# Robust Direct Learning for Causal Data Fusion

**Xinyu Li**[1,2]                                                    XINYU.LI@PKU.EDU.CN
**Yilin Li**[1]                                                         YILINLI@PKU.EDU.CN
**Qing Cui**[2]                                                  CUIQING.CQ@ANTGROUP.COM
**Longfei Li**[2]                                               LONGYAO.LLF@ANTGROUP.COM
**Jun Zhou**[2][*]                                                 JUN.ZHOUJUN@ANTFIN.COM
[1]*School of Mathematical Sciences, Peking University, Beijing, China*
[2]*Ant Group, Hangzhou, China*

**Editors:** Emtiyaz Khan and Mehmet Gönen

## Abstract

In the era of big data, the explosive growth of multi-source heterogeneous data offers many exciting challenges and opportunities for improving the inference of conditional average treatment effects. In this paper, we investigate homogeneous and heterogeneous causal data fusion problems under a general setting that allows for the presence of source-specific covariates. We provide a direct learning framework for integrating multi-source data that separates the treatment effect from other nuisance functions, and achieves double robustness against certain misspecification. To improve estimation precision and stability, we propose a causal information-aware weighting function motivated by theoretical insights from the semiparametric efficiency theory; it assigns larger weights to samples containing more causal information with high interpretability. We introduce a two-step algorithm, the weighted multi-source direct learner, based on constructing a pseudo-outcome and regressing it on covariates under a weighted least square criterion; it offers us a powerful tool for causal data fusion, enjoying the advantages of easy implementation, double robustness and model flexibility. In simulation studies, we demonstrate the effectiveness of our proposed methods in both homogeneous and heterogeneous causal data fusion scenarios.

**Keywords:** data fusion; direct learning; double robustness; heterogeneous treatment effect

## 1. Introduction

To understand the causal mechanism, a classic parameter of interest is the conditional average treatment effects (CATE), also known as heterogeneous treatment effect (Abrevaya et al., 2015; Athey and Imbens, 2016; Künzel et al., 2019), defined by the difference in outcome means between the two treatment groups conditional on a set of background attributes. Learning CATE is one of the fundamental problems in experimental sciences, observational studies, and electric commerce, such as the average effectiveness of medical treatment on patients (Obermeyer and Emanuel, 2016) and benefits of advertising on consumers (Bottou et al., 2013). A vast number of methods have been proposed to estimate the CATE based on flexible machine learning methods, including tree-based methods (Athey and Imbens, 2016; Tang et al., 2022), random forests (Wager and Athey, 2018), boosting (Powers et al., 2018; Nie and Wager, 2021), neural networks (Johansson et al., 2016; Louizos et al., 2017;

---

[*] Corresponding author.

Shi et al., 2019), and meta-learners with any supervised learning method (Künzel et al., 2019; Nie and Wager, 2021).

As pointed out by Kallus and Oprescu (2022), learning the conditional outcome means on two treatment arms separately and taking the difference to obtain the CATE may suffer from error accumulation, especially when the CATE function has a simpler and sparser form. This inspires us to learn the CATE function directly by modeling it as a whole. Qi and Liu (2018); Qi et al. (2020) proposed a one-step method (D-learning) to directly learn the optimal individual treatment rule which is closely related to the CATE. Meng and Qiao (2020) further generalized D-learning by replacing the outcome with the residual of some main effect function to achieve double robustness. The double robustness property is well studied in causal inference, meaning that the estimation is consistent if either the propensity score or the conditional outcome mean model is correct but not necessarily both, see Bang and Robins (2005); Zhang et al. (2012).

Although various methods have been proposed for estimating the CATE on a single dataset, the practical performance might be poor due to the limited sample size. A natural idea is to combine other similar datasets and improve the precision of the estimating procedures. Integrating and leveraging data from multiple sources have received wide attention in recent years. This problem is typically known as *causal data fusion*. Some notable advances focus on the average treatment effect (ATE) (Fan et al., 2014; Bareinboim and Pearl, 2016; Colnet et al., 2020; Li et al., 2021), causal discovery (Claassen and Heskes, 2010; Zhang et al., 2017) across multiple data sources. Recently, a tree-based approach for estimating the CATE is proposed when the individual-level data cannot be pooled (Tan et al., 2021).

To improve efficiency, we propose a novel approach for estimating the CATE on heterogeneous data sources by generalizing the approach from Meng and Qiao (2020). We have the following four concrete contributions: (i) We formulate and investigate the homogeneous and heterogeneous causal data fusions under general settings that allow for the presence of source-specific covariates. (ii) We present a multi-source direct learning framework, and propose a direct, model-flexible and doubly robust algorithm for causal data fusion. (iii) We propose a causal information-aware weighting function based on semiparametric efficiency theory to improve efficiency. (iv) We demonstrate the performance of the proposed methods and show the improvement compared with other methods.

## 2. Preliminaries and notations

### 2.1. Heterogeneous Treatment Effect

As is customary, we use capital letters for random variables and lowercase letters for realized values; in particular, let $Y$ denote the outcome of interest, $A \in \{1, -1\}$ denote a binary treatment indicator, $X \in \mathcal{X}$ denote a vector of features, and $Y(a)$ denote the potential outcome that would be observed when the treatment $A$ had been set to $a \in \{1, -1\}$ (Rubin, 1974; Imbens and Rubin, 2015). We maintain the classic stable unit treatment value assumption (SUTVA) that no interference between units and no hidden variations of treatments occur, and assume that the observed outcome is a realization of the potential outcome under the intervention actually received, i.e., $Y = \sum_a \mathbb{1}(A = a)Y(a)$, where $\mathbb{1}(\cdot)$ denotes the indicator function.

The CATE is defined as the mean difference in potential outcomes between the treated and control groups for individuals with feature $X$, that is, $\tau(X) = \mu_1(X) - \mu_{-1}(X)$, where $\mu_a(X) = E[Y(a) \mid X]$ for $a \in \{1, -1\}$; it captures the heterogeneity of average treatment effect across subpopulations defined by the values of the features.

We consider the following general model for potential outcomes, $Y(a) = \mu_a(X) + \tilde{\varepsilon}_a$, where $\tilde{\varepsilon}_a$ is a random error satisfying $E(\tilde{\varepsilon}_a \mid X) = 0$ and $\text{var}(\tilde{\varepsilon}_a) < \infty$ for $a \in \{1, -1\}$. Then because $Y = \sum_a \mathbb{1}(A = a)Y(a)$, the observed outcome can be naturally modeled by

$$
\begin{aligned}
Y &= \frac{\mu_1(X) + \mu_{-1}(X)}{2} + A\frac{\mu_1(X) - \mu_{-1}(X)}{2} + \varepsilon \\
&\triangleq m(X) + A\delta(X) + \varepsilon,
\end{aligned}
\tag{1}
$$

where $E(\varepsilon \mid X, A) = 0$, $\text{var}(\varepsilon) < \infty$. Hereafter we refer to $m(X)$ as the main effect function and $\delta(X)$ as the treatment effect function. Note that $\tau(X) = 2\delta(X)$, and thus estimating the CATE is equivalent to estimating the treatment effect function $\delta(X)$.

## 2.2. Multiple Heterogeneous Data Sources

Suppose that data are available from $K$ mutually independent data sources; in each of them, an individual either receives the experimental or the control intervention. When data across different sources are samples from the same global population, we refer to the data sources as *homogeneous*, and otherwise as *heterogeneous*; here we consider the latter, a more challenging case that allows for differences in the distribution of baseline features across sources. We introduce $S$ as the source indicator which takes values in $\mathcal{S} = \{1, \cdots, K\}$, and denote the dataset observed from the $s$-th data source as $\mathcal{O}_s$.

We consider the problem where observed covariates differ across multiple datasets; such scenarios are common in practice (e.g., Jia et al., 2006; Li et al., 2022). In particular, we denote the covariates of interest that are shared by all datasets as $X$, and the covariates of the $s$-th dataset other than $X$ as $Z_s$. For each data source $s$, the observed dataset $\mathcal{O}_s = \{(Y_i, A_i, X_i, Z_{s,i}), i = 1, \cdots, n_s\}$, in which the samples are independent and identically distributed according to $f(Y, A, X, Z_s \mid S = s)$, where $f(\cdot)$ denotes the density function. If $Z_s = \varnothing$ for each $s \in \mathcal{S}$, then it degenerates to the classic setting where all covariates are observed across data sources (Colnet et al., 2020). In the following, if not otherwise specified, we denote $E(\cdot)$ the expectation with respect to the joint distribution of observed data $\prod_{s=1}^{K} [f(Y, A, X, Z_s, S = s)]^{\mathbb{1}(S=s)}$, and $\hat{E}(\cdot)$ the empirical expectation.

We are interested in drawing inference about the causal effects conditional on the commonly shared covariates $X$, and the causal quantity of interest is the CATE function $\tau_s(X) = E\{Y(1) - Y(-1) \mid X, S = s\} = E[E\{Y(1) - Y(-1) \mid X, Z_s, S = s\} \mid X, S = s]$ of each data source. To identify the CATE, we impose the following regularity assumptions (Imbens and Rubin, 2015).

**Assumption 1 (Ignorability)** $Y(a) \perp A \mid X, Z_s, S = s$ *for each $s$.*

**Assumption 2 (Positivity)** *There exist a constant $c > 0$ such that $P(A = a \mid X, Z_s, S = s) > c$ almost surely for each $a$ and $s$.*

Assumption 1 excludes the unmeasured confounders between $A$ and $Y(a)$ within each data source, which is plausible as in many cases the assignment of interventions within each data source can be characterized by observed features. Assumption 2 requires that populations of the treated and control group have some overlap. Within each data source, under Assumptions 1–2, the CATE function $\tau_s(X)$ is identified, and a straightforward estimating method is to utilize data from the $s$-th source solely. However, it is appealing to integrate all data from different sources to enhance efficiency, especially when the sample size of each source is relatively small.

## 3. Direct Learning for Causal Data Fusion

### 3.1. Causal Data Fusion

We define the main and the treatment effect function of the source $s$ as

$$m_s(X, Z_s) = \frac{\mu_{1s}(X, Z_s) + \mu_{-1s}(X, Z_s)}{2} \text{ and } \delta_s(X, Z_s) = \frac{\mu_{1s}(X, Z_s) - \mu_{-1s}(X, Z_s)}{2},$$

respectively, where $\mu_{as}(X, Z_s) = E\{Y(a) \mid X, Z_s, S = s\}$. Following the same derivation as for Equation (1), in the scenario of multiple data sources, one can naturally model the observed outcome by

$$Y = m_S(X, Z_S) + A\delta_S(X, Z_S) + \varepsilon_S, \tag{2}$$

where $E(\varepsilon_s \mid X, Z_s, A, S = s) = 0$ and $\text{var}(\varepsilon_s) < \infty$ for $s = 1, \cdots, K$. When there is no confusion, we let $\delta_s(X) = E\{\delta_s(X, Z_s) \mid X, S = s\}$ and also refer to it as the treatment effect function. By definition it follows that $\delta_s(X) = \tau_s(X)/2$, hence estimating the CATE $\tau_s(X)$ is equivalent to estimating the treatment effect function $\delta_s(X)$.

Without loss of generality, we set the first data source as the target one, i.e., its samples are drawn from the population over which we are interested in inferring causal effects. For example, the first dataset $\mathcal{O}_1$ is collected from a randomized controlled trial (RCT) among the overall population of scientific interest, whereas the second dataset $\mathcal{O}_2$ is obtained from an observational study over a specific sub-population.

We classify causal data fusion tasks under multiple heterogeneous data sources into two categories:

- homogeneous causal data fusion, where the conditional average treatment effect across different data sources are homogeneous, see Assumption 3 described below;

- heterogeneous causal data fusion, where at least one of the data sources has a different conditional mean treatment effect than the others, and we will discuss it later.

**Assumption 3 (Homogeneity of CATE)** $\delta_1(X) = \cdots = \delta_K(X) \triangleq \delta(X)$ *almost surely.*

Assumption 3 reveals that for units with the same value of covariates $X$, the expected treatment effect remains the same regardless of the data source. It is plausible in many real cases, and weaker than the mean exchangeability or distribution exchangeability assumptions commonly adopted in causal inference literature when dealing with data fusion (e.g., Rudolph and van der Laan, 2017; Buchanan et al., 2018; Dahabreh et al., 2020; Li et al., 2021).

## 3.2. Direct Learning for Homogeneous Causal Data Fusion

We propose a direct learning approach to estimate the treatment effect function $\delta(X)$ without modeling each $\delta(X, Z_s)$. Unlike traditional methods that model CATE through modeling both $\mu_{1s}$ and $\mu_{-1s}$, the direct learner separates $m_s$ and $\delta$ in the estimation procedure, allowing the use of a flexible model for the quantity of interest CATE; for example, one can characterize CATE with a tree model with better interpretation, but fit the nuisance functions with a more complex and powerful model. Further, we also show the robustness of the direct learner against the failure of nuisance estimators.

In addition to main effect functions, the direct learners may also need to model a few other nuisance functions. We denote the treatment propensity scores $P(A = a \mid X, Z_s, S = s)$ by $p_{a|s}(X, Z_s)$, and selection propensity scores $P(S = s \mid X)$ by $\pi_s(X)$. For simplicity, when there causes no confusion, we may omit the random variables, e.g., $p_{a|s}$ for $p_{a|s}(X, Z_s)$. Note that the notations $P_{A|s} = \sum_a \mathrm{I}(A = a)P_{a|s}$ and $P_{A|S} = \sum_s \sum_a \mathrm{I}(S = s, A = a)P_{a|s}$.

We refer to working propensity score models, say $\widetilde{p}_{a|s}(X, Z_s)$, as satisfying the *partial balance* property with respect to $\delta_s(x, z_s)$, if

$$E\{e_A(X, Z_s)\delta(X, Z_s) \mid X, S = s\} = E\{\delta(X, Z_s) \mid X, S = s\},$$

where $e_A(X, Z_s) = \widetilde{p}_{A|s}^{-1}(X, Z_s)/E\{\widetilde{p}_{A|s}^{-1}(X, Z_s) \mid X, S = s\}$, a standardized inverse probability weighting term. We call this property *partial balance* because it characterizes the ability of the working propensity scores to adjust for the imbalance of $Z_s$ between the treated and control groups given $X$ along the direction of $\delta(X, Z_s)$. It is weaker than requiring the correctness of working propensity scores, because that calls for the ability to fully adjust for the imbalance of all covariates including both $X$ and $Z_s$. In order to train propensity score models that satisfy such partial balance property, one can first obtain an initial $\hat{\delta}(X, Z_s)$ using T-learner ([Künzel et al., 2019](#)) and then add the term $\hat{E}[g(X)\{P_{A|s}^{-1} - 1\}\hat{\delta} \mid S = s]^2$ to the loss function, where $g(X)$ is a vector of user-specified functions. The partial balance property is naturally satisfied in many cases, for example, (a) the working propensity scores are correct, i.e., $\widetilde{p}_{a|s}(X, Z_s) = p_{a|s}(X, Z_s)$; (b) there are no source-specific covariates, i.e., $Z_s = \varnothing$; (c) only the covariates $X$ lead to the heterogeneity of causal effects, i.e., $\delta(X, Z_s) = \delta(X)$. Case (c) holds especially when many variables predict potential outcomes but only a few has a strong modulate effect ([Kallus and Oprescu, 2022](#)).

The following result establishes the double robustness property of the direct learning approach for homogeneous causal data fusion, which is a generalization of the result in [Meng and Qiao (2020)](#). It offers us two chances to obtain a correct estimate of the CATE, thus more robust than methods that only rely on a single nuisance model.

**Theorem 4** *Suppose that Assumptions [1–3](#) hold and $\{w_s(X)\}_{s=1}^K$ are arbitrary positive integrable functions. Then*

$$\delta \in \underset{l \in \{\mathcal{X} \to \mathbb{R}\}}{\arg\min} E\left[\frac{w_S(X)}{\widetilde{p}_{A|S}(X, Z_S)}\left\{\frac{Y - \widetilde{m}_S(X, Z_S)}{A} - l(X)\right\}^2\right], \tag{3}$$

*if for each data source $s$, either one of the following conditions holds:*

*1. the working propensity scores $\widetilde{p}_{a|s}(X, Z_s) = p_{a|s}(X, Z_s)$, or*

*2. the working main effect functions $\widetilde{m}_s(X, Z_s) = m_s(X, Z_s)$, and the working propensity scores $\widetilde{p}_{a|s}(X, Z_s)$ satisfies the partial balance property with respect to $\delta(X, Z_s)$.*

Theorem 4 suggests that we can consistently estimate $\delta(X)$ through the empirical version of Equation (3) if the limiting functions of working nuisance models satisfy either one of the above conditions. This method avoids modeling the treatment effect function $\delta(X, Z_s)$ on each dataset. Instead, we integrate all the datasets to solve directly for the objective treatment effect function of $X$. At the same time, it allows us to separate the treatment effect from other variation independent nuisance functions, and achieve double robustness against misspecification of them.

**Remark 5** *In many cases, such as RCTs, the propensity scores are known and often the treatment assignment mechanism only depends on the commonly shared covariates $X$. With such prior knowledge, we simply need to fit propensity scores on $X$, and under the same condition as Theorem 4,*

$$\delta \in \underset{l \in \{\mathcal{X} \to \mathbb{R}\}}{\arg\min} E\left[\frac{w_S(X)}{\widetilde{p}_{A|S}(X)}\{Y - \widetilde{m}_S(X, Z_S) - Al(X)\}^2\right],$$

*if for each $s$, either $\widetilde{p}_{a|s}(X) = p_{a|s}(X)$ or $\widetilde{m}_s(X, Z_s) = m_s(X, Z_s)$ holds.*

### 3.3. Causal Information-aware Weighting Function

As shown in Theorem 4, the weight function $w_s(X)$ plays an important role in the direct learning; the choice of $w_s(X)$ affects the efficiency of estimators, but not the consistency. Users may choose a constant weight, but ideal weights should be interpretable as well as beneficial for improving estimation precision and stability. Towards this end, we propose a weighting function motivated by the semiparametric efficiency bound (also referred to as semiparametric Cramér-Rao bound, see Newey, 1990; Bickel et al., 1993; Tsiatis, 2006; Kennedy, 2016), which characterizes the amount of information contained in the observed data for inferring target parameters.

To illustrate, suppose all datasets share the same covariates and consider the case where covariates are discrete-valued for simplicity. At a fixed $X = \boldsymbol{x}$, in the statistical model $\mathcal{F}_{s,\boldsymbol{x}} = \{f(Y, A, X = \boldsymbol{x} \mid S = s) : \text{Assumptions 1–2 holds with no other restrictions}\}$, the amount of information for estimating CATE carried in each observation of the $s$-th dataset can be measured by the semiparametric efficiency bound of $\tau_s(\boldsymbol{x})$. We refer to the bound that corresponds to $s$ and $\boldsymbol{x}$ as $\mathcal{B}_{s,\boldsymbol{x}}$; the asymptotic variance of any regular and asymptotic linear estimator of $\tau_s(\boldsymbol{x})$ based solely on $\mathcal{O}_s$ can be no smaller than this bound. Then the relative information between $\mathcal{O}_{s_1}$ and $\mathcal{O}_{s_2}$ with respect to $\boldsymbol{x}$ is

$$\mathcal{B}_{s_1,\boldsymbol{x}}/\mathcal{B}_{s_2,\boldsymbol{x}} = \left\{\frac{V_{1|s_2}(\boldsymbol{x})}{p_{1|s_2}(\boldsymbol{x})} + \frac{V_{-1|s_2}(\boldsymbol{x})}{p_{-1|s_2}(\boldsymbol{x})}\right\}^{-1} \left\{\frac{V_{1|s_1}(\boldsymbol{x})}{p_{1|s_1}(\boldsymbol{x})} + \frac{V_{-1|s_1}(\boldsymbol{x})}{p_{-1|s_1}(\boldsymbol{x})}\right\},$$

where $s_1, s_2 \in \mathcal{S}$, $V_{a|s}(X) = \text{var}(Y_a \mid X, S = s) = \text{var}(Y \mid X, A = a, S = s)$ and $p_{a|s}(X) = P(A = a \mid X, S = s)$ for $a \in \{1, -1\}$.

In light of this, we propose an information-aware weighting function as

$$w_s(X) \propto \frac{P(S = s)}{P(S = 1)} \frac{\pi_1(X)}{\pi_s(X)} \left\{\frac{V_{1|s}(X)}{p_{1|s}(X)} + \frac{V_{-1|s}(X)}{p_{-1|s}(X)}\right\}^{-1} \tag{4}$$

$$= \frac{f(X \mid S = 1)}{f(X \mid S = s)} \left\{\frac{V_{1|s}(X)}{p_{1|s}(X)} + \frac{V_{-1|s}(X)}{p_{-1|s}(X)}\right\}^{-1}.$$

(a) $f(X \mid S = 1)$        (b) $f(X \mid S = 2)$

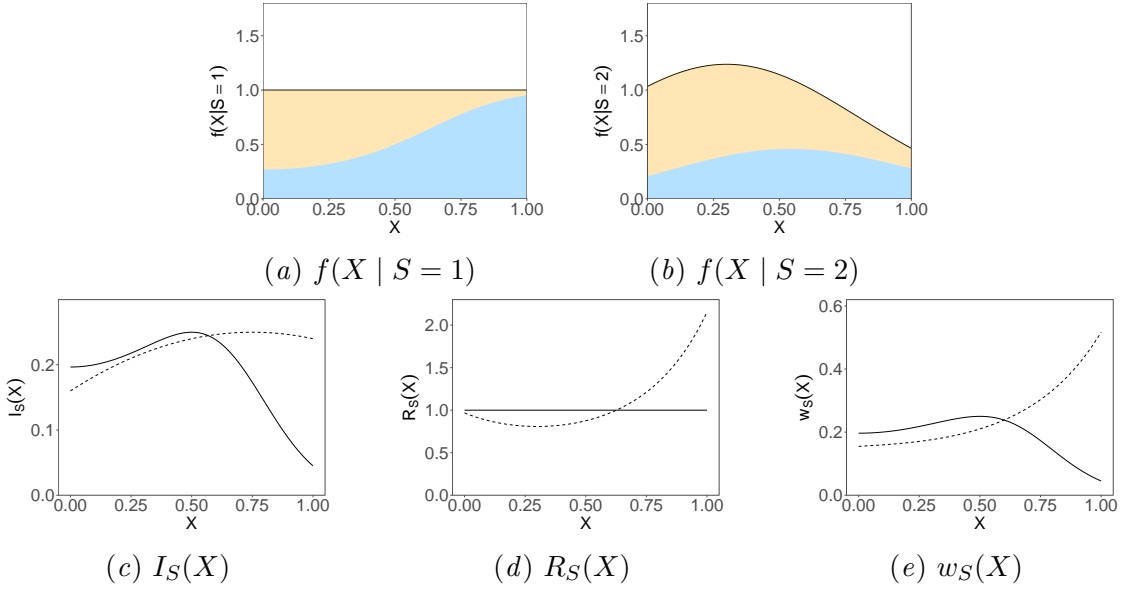(c) $I_S(X)$        (d) $R_S(X)$        (e) $w_S(X)$

Figure 1: Illustration of the weighting function. Suppose that $V_{a|s}(\boldsymbol{x})$ is a constant, and $f(X \mid S = s)$ is shown in (a) and (b). The blue and yellow areas indicate the proportions of the treated and control units at each value of $X$, respectively. (c), (d) and (e) show the variation of $I_s(X)$, $R_s(X)$ and $w_s(X) = I_s(X)R_s(X)$ with $X$. The solid and dashed lines stand for the first and second sources, respectively. Here $I_s(X)$ characterizes the imbalance of treatment, while $R_s(X)$ characterizes the imbalance of populations between two sources. When $X$ is close to one, $R_2(X) > R_1(X)$ which enables the second source population to match the first one, and $I_2(X) > I_1(X)$ because the first source suffers from a severe treatment imbalance; in this case, the weighting function $w_s(X)$ assigns more weight to the samples from the second source.

Our proposed weighting function $w_s(X)$ can be decomposed into the product of two components: the transfer term $R_s(X) = f(X \mid S = 1)/f(X \mid S = s)$ and the information term denoted by $I_s(X)$. The transfer term characterizes the imbalance of populations between the $s$-th source and the target one. The information term characterizes the amount of causal information satisfying that $I_{s_1}(\boldsymbol{x})/I_{s_2}(\boldsymbol{x}) = \mathcal{B}_{s_2,\boldsymbol{x}}/\mathcal{B}_{s_1,\boldsymbol{x}}$, therefore assigns larger weights to the samples containing more causal information. To match with intuition, for observations with covariates $X = \boldsymbol{x}$ in the $s$-th dataset, their weights are determined by:

- Density ratio corresponding to $f(\boldsymbol{x} \mid S = 1)/f(\boldsymbol{x} \mid S = s)$. This term enables us to pay more attention to the subpopulations that account for a larger share of the target population. Individuals that do not belong to the target population are discarded because their weights become zero. Within the target dataset $S = 1$, each observation is equally treated on this term.

- Degree of imbalance in treatment assignment corresponding to $p_{1|s}(\boldsymbol{x})$. The weighting function prefers to reward a more balanced treatment assignment mechanism. To illustrate, suppose that $V_{a|s}(\boldsymbol{x}) = 1$ and then the term $[1/p_{1|s}(\boldsymbol{x}) + 1/\{1 - p_{1|s}(\boldsymbol{x})\}]^{-1}$ degenerates to $p_{1|s}(\boldsymbol{x})\{1 - p_{1|s}(\boldsymbol{x})\}$, which is maximized at $p_{1|s}(\boldsymbol{x}) = 0.5$.

- Noise corresponding to $V_{1|s}(\boldsymbol{x})$ and $V_{-1|s}(\boldsymbol{x})$. The noise of potential outcomes varies across subpopulations and datasets, and our weighting function assigns larger weights to those with less noise in a subtle way that cooperates with the information of treatment assignment.

**Remark 6** *When $K = 1$ and $Z_1 = \varnothing$, the setting degenerates to the classic problem of estimating CATE on a single dataset. In this case, our proposed direct learning method can be regarded as a weighted version of Meng and Qiao (2020), utilizing a weighting function $w(X) \propto [V_1(X)/p_1(X) + V_{-1}(X)/\{1 - p_1(X)\}]^{-1}$. In the numerical experiments reported below, we show that the use of our proposed weight can significantly improve the performance of direct learning.*

## 3.4. Algorithm

We propose a two-step procedure called Weighted Multi-source Direct Learner (WMDL) for estimating treatment effect functions from heterogeneous multi-source data, which is summarized in Algorithm 1. As a first step, we construct a pseudo-outcome and a weight value for each data point by plugging in the estimators of nuisance functions. The nuisances are variation independent to CATE, including the main effect functions $m_s(X, Z_s)$, treatment propensity scores $p_{a|s}(X, Z_s)$ and weighting functions $w_s(X)$. By Equation (4), estimating weighting function is equivalent to estimating selection propensity scores $\pi_s(X)$, conditional treatment propensity scores $p_{a|s}(X)$ and conditional outcome variances $V_{a|s}(X)$. Alternatively, one can also choose to model the conditional density ratios instead of the selection propensity scores. The working models used to construct weighting functions only affect the efficiency, but not the consistency. To obtain nuisance estimators, we split each dataset into $G$ even folds, and estimates are fit on data excluding the fold where the data point lies. The data splitting technique is commonly used when learning nuisance estimators (Kallus and Oprescu, 2022).

---

**Algorithm 1** Weighted Multi-source Direct Learner

---

**Input:** Datasets $\mathcal{O}_s, s = 1, \cdots, K$, nuisance function learners, regression learner

**Output:** $\hat{\delta}(X)$

**for** $s = 1, \cdots, K$ **do**

    Use dataset $\mathcal{O}_s$ to construct nuisance estimators $\hat{m}_s$, $\hat{p}_{a|s}$ and $\hat{w}_s$

    Pseudo-outcome $\widetilde{Y}_i = A_i\{Y_i - \hat{m}_s(X_i, Z_{s,i})\}$ and $\widetilde{W}_i = \hat{w}_s(X_i)/\hat{p}_{A_i|s}(X_i, Z_{s,i})$

    Set $\mathcal{P}_s = \{\widetilde{Y}_i, \widetilde{W}_i, X_i\}_{i \in \mathcal{O}_s}$

**end**

Set $\mathcal{P} = \mathcal{P}_1 \cup \cdots \cup \mathcal{P}_K$

Fit regression learner on $\mathcal{P}$ to obtain $\hat{\delta}(X) = \operatorname{argmin}_l \sum_{i \in \mathcal{P}} \left[ \widetilde{W}_i \left\{ \widetilde{Y}_i - l(X_i) \right\}^2 \right]$

---

Then given a regression learner, for example, random forest, we regress the pseudo-outcomes on the covariates $X$ by solving a weighted least square problem. The regression learner is specified independently of those nuisance function learners, which affords us significant generality. The algorithm offers us a powerful tool for causal data fusion, enjoying the advantages of easy implementation, double robustness, and model flexibility.

### 3.5. Direct Learning for Heterogeneous Causal Data Fusion

In this part, we turn to investigate heterogeneous causal data fusion, where Assumption 3 does not hold anymore. We set the first data source $S = 1$ as the target one, and hence the causal quantity of interest is $\delta_1(X)$. We rewrite $\delta_S(X)$ as $\delta(X, S)$, then analogous to the Theorem 4, under Assumptions 1 and 2 one can obtain the function $\delta(X, S)$ by solving

$$\underset{l \in \{\mathcal{X} \times \mathcal{S} \to \mathbb{R}\}}{\arg\min} \ E \left[ \frac{w_S(X)}{p_{A|S}(X, Z_S)} \left\{ \frac{Y - m_S(X, Z_S)}{A} - l(X, S) \right\}^2 \right].$$

The inclusion of variable $S$ assists us in both distinguishing the heterogeneity of causal effects between different datasets and capturing the commonly shared structures. Therefore, for heterogeneous causal data fusion, we can still apply the proposed WMDL to estimate CATE in a direct, model-flexible, and robust way. To be specific, the nuisance estimators remains the same as the homogeneous causal data fusion setting, and the estimation procedure for CATE is similar to Algorithm 1, except that we turn to solve $\hat{\delta}(X, S) = \operatorname{argmin}_l \sum_{i \in \mathcal{P}} [\widetilde{W}_i \{\widetilde{Y}_i - l(X_i, S_i)\}^2]$ in the last step. We take $\hat{\delta}(x, 1)$ as the resulted estimator of the treatment effect function $\delta_1(x)$; it enjoys the property of double robustness against misspecification of nuisance models corresponding to $\mathcal{O}_1$. In the numerical experiments, we demonstrate that WMDL can effectively integrate information from multiple data sources, resulting in more accurate estimates than a single data source.

## 4. Experiments

### 4.1. Settings

In this part, we conduct a broad simulation comparing the weighted multi-source direct learner with other popular methods under both homogeneous and heterogeneous causal data fusions. We set the number of data sources $K = 10$, and draw the same number of independent observations for each data source from the following data generating process:

$$X \sim \text{Unif}\left([-1, 1]^4\right) \text{ for } S = 1, \quad X \sim N_{[-1,1]^4}(\mu_S, I) \text{ for } S = 2, \dots, K,$$

where $\mu_s \sim N(0, \sigma^2 I)$ with $\sigma = 0.3$, $N_{[-1,1]^4}$ stands for the four-dimensional truncated normal distribution and $I$ denotes the identity matrix. For each causal data fusion, we consider two scenarios: the Scenario I in which there are no source-specific covariates with $Z_s = \varnothing$, and the Scenario II in which each data source has a source-specific covariate $Z_s \sim N_{[-1,1]}(0, 1)$ for $s = 1, \dots, K$.

Under both homogeneous and heterogeneous causal data fusions, we generate the treatment $A$ in each dataset by Bernoulli$\left(\text{expit}\{(X, Z_s)^\mathrm{T} \beta\}\right)$, where $\text{expit}(\cdot) = \exp(\cdot)/\{1 + \exp(\cdot)\}$, and the parameter $\beta$ is randomly generated by the normal distribution $N(0, I)$.

Then we generate the outcome $Y$ by Equation (2); the main effect function $m_s$'s detailed form is provided in the supplementary material, the random noise $\varepsilon_s \sim N(0, \sigma_s^2)$ with $\sigma_s = 0.1$, and the treatment effect functions

$$
\begin{aligned}
\delta(X) =& (X_1 + X_2 + X3) \cdot \mathbb{1}(X_1 < 0.5) + X_4, \\
\delta_s(X) =& X_1 \cdot \mathbb{1}(X_1 < 0.5)\mathbb{1}(\mathcal{S}_1) + X_2 \cdot \mathbb{1}(X_2 < 0.5)\mathbb{1}(s \le 7) \\
& + X_3 \cdot \mathbb{1}(X_3 < 0.5) + X_4 \cdot \mathbb{1}(\mathcal{S}_1) + 2 \cdot \mathbb{1}(X_1 < 0)\mathbb{1}(\mathcal{S}_2),
\end{aligned}
$$

for homogeneous and heterogeneous causal data fusions, respectively, where $\mathcal{S}_1 = \{s : s \bmod 2 = 1\}$ and $\mathcal{S}_2 = \{s : s \bmod 2 = 0\}$.

We demonstrate the numerical performance by comparing it with the following widely-used methods. The meta-learner methods decompose the CATE into several regression and classification problems solved by any supervised learning method, such as T-learner (TL), S-learner (SL), and X-learner (XL), see Künzel et al. (2019) for details. Causal forest (CF) is another popular approach for estimating the CATE based on random forest (Wager and Athey, 2018). For each method, we take the following two strategies to integrate datasets as the benchmark: (i) directly combine all the data and then apply the method to learn CATE, and (ii) include the source indicator as an additional predictor. We add "$-s$" at the end of the abbreviation (e.g., XL-s) to indicate the inclusion of $S$. To assess the role of proposed weighting function $w_S(X)$, we also take the multi-source direct learner using $w_S(X) = 1$ in Algorithm 1 as a benchmark, referred to as MDL. Furthermore, to reveal the benefits of integrating multiple sources, we apply direct learners only on the target dataset, then MDL/WMDL degenerates to the robust direct learning (DL) by Meng and Qiao (2020) and the weighted robust direct learning (WDL). For the direct learners (WMDL, MDL, WDL and DL) and meta-learners (TL, SL, XL), we all use XGBoost for building the CATE and conditional mean models. We set hyperparameters for XGBoost as follows: the learning rate is 0.01, the maximum depth of a tree is 6 and the max number of boosting iterations is 20000. We evaluate by replicating 100 times, each time on an independent test dataset of sample size 1000 generated from the target source. We use the mean square error (MSE), an empirical version of the $L_2$ distance,

$$
\text{MSE} = n^{-1} \sum_{i=1}^{n} \left\{ \hat{\delta}_1(X_i) - \delta_1(X_i) \right\}^2,
$$

as the performance metric.

## 4.2. Results

We summarize the average of mean squared error and the corresponding standard deviation of the above experiments in Table 1.

From Table 1, we reach the following conclusions:

- Compared with the benchmarks, the proposed WMDL results in the smallest mean squared error on average in all scenarios, demonstrating its advantages in terms of precision under both homogeneous and heterogeneous causal data fusions. Also, MDL and WMDL typically have a smaller standard deviation than other methods, which implies that multi-source direct learning may lead to a more stable performance.

Table 1: Mean and standard deviation of MSE under sample sizes 3000 and 5000

| | Homogeneous | | | | Heterogeneous | | | |
|---|---|---|---|---|---|---|---|---|
| | I | | II | | I | | II | |
| | 3000 | 5000 | 3000 | 5000 | 3000 | 5000 | 3000 | 5000 |
| WMDL | **0.084** | **0.045** | **0.106** | **0.064** | **0.163** | **0.102** | **0.289** | **0.192** |
| | (0.008) | (0.003) | (0.009) | (0.006) | (0.035) | (0.014) | (0.066) | (0.030) |
| MDL | 0.104 | 0.061 | 0.143 | 0.089 | 0.198 | 0.121 | 0.380 | 0.251 |
| | (0.009) | (0.005) | (0.013) | (0.009) | (0.037) | (0.016) | (0.084) | (0.045) |
| WDL | 0.227 | 0.148 | 0.458 | 0.303 | 0.226 | 0.143 | 0.451 | 0.303 |
| | (0.060) | (0.026) | (0.096) | (0.051) | (0.053) | (0.020) | (0.089) | (0.049) |
| DL | 0.275 | 0.175 | 0.532 | 0.346 | 0.275 | 0.166 | 0.532 | 0.350 |
| | (0.064) | (0.030) | (0.112) | (0.057) | (0.058) | (0.023) | (0.094) | (0.059) |
| CF | 0.171 | 0.159 | 0.198 | 0.155 | 0.959 | 0.948 | 0.819 | 0.780 |
| | (0.011) | (0.009) | (0.014) | (0.010) | (0.038) | (0.039) | (0.043) | (0.038) |
| XL | 0.394 | 0.356 | 0.478 | 0.359 | 1.225 | 1.141 | 1.136 | 0.988 |
| | (0.024) | (0.029) | (0.039) | (0.026) | (0.049) | (0.056) | (0.077) | (0.054) |
| TL | 0.591 | 0.500 | 0.781 | 0.593 | 1.448 | 1.331 | 1.488 | 1.264 |
| | (0.039) | (0.037) | (0.060) | (0.039) | (0.071) | (0.066) | (0.093) | (0.067) |
| SL | 0.267 | 0.226 | 0.248 | 0.202 | 1.058 | 1.005 | 0.893 | 0.831 |
| | (0.028) | (0.023) | (0.023) | (0.020) | (0.044) | (0.047) | (0.055) | (0.042) |
| CF-s | 0.137 | 0.105 | 0.160 | 0.129 | 0.637 | 0.474 | 0.625 | 0.425 |
| | (0.009) | (0.006) | (0.010) | (0.009) | (0.067) | (0.047) | (0.102) | (0.042) |
| XL-s | 0.119 | 0.080 | 0.708 | 0.649 | 0.213 | 0.131 | 0.858 | 0.612 |
| | (0.028) | (0.016) | (0.215) | (0.168) | (0.038) | (0.020) | (0.280) | (0.191) |
| TL-s | 0.285 | 0.215 | 1.237 | 1.066 | 0.350 | 0.231 | 1.372 | 1.049 |
| | (0.074) | (0.047) | (0.284) | (0.209) | (0.071) | (0.042) | (0.343) | (0.195) |
| SL-s | 0.158 | 0.155 | 0.438 | 0.464 | 0.420 | 0.355 | 0.560 | 0.503 |
| | (0.019) | (0.016) | (0.194) | (0.155) | (0.062) | (0.047) | (0.170) | (0.134) |

- By comparing WMDL/WDL to MDL/DL, as expected, the use of our proposed information-aware weighting function significantly improves both the accuracy and stability of the CATE estimates, regardless of whether multiple sources or a single source is analyzed.

- The multi-source approach (MDL/WMDL) outperforms the single-source approach (DL/WDL) substantially, where the MSE decreases considerably in homogeneous cases and reduces to a certain degree in heterogeneous cases. This indicates that our proposed method can effectively integrate causal information from multiple sources, and may lead to efficiency gains even when the CATE differs across data sources.

- The inclusion of the source indicator as a covariate in CF, TL, SL and XL can generally promote the performance in the Scenario I, but may lead to larger MSEs in the Scenario II. In contrast, the WMDL behaves well across scenarios. It suggests that

our method can provide a simple but powerful way to make effective use of source-specific covariates, which only requires these covariates when modeling the nuisance functions but not the CATE.

Figure 2 presents the MSE of MDL and WMDL with increasing sample size. The WMDL outperforms the MDL under both homogeneous and heterogeneous settings, which highlights the effectiveness of our proposed causal information-aware function. The box plots in Figure 2 also show the variability of MD-learning can be reduced with the weighting function. These numerical results show the potential application of efficiency theory in Section 3.3 to CATE problems.
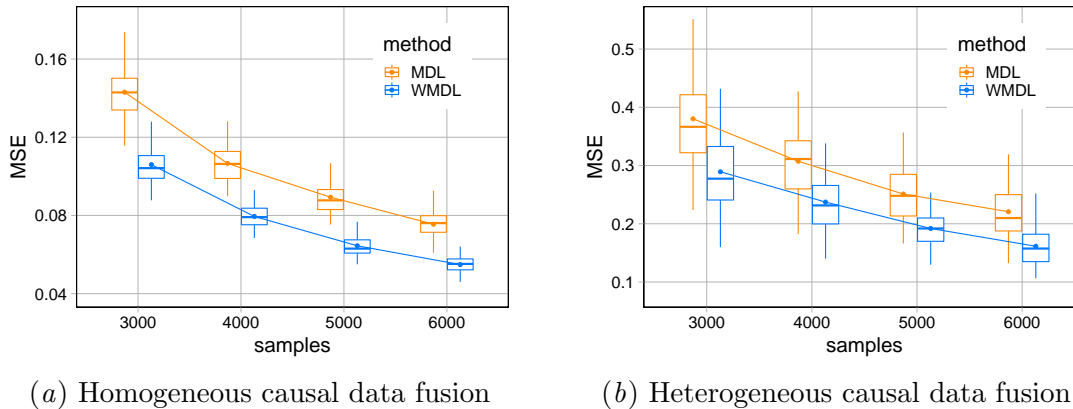


$(a)$ Homogeneous causal data fusion   $(b)$ Heterogeneous causal data fusion

Figure 2: Boxplots of mean square error for multi-source direct learning with and without the causal information-aware weighting function.

## 5. Discussion

We provide a simple yet powerful algorithm, the WMDL, for causal data fusion. In this paper, the target dataset contains the outcome data under different treatments. However, in some cases, only the covariates of the target population are available (Dahabreh et al., 2020). Therefore transferring the causal inference from other sources to the target population is also of interest. We highlight that our proposed WMDL can also be applied to *causal transfer learning*, see details in Appendix A. Our approach is double-weighted, with treatment propensity scores adjusting for bias and enhancing robustness, and causal information-aware weighting functions improving efficiency. In the future, we may consider adding the uncertainty of models to weighting functions to achieve more stable performance.

## Acknowledgements

The proof of Theorem 4 and more simulation details are available in the supplementary material.

## References

Jason Abrevaya, Yu-Chin Hsu, and Robert P Lieli. Estimating conditional average treatment effects. *Journal of Business & Economic Statistics*, 33(4):485–505, 2015.

Susan Athey and Guido Imbens. Recursive partitioning for heterogeneous causal effects. *Proceedings of the National Academy of Sciences*, 113(27):7353–7360, 2016.

Heejung Bang and James M Robins. Doubly robust estimation in missing data and causal inference models. *Biometrics*, 61(4):962–973, 2005.

Elias Bareinboim and Judea Pearl. Causal inference and the data-fusion problem. *Proceedings of the National Academy of Sciences*, 113(27):7345–7352, 2016.

Peter J Bickel, Chris A J Klaassen, Ya'acov Ritov, and Jon A Wellner. *Efficient and Adaptive Estimation for Semiparametric Models*. Baltimore: Johns Hopkins University Press, 1993.

Léon Bottou, Jonas Peters, Joaquin Quiñonero-Candela, Denis X Charles, D Max Chickering, Elon Portugaly, Dipankar Ray, Patrice Simard, and Ed Snelson. Counterfactual reasoning and learning systems: the example of computational advertising. *Journal of Machine Learning Research*, 14(11), 2013.

Ashley L Buchanan, Michael G Hudgens, Stephen R Cole, Katie R Mollan, Paul E Sax, Eric S Daar, Adaora A Adimora, Joseph J Eron, and Michael J Mugavero. Generalizing evidence from randomized trials using inverse probability of sampling weights. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 181(4):1193–1209, 2018.

Tom Claassen and Tom Heskes. Causal discovery in multiple models from different experiments. In *Advances in Neural Information Processing Systems*, volume 23, pages 415–423, 2010.

Bénédicte Colnet, Imke Mayer, Guanhua Chen, Awa Dieng, Ruohong Li, Gaël Varoquaux, Jean-Philippe Vert, Julie Josse, and Shu Yang. Causal inference methods for combining randomized trials and observational studies: a review. *arXiv preprint arXiv:2011.08047*, 2020.

Issa J Dahabreh, Lucia C Petito, Sarah E Robertson, Miguel A Hernán, and Jon A Steingrimsson. Towards causally interpretable meta-analysis: Transporting inferences from multiple randomized trials to a new target population. *Epidemiology*, 31(3):334, 2020.

Yanqin Fan, Robert Sherman, and Matthew Shum. Identifying treatment effects under data combination. *Econometrica*, 82(2):811–822, 2014.

Guido W Imbens and Donald B Rubin. *Causal Inference in Statistics, Social, and Biomedical Sciences*. New York: Cambridge University Press, 2015.

Jinzhu Jia, Zhi Geng, and Mingfeng Wang. Identifiability and estimation of probabilities from multiple databases with incomplete data and sampling selection. In *Joint IAPR International Workshops on Statistical Techniques in Pattern Recognition (SPR) and Structural and Syntactic Pattern Recognition (SSPR)*, pages 792–798, 2006.

Fredrik Johansson, Uri Shalit, and David Sontag. Learning representations for counterfactual inference. In *International Conference on Machine Learning*, pages 3020–3029, 2016.

Nathan Kallus and Miruna Oprescu. Robust and agnostic learning of conditional distributional treatment effects. *arXiv preprint arXiv:2205.11486*, 2022.

Edward H Kennedy. Semiparametric theory and empirical processes in causal inference. In *Statistical Causal Inferences and Their Applications in Public Health Research*, pages 141–167. 2016.

Sören R Künzel, Jasjeet S Sekhon, Peter J Bickel, and Bin Yu. Metalearners for estimating heterogeneous treatment effects using machine learning. *Proceedings of the National Academy of Sciences*, 116(10):4156–4165, 2019.

Hongkai Li, Jinzhu Jia, Ran Yan, Fuzhong Xue, and Zhi Geng. A causal data fusion method for the general exposure and outcome. *Statistics in Medicine*, 41(2):328–339, 2022.

Xinyu Li, Wang Miao, Fang Lu, and Xiao-Hua Zhou. Improving efficiency of inference in clinical trials with external control data. *Biometrics*, 2021. In press.

Christos Louizos, Uri Shalit, Joris M Mooij, David Sontag, Richard Zemel, and Max Welling. Causal effect inference with deep latent-variable models. In *Advances in Neural Information Processing Systems*, volume 30, pages 6446–6456, 2017.

Haomiao Meng and Xingye Qiao. Doubly robust direct learning for estimating conditional average treatment effect. *arXiv preprint arXiv:2004.10108*, 2020.

Whitney K Newey. Semiparametric efficiency bounds. *Journal of Applied Econometrics*, 5 (2):99–135, 1990.

Xinkun Nie and Stefan Wager. Quasi-oracle estimation of heterogeneous treatment effects. *Biometrika*, 108(2):299–319, 2021.

Ziad Obermeyer and Ezekiel J Emanuel. Predicting the future—big data, machine learning, and clinical medicine. *The New England Journal of Medicine*, 375(13):1216, 2016.

Scott Powers, Junyang Qian, Kenneth Jung, Alejandro Schuler, Nigam H Shah, Trevor Hastie, and Robert Tibshirani. Some methods for heterogeneous treatment effect estimation in high dimensions. *Statistics in Medicine*, 37(11):1767–1787, 2018.

Zhengling Qi and Yufeng Liu. D-learning to estimate optimal individual treatment rules. *Electronic Journal of Statistics*, 12(2):3601–3638, 2018.

Zhengling Qi, Dacheng Liu, Haoda Fu, and Yufeng Liu. Multi-armed angle-based direct learning for estimating optimal individualized treatment rules with various outcomes. *Journal of the American Statistical Association*, 115(530):678–691, 2020.

Donald B Rubin. Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of Educational Psychology*, 66(5):688, 1974.

Kara E Rudolph and Mark J van der Laan. Robust estimation of encouragement design intervention effects transported across sites. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 79(5):1509–1525, 2017.

Claudia Shi, David Blei, and Victor Veitch. Adapting neural networks for the estimation of treatment effects. In *Advances in Neural Information Processing Systems*, volume 32, pages 2503–2513, 2019.

Xiaoqing Tan, Chung-Chou H Chang, and Lu Tang. A tree-based federated learning approach for personalized treatment effect estimation from heterogeneous data sources. *arXiv preprint arXiv:2103.06261*, 2021.

Caizhi Tang, Huiyuan Wang, Xinyu Li, Qing Cui, Ya-Lin Zhang, Feng Zhu, Longfe Li, Jun Zhou, and Linbo Jiang. Debiased causal tree: heterogeneous treatment effects estimation with unmeasured confounding. In *Advances in Neural Information Processing Systems*, volume 36, 2022. In press.

Anastasios A Tsiatis. *Semiparametric Theory and Missing Data*. New York: Springer, 2006.

Stefan Wager and Susan Athey. Estimation and inference of heterogeneous treatment effects using random forests. *Journal of the American Statistical Association*, 113(523):1228–1242, 2018.

Baqun Zhang, Anastasios A Tsiatis, Eric B Laber, and Marie Davidian. A robust method for estimating optimal treatment regimes. *Biometrics*, 68(4):1010–1018, 2012.

Kun Zhang, Biwei Huang, Jiji Zhang, Clark Glymour, and Bernhard Schölkopf. Causal discovery from nonstationary/heterogeneous data: skeleton estimation and orientation determination. In *IJCAI: Proceedings of the Conference*, volume 2017, pages 1347–1353, 2017.

## Appendix A. Discussions on Causal Transfer Learning

Suppose we have individual-level data under different treatments from multiple sources as well as covariate data from a target population. We aim to synthesize findings across multiple observational datasets and transport causal inferences to the target population. Such scenarios are common in real applications, see Dahabreh et al. (2020). In this section, we provide a detailed discussion on how to apply our proposed framework to transfer causal inference.

To formalize the problem, we keep the notation consistent with the main text and let the source indicator $S = 0$ to indicate the target population. Individual-level data under different treatments are collected in $K$ mutually independent trial datasets, but for the target population, only covariate information is available. We summarize the observed datasets as $\mathcal{O}_0 = \{X_i, i = 1, \cdots, n_0\}$ and $\mathcal{O}_s = \{(Y_i, A_i, X_i, Z_{s,i}), i = 1, \cdots, n_s\}$ for $s = 1, \cdots, K$. Then the joint distribution of observed data is $f(X, S = 0)^{\mathbb{1}(S=0)} \prod_{s=1}^{K} [f(Y, A, X, Z_s, S = s)]^{\mathbb{1}(S=s)}$. Our parameter of interest is the CATE function on the target population $\tau_0(X) = E\{Y(1) - Y(-1)|X, S = 0\}$. To identify $\tau_0(X)$, we impose the following assumptions.

**Assumption A.1 (Ignorability)** $Y(a) \perp A \mid X, Z_s, S = s$ *for each* $s \in \{1, \cdots, K\}$.

**Assumption A.2 (Treatment Positivity)** $P(A = a \mid X = \boldsymbol{x}, Z_s, S = s) > 0$ *for all* $\boldsymbol{x}$ *such that* $f(X = \boldsymbol{x} \mid S = s)f(X = \boldsymbol{x} \mid S = 0) > 0$, *where* $a \in \{1, -1\}$ *and* $s \in \{1, \cdots, K\}$.

**Assumption A.3 (Population Positivity)** $\sum_{s=1}^{K} f(X = \boldsymbol{x} \mid S = s) > 0$ *for all* $\boldsymbol{x}$ *such that* $f(X = \boldsymbol{x} \mid S = 0) > 0$.

**Assumption A.4 (Homogeneity of CATE)** $\delta_0(X) = \cdots = \delta_K(X)$ *almost surely.*

Notably, in addition to the assumptions required in homogeneous causal data fusion, we impose the *population positivity* assumption, which implies that similar individuals exist in at least one trial dataset for those in the target population. Under Assumptions A.1–A.4, the CATE function $\tau_0(X)$ is identified.

The estimation procedures are the same as in Algorithm 1 of the main text, except that the information-aware weighting function

$$w_s(X) \propto \frac{P(S = s)}{P(S = 0)} \frac{\pi_0(X)}{\pi_s(X)} \left\{ \frac{V_{1|s}(X)}{p_{1|s}(X)} + \frac{V_{-1|s}(X)}{p_{-1|s}(X)} \right\}^{-1} \tag{A.1}$$

$$= \frac{f(X \mid S = 0)}{f(X \mid S = s)} \left\{ \frac{V_{1|s}(X)}{p_{1|s}(X)} + \frac{V_{-1|s}(X)}{p_{-1|s}(X)} \right\}^{-1}. \tag{A.2}$$

The density ratio $f(X \mid S = 0)/f(X \mid S = s)$ enables us to pay more attention to the subpopulations that account for a larger share of the target population. We assign a weight of zero to individuals not belonging to the target population, i.e., those with covariate $\boldsymbol{x}$ such that $f(X = \boldsymbol{x} \mid S = 0) = 0$. Analogous to the causal data fusion, the specification of $w_s(X)$ only relates to the efficiency, but not the consistency. One can either model propensity scores and then estimate $w_s(X)$ via Equation (A.1), or model densities and then estimate $w_s(X)$ via Equation (A.2). Under Assumptions A.1–A.4, the obtained $\hat{\delta}(X)$ shares the same properties as those described in the main text.