_1(X_j, U_j) - \mu_0(X_j, U_j) - \tau\} \\
&+ n_2^{-1/2} \sum_{j \in \mathcal{S}_2} (2A_j - 1)\left\{1 + M^{-1}K_{2,(X,U),j}\right\}\left\{Y_j - \mu_{A_j}(X_j, U_j)\right\} + o_P(1).
\end{aligned}
$$

Second, let

$$
\begin{aligned}
T^\mu &\equiv n_2^{-1/2} \sum_{j \in \mathcal{S}_2} \{\mu_1(X_j, U_j) - \mu_0(X_j, U_j) - \tau\}, \\
T^e &\equiv n_2^{-1/2} \sum_{j \in \mathcal{S}_2} (2A_j - 1)\{1 + M^{-1}K_{2,(X,U),j}\}\{Y_j - \mu_{A_j}(X_j, U_j)\}, \\
\mathcal{F}_{X,U} &\equiv \{(X_j, U_j) : j \in \mathcal{S}_2\}.
\end{aligned}
$$

We verify that the covariance between $T^\mu$ and $T^e$ is zero:

$$
\begin{aligned}
\text{cov}\left(T^\mu, T^e\right) &= E\left\{\text{cov}\left(T^\mu, T^e \mid \mathcal{F}_{X,U}\right)\right\} + \text{cov}\left\{E(T^\mu \mid \mathcal{F}_{X,U}), E(T^e \mid \mathcal{F}_{X,U})\right\} \\
&= E(0) + \text{cov}\left\{E(T^\mu \mid \mathcal{F}_{X,U}), 0\right\} \\
&= 0.
\end{aligned}
$$

Ignoring the $o_P(1)$ term, the variance of $n_2^{1/2}(\hat{\tau}_{\text{mat},2} - \tau)$ is

$$
\begin{aligned}
\text{var}&\left[n_2^{-1/2} \sum_{j \in \mathcal{S}_2} \left\{\mu_1(X_j, U_j) - \mu_0(X_j, U_j) - \tau\right\}\right] \\
&+ \text{var}\left[n_2^{-1/2} \sum_{j \in \mathcal{S}_2} (2A_j - 1)\left\{1 + M^{-1}K_{2,(X,U),j}\right\}\left\{Y_j - \mu_{A_j}(X_j, U_j)\right\}\right].
\end{aligned}
$$

As $n_2 \to \infty$, the first term becomes

$$
v_2^\tau = \text{var}\left\{\tau(X, U)\right\},
$$

and the second term becomes

$$
v_2^e = \text{plim}\left[n_2^{-1} \sum_{j \in \mathcal{S}_2} \left\{1 + M^{-1}K_{2,(X,U),j}\right\}^2 \sigma_{A_j}^2(X_j, U_j)\right].
$$

Under Assumption S1, $K_{2,(X,U),j} = O_P(1)$ and $E\{K_{2,(X,U),j}\}$ and $E\{K_{2,(X,U),j}^2\}$ are uniformly bounded over $n_2$ (Abadie and Imbens, 2006, Lemma 3). Therefore, a simple algebra yields $v_2^e = O(1)$. Combining all results, the asymptotic variance of $n_2^{1/2}(\hat{\tau}_{\text{mat},2} - \tau)$ is $v_2^\tau + v_2^e$.

The derivations for $\Gamma$ and $V$ are similar and thus omitted.

## S10 Proof of Theorem 3

For the matching estimators, Otsu and Rai (2016) showed that the distribution of $\hat{\tau}_2^* - \hat{\tau}_2$ given the observed data approximates the sampling distribution of $\hat{\tau}_2$. In what follows, we prove that $\widehat{\text{var}}(\hat{\tau}_2)$ is consistent for $\text{var}(\hat{\tau}_2)$, which covers both cases with RAL and matching estimators and is simpler than Otsu and Rai (2016) for the matching estimators.

Let $(M_1, \ldots, M_{n_2})$ be a multinomial random vector with $n_2$ draws on $n_2$ cells with equal probabilities. Let $W_j^* = n_2^{-1/2}M_j$ for $j = 1, \ldots, n_2$, and $\bar{W}^* = n_2^{-1}\sum_{j \in \mathcal{S}_2} W_j^*$. Then, the bootstrap weights $\{W_j^* : j = 1, \ldots, n_2\}$ satisfy that as $n_2 \to \infty$,

$$
\max_{j \in \mathcal{S}_2} |W_j^* - \bar{W}^*| \xrightarrow{P} 0, \tag{S10}
$$

$$
\sum_{j \in \mathcal{S}_2} (W_j^* - \bar{W}^*)^2 \xrightarrow{P} 1. \tag{S11}
$$

See, e.g., Mason and Newton (1992). The bootstrap replicate of $(\hat{\tau}_2 - \tau)$ can be written as

$$
\hat{\tau}_2^* - \hat{\tau}_2 = n_2^{-1/2} \sum_{j \in \mathcal{S}_2} \left(W_j^* - \bar{W}^*\right) \hat{\psi}_j. \tag{S12}
$$

For RAL estimators, following (S12), we have $\hat{\tau}_2^* - \hat{\tau}_2 = T_1^* + T_2^*$, where

$$T_1^* = n_2^{-1/2} \sum_{j \in \mathcal{S}_2} \left( W_j^* - \bar{W}^* \right) \psi(A_j, X_j, U_j, Y_j),$$

$$T_2^* = n_2^{-1/2} \sum_{j \in \mathcal{S}_2} \left( W_j^* - \bar{W}^* \right) \left\{ \hat{\psi}(A_j, X_j, U_j, Y_j) - \psi(A_j, X_j, U_j, Y_j) \right\}.$$

By (S10) and the fact that $\hat{\psi}(A_j, X_j, U_j, Y_j)$ is root-$n$ consistent for $\psi(A_j, X_j, U_j, Y_j)$, we have $T_2^* = o_P\left(n_2^{-1}\right)$. By (S11), we have

$$\hat{\tau}_2^* - \hat{\tau}_2 = T_1^* + o_P(1) = n_2^{-1/2} \sum_{j \in \mathcal{S}_2} \left( W_j^* - \bar{W}^* \right) \psi(A_j, X_j, U_j, Y_j) + o_P(1). \tag{S13}$$

For matching estimators, following (S12), we have

$$\hat{\tau}_2^* - \hat{\tau}_2 = n_2^{-1/2} \sum_{j \in \mathcal{S}_2} \left( W_j^* - \bar{W}^* \right) \hat{\psi}_{\mathrm{mat},j}$$

$$\equiv T_{\mathrm{mat},1}^* + T_{\mathrm{mat},2}^* + T_{\mathrm{mat},3}^*,$$

where

$$T_{\mathrm{mat},1}^* = n_2^{-1/2} \sum_{j \in \mathcal{S}_2} \left( W_j^* - \bar{W}^* \right) \psi_{\mathrm{mat},j},$$

$$T_{\mathrm{mat},2}^* = n_2^{-1/2} \sum_{j \in \mathcal{S}_2} \left( W_j^* - \bar{W}^* \right) \left[ \{ \hat{\mu}_1(X_j, U_j) - \hat{\mu}_0(X_j, U_j) - \hat{\tau}_2 \} - \{ \mu_1(X_j, U_j) - \mu_0(X_j, U_j) - \tau \} \right],$$

$$T_{\mathrm{mat},3}^* = n_2^{-1/2} \sum_{j \in \mathcal{S}_2} \left( W_j^* - \bar{W}^* \right) (2A_j - 1) \left\{ 1 + M^{-1} K_{2,(X,U),j} \right\} \left\{ \mu_{A_j}(X_j, U_j) - \hat{\mu}_{A_j}(X_j, U_j) \right\}.$$

Under Assumption S1, for any $\zeta > 0$, $E\{K_{2,(X,U),j}^{\zeta}\}$ is uniformly bounded over $n_2$ (Abadie and Imbens, 2006). Together with Assumption S2 and the property of the bootstrap weights that $\max_{j \in \mathcal{S}_2} |W_j^* - \bar{W}^*| \to 0$ in probability, as $n_2 \to \infty$, we have $T_{\mathrm{mat},2}^* = o_P\left\{ n_2^{-1+1/\dim(x,u)} \right\}$ and $T_{\mathrm{mat},3}^* = o_P\left\{ n_2^{-1+1/\dim(x,u)} \right\}$. By (S11), we have

$$\hat{\tau}_2^* - \hat{\tau}_2 = T_{\mathrm{mat},1}^* + o_P(1) = n_2^{-1/2} \sum_{j \in \mathcal{S}_2} \left( W_j^* - \bar{W}^* \right) \psi_{\mathrm{mat},j} + o_P(1). \tag{S14}$$

Unifying (S13) and (S14), $\hat{\tau}_2^* - \hat{\tau}_2 = n_2^{-1/2} \sum_{j \in \mathcal{S}_2} \left( W_j^* - \bar{W}^* \right) \psi_j + o_P(1)$. Let $\psi_{(i)}$ be the $i$th order statistic of $\{ \psi_j : j \in \mathcal{S}_2 \}$. Because $E\left( \psi_j^{\gamma} \right) < \infty$ for $0 \leq \gamma \leq 4$, we have

$$\frac{|\hat{\tau}_2^* - \hat{\tau}_2|}{n_2^{1/\gamma}} \leq 2 \frac{|\psi_{(1)}| + |\psi_{(n_2)}|}{n_2^{1/\gamma}} \to 0,$$

almost surely, as $n_2 \to \infty$, leading to $\max_{\{W_j^* : j \in \mathcal{S}_2\}} |\hat{\tau}_2^* - \hat{\tau}_2|/n_2^{1/\gamma} \to 0$, almost surely, as $n_2 \to \infty$, where the maximum is taken over all possible bootstrap replicates. By Theorem 3.8 of Shao and Tu (2012),

$$\frac{\widehat{\mathrm{var}}(\hat{\tau}_2)}{\mathrm{var}(\hat{\tau}_2)} \to 1,$$

36

almost surely, as $n_2 \to \infty$. This proves that $\widehat{\mathrm{var}}(\hat\tau_2)$ is consistent for $\mathrm{var}(\hat\tau_2)$.

The proofs for the consistency of $\hat\Gamma$ and $\hat V$ for $\Gamma$ and $V$ are similar and thus omitted. Therefore, $\widehat{\mathrm{var}}(\hat\tau)$ is consistent for $\mathrm{var}(\hat\tau)$.

## S11   Derivation of the optimal sample allocation

The goal is to minimize (20) subject to the constraint (19). By the Lagrange multipliers technique, it suffices to find the minimizer of

$$\mathcal{L}(n_1, n_2, \lambda) = n_2^{-1}v_2 - (n_2^{-1} - n_1^{-1})\gamma - \lambda(C - n_1 C_1 - n_2 C_2).$$

Solving $\partial\mathcal{L}(n_1, n_2, \lambda)/\partial n_1 = 0$ and $\partial\mathcal{L}(n_1, n_2, \lambda)/\partial n_2 = 0$, we have

$$n_1 = \left(\frac{\gamma}{\lambda C_1}\right)^{1/2}, \qquad n_2 = \left(\frac{v_2 - \gamma}{\lambda C_2}\right)^{1/2},$$

and therefore, (21) follows.

## S12   Proof of Theorem 4

Following the discussion in Section 7.1, we can express

$$n_2^{1/2}(\hat\tau_2 - \tau) = \rho^{1/2} n_1^{-1/2} \sum_{j \in \mathcal{S}_1} \pi_j^{-1} I_j \psi(A_j, X_j, U_j, Y_j) + o_P(1),$$

$$n_2^{1/2}(\hat\tau_{2,\mathrm{ep}} - \hat\tau_{1,\mathrm{ep}}) = \rho^{1/2} n_1^{-1/2} \sum_{j \in \mathcal{S}_1} (\pi_j^{-1} I_j - 1)\phi(A_j, X_j, Y_j) + o_P(1).$$

Then, the asymptotic normality of $\{n_2^{1/2}(\hat\tau_2 - \tau), n_2^{1/2}(\hat\tau_{2,\mathrm{ep}} - \hat\tau_{1,\mathrm{ep}})\}^{\mathrm{T}}$ follows from the multivariate central limit theorem. The corresponding asymptotic covariance matrix can be obtained by the following calculation. First, $v_2$ is

$$\mathrm{var}\left\{\rho^{1/2} n_1^{-1/2} \sum_{j \in \mathcal{S}_1} \pi_j^{-1} I_j \psi(A_j, X_j, U_j, Y_j)\right\}$$
$$= \rho \times \mathrm{var}\{\pi^{-1} I \psi(A, X, U, Y)\}.$$

Second, $\Gamma$ is

$$\mathrm{cov}\left\{\rho^{1/2} n_1^{-1/2} \sum_{j \in \mathcal{S}_1} \pi_j^{-1} I_j \psi(A_j, X_j, U_j, Y_j), \rho^{1/2} n_1^{-1/2} \sum_{j \in \mathcal{S}_1} (\pi_j^{-1} I_j - 1)\phi(A_j, X_j, Y_j)\right\}$$
$$= \rho \times \mathrm{cov}\{\pi^{-1} I \psi(A, X, U, Y), (\pi^{-1} I - 1)\phi(A, X, Y)\}.$$

Third, $V$ is

$$\mathrm{var}\left\{\rho^{1/2} n_1^{-1/2} \sum_{j \in \mathcal{S}_1} (\pi_j^{-1} I_j - 1)\phi(A_j, X_j, Y_j)\right\}$$
$$= \rho \times \mathrm{var}\{(\pi^{-1} I - 1)\phi(A, X, Y)\}.$$

## S13  Error-prone matching estimators in Section 7.2

Let the error-prone matching estimators be

$$
\hat{\tau}^{(0)}_{\mathrm{mat},1,\mathrm{ep}} \;=\; n_1^{-1} \sum_{i\in\mathcal{S}_1}(2A_i-1)\left(Y_i - \frac{1}{M}\sum_{l\in\mathcal{J}_{1,X,i}} Y_l\right),
$$

$$
\hat{\tau}^{(0)}_{\mathrm{mat},2,\mathrm{ep}} \;=\; \left(\sum_{j\in\mathcal{S}_2}\pi_j^{-1}\right)^{-1}\sum_{j\in\mathcal{S}_2}\pi_j^{-1}(2A_j-1)\left(Y_j - \frac{1}{M}\sum_{l\in\mathcal{J}_{2,X,j}} Y_l\right).
$$

Let the bias-corrected matching estimators be

$$
\hat{\tau}_{\mathrm{mat},d,\mathrm{ep}} = \hat{\tau}^{(0)}_{\mathrm{mat},d,\mathrm{ep}} - n_1^{-1/2}\hat{B}_{d,\mathrm{ep}}, \quad (d=1,2),
$$

where

$$
\hat{B}_{d,\mathrm{ep}} = n_1^{-1/2}\sum_{j\in\mathcal{S}_2}\pi_j^{-1}(2A_j-1)\left[\frac{1}{M}\sum_{l\in\mathcal{J}_{d,X,j}}\left\{\hat{\mu}_{1-A_j}(X_j)-\hat{\mu}_{1-A_j}(X_l)\right\}\right].
$$

## S14  Proof of Theorem 5

First, with the new definition (26), we can express

$$
\begin{aligned}
&n_2^{1/2}(\hat{\tau}_{\mathrm{mat},2}-\tau)\\
=\;&n_2^{1/2}(\hat{\tau}^{(0)}_{\mathrm{mat},2}-n_1^{-1/2}\hat{B}_2-\tau)\\
=\;&n_2^{1/2}\left\{\left(\sum_{j\in\mathcal{S}_2}\pi_j^{-1}\right)^{-1}\sum_{j\in\mathcal{S}_2}\pi_j^{-1}(2A_j-1)\left\{1+M^{-1}K_{2,(X,U),j}\right\}Y_j-\tau\right\}\\
&-n_2^{1/2}n_1^{-1}\sum_{j\in\mathcal{S}_2}\pi_j^{-1}(2A_j-1)\left[M^{-1}\sum_{l\in\mathcal{J}_{2,(X,U),j}}\left\{\hat{\mu}_{1-A_j}(X_j,U_j)-\hat{\mu}_{1-A_j}(X_l,U_l)\right\}\right]\\
=\;&\rho^{1/2}n_1^{-1/2}\sum_{j\in\mathcal{S}_2}\pi_j^{-1}\left\{\mu_1(X_j,U_j)-\mu_0(X_j,U_j)-\tau\right\}\\
&+\rho^{1/2}n_1^{-1/2}\sum_{j\in\mathcal{S}_2}\pi_j^{-1}(2A_j-1)\left\{1+M^{-1}K_{2,(X,U),j}\right\}\left\{Y_j-\mu_{A_j}(X_j,U_j)\right\}+o_P(1),
\end{aligned}
$$

where the third equality follows by some algebra.

Second, let

$$
\begin{aligned}
T^\mu &\equiv n_1^{-1/2}\sum_{j\in\mathcal{S}_2}\pi_j^{-1}\{\mu_1(X_j,U_j)-\mu_0(X_j,U_j)-\tau\},\\
T^e &\equiv n_1^{-1/2}\sum_{j\in\mathcal{S}_2}\pi_j^{-1}(2A_j-1)\{1+M^{-1}K_{2,(X,U),j}\}\{Y_j-\mu_{A_j}(X_j,U_j)\},\\
\mathcal{F}_{X,U} &\equiv \{(X_j,U_j):j\in\mathcal{S}_2\}.
\end{aligned}
$$

We verify that the covariance between $T^\mu$ and $T^e$ is zero:

$$
\begin{aligned}
\operatorname{cov}\left(T^\mu, T^e\right) &= E\left\{\operatorname{cov}\left(T^\mu, T^e \mid \mathcal{F}_{X,U}\right)\right\} + \operatorname{cov}\left\{E(T^\mu \mid \mathcal{F}_{X,U}), E(T^e \mid \mathcal{F}_{X,U})\right\} \\
&= E(0) + \operatorname{cov}\left\{E(T^\mu \mid \mathcal{F}_{X,U}), 0\right\} \\
&= 0.
\end{aligned}
$$

Ignoring the $o_P(1)$ term, the variance of $n_2^{1/2}(\hat{\tau}_{\mathrm{mat},2} - \tau)$ is

$$
\rho \times \operatorname{var}\left[n_1^{-1/2} \sum_{j \in \mathcal{S}_2} \pi_j^{-1}\left\{\mu_1(X_j, U_j) - \mu_0(X_j, U_j) - \tau\right\}\right]
$$

$$
+ \rho \times \operatorname{var}\left[n_1^{-1/2} \sum_{j \in \mathcal{S}_2} \pi_j^{-1}(2A_j - 1)\left\{1 + M^{-1}K_{2,(X,U),j}\right\}\left\{Y_j - \mu_{A_j}(X_j, U_j)\right\}\right].
$$

As $n_2 \to \infty$, the first term becomes

$$
\tilde{v}_2^\tau = \rho \times E\left[\frac{1 - \pi}{\pi}\left\{\tau(X, U) - \tau\right\}^2\right],
$$

and the second term becomes

$$
\tilde{v}_2^e = \rho \times \operatorname{plim}\left[n_1^{-1} \sum_{j \in \mathcal{S}_1} \frac{1 - \pi_j}{\pi_j}\left\{1 + M^{-1}K_{2,(X,U),j}\right\}^2 \sigma_{A_j}^2(X_j, U_j)\right].
$$

Because $\pi_j$'s are bounded away from zero, under Assumption S1, we have $K_{2,(X,U),j} = O_P(1)$, and $E\{K_{2,(X,U),j}\}$ and $E\{K_{2,(X,U),j}^2\}$ are uniformly bounded over $n_2$ (Abadie and Imbens, 2006, Lemma 3). Therefore, a simple algebra yields $\tilde{v}_2^e = O(1)$. Combining all results, the asymptotic variance of $n_2^{1/2}(\hat{\tau}_{\mathrm{mat},2} - \tau)$ is $\tilde{v}_2^\tau + \tilde{v}_2^e$.

The derivations for $\Gamma$ and $V$ are similar and thus omitted.

## S15   Proof of (23)

We write $\hat{\tau}_2 = \hat{\tau}_2(\hat{\alpha}, \hat{\beta}_0, \hat{\beta}_1)$ and $\hat{\tau}_{\mathrm{AIPW},2,j} = \hat{\tau}_{\mathrm{AIPW},2,j}(\hat{\alpha}, \hat{\beta}_0, \hat{\beta}_1)$ to emphasize its dependence on the parameter estimates $(\hat{\alpha}, \hat{\beta}_0, \hat{\beta}_1)$.

First, we write

$$
\begin{aligned}
\hat{\tau}_2(\alpha^*, \beta_0^*, \beta_1^*) &= \frac{\sum_{j \in \mathcal{S}_2} \pi_j^{-1}\hat{\tau}_{\mathrm{AIPW},2,j}(\alpha^*, \beta_0^*, \beta_1^*)}{\sum_{j \in \mathcal{S}_2} \pi_j^{-1}} \\
&= \frac{n_1^{-1}\sum_{j \in \mathcal{S}_1} \pi_j^{-1}I_j\hat{\tau}_{\mathrm{AIPW},2,j}(\alpha^*, \beta_0^*, \beta_1^*)}{n_1^{-1}\sum_{j \in \mathcal{S}_1} \pi_j^{-1}I_j} \equiv \frac{\hat{T}_1}{\hat{T}_2},
\end{aligned}
$$

where $\hat{T}_1 \equiv n_1^{-1}\sum_{j \in \mathcal{S}_1} \pi_j^{-1}I_j\hat{\tau}_{\mathrm{AIPW},2,j}(\alpha^*, \beta_0^*, \beta_1^*)$ and $\hat{T}_2 \equiv n_1^{-1}\sum_{j \in \mathcal{S}_1} \pi_j^{-1}I_j$. Then, $\hat{T}_1$ is

consistent for $T_1 \equiv \tau$, and $\hat{T}_2$ is consistent for $T_2 \equiv 1$. By the Taylor expansion, we have

$$
\begin{aligned}
\hat{\tau}_2(\alpha^*, \beta_0^*, \beta_1^*) - \tau &\cong \frac{T_1}{T_2} + \frac{1}{T_2}(\hat{T}_1 - T_1) - \frac{T_1}{T_2^2}(\hat{T}_2 - T_2) - \tau, \\
&= \left\{ n_1^{-1} \sum_{j \in \mathcal{S}_1} \pi_j^{-1} I_j \hat{\tau}_{\mathrm{AIPW},2,j}(\alpha^*, \beta_0^*, \beta_1^*) - \tau \right\} - \tau \left( n_1^{-1} \sum_{j \in \mathcal{S}_1} \pi_j^{-1} I_j - 1 \right) \\
&= n_1^{-1} \sum_{j \in \mathcal{S}_1} \pi_j^{-1} I_j \left\{ \hat{\tau}_{\mathrm{AIPW},2,j}(\alpha^*, \beta_0^*, \beta_1^*) - \tau \right\}.
\end{aligned}
$$

Second, by the Taylor expansion,

$$
\begin{aligned}
0 &= n_1^{-1} \sum_{j \in \mathcal{S}_2} \pi_j^{-1} S(A_j, X_j, U_j; \hat{\alpha}) \\
&\cong n_1^{-1} \sum_{j \in \mathcal{S}_2} \pi_j^{-1} S(A_j, X_j, U_j; \alpha^*) + \left\{ n_1^{-1} \sum_{j \in \mathcal{S}_2} \pi_j^{-1} \frac{\partial S(A_j, X_j, U_j; \alpha^*)}{\partial \alpha^{\mathrm{T}}} \right\} (\hat{\alpha} - \alpha^*) \\
&\cong n_1^{-1} \sum_{j \in \mathcal{S}_2} \pi_j^{-1} S(A_j, X_j, U_j; \alpha^*) + E \left\{ \frac{\partial S(A, X, U; \alpha^*)}{\partial \alpha^{\mathrm{T}}} \right\} (\hat{\alpha} - \alpha^*),
\end{aligned}
$$

where the last line follows because by Assumption 6, $n_1^{-1} \sum_{j \in \mathcal{S}_2} \pi_j^{-1} \partial S(A_j, X_j, U_j; \alpha^*) / \partial \alpha^{\mathrm{T}}$ converges to $E \{ \partial S(A, X, U; \alpha^*) / \partial \alpha^{\mathrm{T}} \}$ in probability. Therefore, we can establish

$$
\hat{\alpha} - \alpha^* \cong n_1^{-1} \sum_{j \in \mathcal{S}_2} E \left\{ \frac{\partial S(A, X, U; \alpha^*)}{\partial \alpha^{\mathrm{T}}} \right\}^{-1} \pi_j^{-1} S(A_j, X_j, U_j; \alpha^*). \tag{S15}
$$

Similarly, we can establish

$$
\hat{\beta}_a - \beta_a^* \cong n_1^{-1} \sum_{j \in \mathcal{S}_2} E \left\{ \frac{\partial S_a(A, X, U, Y; \beta_a^*)}{\partial \beta_a^{\mathrm{T}}} \right\}^{-1} \pi_j^{-1} S_a(A_j, X_j, U_j, Y_j; \beta_a^*) \quad (a = 0, 1). \tag{S16}
$$

Third, by the Taylor expansion and (S15) and (S16), we obtain

$$
\begin{aligned}
&\hat{\tau}_2(\hat{\alpha}, \hat{\beta}_0, \hat{\beta}_1) \\
\cong\ &\hat{\tau}_2(\alpha^*, \beta_0^*, \beta_1^*) + E \left\{ \frac{\partial \hat{\tau}_2(\alpha^*, \beta_0^*, \beta_1^*)}{\partial \alpha^{\mathrm{T}}} \right\} (\hat{\alpha} - \alpha^*) \\
&+ E \left\{ \frac{\partial \hat{\tau}_2(\alpha^*, \beta_0^*, \beta_1^*)}{\partial \beta_0^{\mathrm{T}}} \right\} (\hat{\beta}_0 - \beta_0^*) + E \left\{ \frac{\partial \hat{\tau}_2(\alpha^*, \beta_0^*, \beta_1^*)}{\partial \beta_1^{\mathrm{T}}} \right\} (\hat{\beta}_1 - \beta_1^*) \\
\cong\ &\hat{\tau}_2(\alpha^*, \beta_0^*, \beta_1^*) \\
&+ n_1^{-1} \sum_{j \in \mathcal{S}_1} E \left\{ \frac{\partial \hat{\tau}_2(\alpha^*, \beta_0^*, \beta_1^*)}{\partial \alpha^{\mathrm{T}}} \right\} E \left\{ \frac{\partial S(A, X, U; \alpha^*)}{\partial \alpha^{\mathrm{T}}} \right\}^{-1} \pi_j^{-1} I_j S(A_j, X_j, U_j; \alpha^*) \\
&- n_1^{-1} \sum_{j \in \mathcal{S}_1} E \left\{ \frac{\partial \hat{\tau}_2(\alpha^*, \beta_0^*, \beta_1^*)}{\partial \beta_0^{\mathrm{T}}} \right\} E \left\{ \frac{\partial S_0(A, X, U, Y; \beta_0^*)}{\partial \beta_0^{\mathrm{T}}} \right\}^{-1} \pi_j^{-1} I_j S_0(A_j, X_j, U_j, Y_j; \beta_0^*) \\
&- n_1^{-1} \sum_{j \in \mathcal{S}_1} E \left\{ \frac{\partial \hat{\tau}_2(\alpha^*, \beta_0^*, \beta_1^*)}{\partial \beta_1^{\mathrm{T}}} \right\} E \left\{ \frac{\partial S_1(A, X, U, Y; \beta_0^*)}{\partial \beta_1^{\mathrm{T}}} \right\}^{-1} \pi_j^{-1} I_j S_1(A_j, X_j, U_j, Y_j; \beta_1^*).
\end{aligned}
\tag{S17}
$$

Moreover, we have the following calculations:

$$E\left\{\frac{\partial\hat{\tau}_2(\alpha^*,\beta_0^*,\beta_1^*)}{\partial\alpha}\right\} = H_{\text{AIPW}},$$

$$E\left\{\frac{\partial\hat{\tau}_2(\alpha^*,\beta_0^*,\beta_1^*)}{\partial\beta_0}\right\} = -E\left[\left\{1-\frac{1-A}{1-e(X,U;\alpha^*)}\right\}\frac{\partial\mu_0(X;\beta_0^*)}{\partial\beta_0}\right],$$

$$E\left\{\frac{\partial\hat{\tau}_2(\alpha^*,\beta_0^*,\beta_1^*)}{\partial\beta_1}\right\} = E\left[\left\{1-\frac{A}{e(X,U;\alpha^*)}\right\}\frac{\partial\mu_1(X,U;\beta_1^*)}{\partial\beta_1}\right].$$

Under Assumption 5,

$$E\left\{\frac{\partial S(A,X,U;\alpha^*)}{\partial\alpha^{\mathrm{T}}}\right\} = E\left\{S(A,X,U;\alpha^*)^{\otimes2}\right\} = \Sigma_{\alpha\alpha}.$$

Under Assumption 4, we can still replace $E\left\{\partial S(A,X,U;\alpha^*)/\partial\alpha^{\mathrm{T}}\right\}$ by $\Sigma_{\alpha\alpha}$, because $H_{\text{AIPW}} = 0$ and the term in (S17) is zero. Therefore, combining the above results, we have $\hat{\tau}_2 - \tau = \hat{\tau}_2(\hat{\alpha},\hat{\beta}_0,\hat{\beta}_1) - \tau = n_1^{-1}\sum_{j\in\mathcal{S}_1}\pi_j^{-1}I_j\psi(A_j,X_j,U_j,Y_j)$, where $\psi(A,X,U,Y)$ is given by (11). The result follows.

## S16  Proof of Proposition 2

Based on (23) and (25) for RAL estimators or (27) and (28) for the matching estimators, we have

$$\begin{aligned}\hat{\tau}-\tau &= \hat{\tau}_2-\tau-\hat{\Gamma}^{\mathrm{T}}\hat{V}^{-1}(\hat{\tau}_{2,\text{ep}}-\hat{\tau}_{1,\text{ep}}) \\ &\cong n_1^{-1}\sum_{i\in\mathcal{S}_1}\pi_i^{-1}I_i\psi_i - \Gamma^{\mathrm{T}}V^{-1}\left(n_1^{-1}\sum_{i\in\mathcal{S}_1}\pi_i^{-1}I_i\phi_i - n_1^{-1}\sum_{i\in\mathcal{S}_1}\phi_i\right) \\ &\cong n_1^{-1}\sum_{i\in\mathcal{S}_1}\left\{\frac{I_i}{\pi_i}\psi_i - \left(\frac{I_i}{\pi_i}-1\right)\Gamma V^{-1}\phi_i\right\}.\end{aligned}$$

## S17  Proof of Theorem 6

Let $(M_1,\ldots,M_{n_1})$ be a multinomial random vector with $n_1$ draws on $n_1$ cells with equal probabilities. Let $W_i^* = n_1^{-1/2}M_i$ for $i=1,\ldots,n_1$, and $\bar{W}^* = n_1^{-1}\sum_{i\in\mathcal{S}_1}W_i^*$. Then, the bootstrap replicate of $(\hat{\tau}_2-\tau)$ can be written as

$$\begin{aligned}\hat{\tau}_2^*-\hat{\tau} &= n_1^{-1/2}\sum_{i\in\mathcal{S}_1}\left(W_i^*-\bar{W}^*\right)\pi_i^{-1}I_i\hat{\psi}_i \\ &\cong n_1^{-1/2}\sum_{i\in\mathcal{S}_1}\left(W_i^*-\bar{W}^*\right)\pi_i^{-1}I_i\psi_i,\end{aligned}$$

where the last line follows by a similar argument as in Section S10 for both RAL and matching estimators.

Let $l_i$ denote $\pi_i^{-1}I_i\psi_i$ and $l_{(i)}$ be the $i$th order statistic of $\{l_i : i\in\mathcal{S}_1\}$. Because $E\left(l_j^\gamma\right) < \infty$ for $0\le\gamma\le4$, we have

$$\frac{|\hat{\tau}_2^*-\hat{\tau}|}{n_1^{1/\gamma}} \le 2\frac{|l_{(1)}|+|l_{(n_1)}|}{n_1^{1/\gamma}} \to 0,$$

almost surely, as $n_1 \to \infty$, leading to $\max_{\{W_i^*:i \in \mathcal{S}_1\}} |\hat{\tau}_2^* - \hat{\tau}|/n_1^{1/\gamma} \to 0$, almost surely, as $n_1 \to \infty$, where the maximum is taken over all possible bootstrap replicates. By Theorem 3.8 of Shao and Tu (2012),

$$\frac{\widehat{\mathrm{var}}(\hat{\tau}_2)}{\mathrm{var}(\hat{\tau}_2)} \to 1,$$

almost surely, as $n_1 \to \infty$. This proves that $\widehat{\mathrm{var}}(\hat{\tau}_2)$ is consistent for $\mathrm{var}(\hat{\tau}_2)$.

The proofs for the consistency of $\hat{\Gamma}$ and $\hat{V}$ for $\Gamma$ and $V$ are similar and thus omitted. Therefore, $\widehat{\mathrm{var}}(\hat{\tau})$ is consistent for $\mathrm{var}(\hat{\tau})$.