## **Estimating Joint Treatment Effects by Combining Multiple Experiments**

Yonghan Jung<sup>1</sup> Jin Tian<sup>2</sup> Elias Bareinboim<sup>3</sup>

## Abstract

Estimating the effects of multi-dimensional treatments (i.e., joint treatment effects) is critical in many data-intensive domains, including genetics and drug evaluation. The main challenges for studying the joint treatment effects include the need for large sample sizes to explore different treatment combinations as well as potentially unsafe treatment interactions. In this paper, we develop machinery for estimating joint treatment effects by combining data from multiple experimental datasets. In particular, first, we develop new identification conditions for determining whether a joint treatment effect can be computed in terms of multiple interventional distributions under various scenarios. Further, we develop estimators with statistically appealing properties, including consistency and robustness to model misspecification and slow convergence. Finally, we perform simulation studies, which corroborate the effectiveness of the proposed methods.

## 1. Introduction

A large body of scientific research is concerned with estimating the effect of multi-dimensional treatments. For example, Genome-Wide Association Studies (GWAS) in computational biology applications study the effect of multiple combinations of genes (Tam et al., 2019). As another example, estimating the multi-dimensional treatment effects is essential in the pharmaceutical industry because potential treatment-treatment interactions can lead to harmful effects to patients, potentially lethal in some situations. Consider two real-world scenarios in which understanding the treatment-treatment interaction is critical:

**Example TTI** (Treatment-Treatment-Interaction (Lee et al., 2019)). Many experimental studies have been conducted on the effects of antihypertensive drugs  $(X_1)$  on

blood pressure (W) with baseline characteristics ( $C_1$ ) (e.g., (Hansson et al., 1999)) and on the effects of anti-diabetic drugs ( $X_2$ ) on cardiovascular disease (Y) with baseline characteristics ( $C_2$ ) (e.g., (Ajjan & Grant, 2006)). Other studies reported that simultaneously taking both durgs was harmful to the population (Ferrannini & Cushman, 2012). This leaves open the question on how to evaluate the joint effect of antihypertensive and anti-diabetic medications from data coming from individual randomized studies.

**Example MTI** (Multiple Treatments Interactions). Many experimental studies have been conducted on the effects of (1) taking an aspirin  $(X_1)$  on blood pressure  $(W_1)$  (e.g., (Hansson et al., 1998); (2) taking acetaminophen  $(X_2)$  on blood coagulation  $(W_2)$  (e.g., (Gazzard et al., 1974)); and (3) taking the ibuprofen  $(X_3)$  on the gastrointestinal disease (Y) (e.g., (Lesko & Mitchell, 1995)). Other more recent studies reported adverse drug reactions to taking ibuprofen with aspirins and acetaminophen (Moore et al., 2015). What are the causal effects of the combinations of such drugs?

Despite their critical importance, the analysis of multidimensional effects remains underrepresented compared to the vast literature on single-treatment experiments. This is primarily due to two major challenges: the need for large sample sizes to investigate all possible treatment combinations and the possibility of unsafe or unethical treatment interactions (Examples TTI and MTI). It is, therefore, of great importance to investigate the possibility of estimating joint treatment effects by combining data from multiple *marginal experiments*, which refer to experiments on a subset of treatments (e.g., a single treatment). In this paper, we present novel methods for estimating joint effects given data from multiple marginal experiments and a qualitative description of the underlying causal system articulated in the form of a causal graph. Specifically,

1. We develop nonparametric identification criteria determining whether a joint treatment effect can be expressed through an adjustment formula using distributions from marginal experiments.

2. We construct estimators for the joint treatment effects using samples from marginal experiments and provide learning guarantees for the estimators. We illustrate the empirical validity of the estimators through simulations.

<sup>&</sup>lt;sup>1</sup>Purdue University <sup>2</sup>Iowa State University <sup>3</sup>Columbia University. Correspondence to: Yonghan Jung <jung222@purdue.edu>.

Proceedings of the 40<sup>th</sup> International Conference on Machine Learning, Honolulu, Hawaii, USA. PMLR 202, 2023. Copyright 2023 by the author(s).

The proofs are provided in Appendix C in suppl. material.

## 1.1. Related Work

**Causal Effect Identification and Estimation.** Recent advances in the literature of *generalized causal effect identification* (g-ID) lead to algorithmic solutions for determining the identification of a causal effect from a set of observational and experimental studies given a causal graph (Bareinboim & Pearl, 2012a; 2016; Lee et al., 2019; Lee & Bareinboim, 2020; Lee et al., 2020; Correa et al., 2021). In addition, recent progress has been made in developing estimators for any causal effects identifiable from observational data in a causal graph (Jung et al., 2020; 2021a;b; Bhattacharya et al., 2022; Jung et al., 2022). However, these estimators are not applicable to g-ID settings that involve multiple experimental distributions.

On a different thread, estimating causal effects from multiple experiments and observations has been investigated for some specific settings. For example, the problems of estimating the long-term effect of a single treatment by combining multiple short-term experimental studies into a surrogate variable have been recently studied (e.g., Bareinboim & Pearl (2012b); Athey et al. (2019; 2020); Imbens et al. (2022)). In epidemiology, estimators for causal effects in a target domain by combining multiple experiments in different source domains have been developed (e.g., Dahabreh et al. (2019); Colnet et al. (2020); Degtiar & Rose (2021); Shi et al. (2022)). However, these methods are not applicable when the goal is to estimate the joint treatment effects from multiple marginal experiments.

**Treatment Combinations.** The aforementioned examples are related to the analysis of treatment combinations which aims to attribute the joint treatment effects to either the effect of treatment combination or marginal treatment effects (e.g., VanderWeele & Knol (2014); Egami & Imai (2018); Parbhoo et al. (2021)). Existing literature commonly relies on the back-door criterion. Such assumptions, however, are not satisfied when latent confounders exist, as illustrated in Examples (TTI, MTI) and Figs. (1a, 2a).

Closer to our work is Saengkyongam & Silva (2020), which investigates the identifiability of joint effects in the additive models with Gaussian noises and continuous treatments by entangling observations and marginal experiments. However, their approach is inapplicable when the treatment variables are discrete, which is common in many applications. In contrast to their methods, we provide nonparametric identifiability criteria for the joint effects from marginal experiments based on a causal graph without imposing constraints on the data-generating processes. Additionally, we develop estimators for joint treatment effects having statistically desirable properties.

## 2. Preliminaries

Notations. Each variable is represented with a capital letter (X) and its realized value with a small letter (x). We use bold letters (X) to denote a random vector. Given an ordered set  $\mathbf{X} = (X_1, \dots, X_n)$  such that  $X_i \prec X_j$  for i < j, we denote  $\mathbf{X}^{(i)} = \{X_1, \dots, X_i\}$ . For a graph G over  $\mathbf{V}$  and disjoint vectors  $\mathbf{X}_1, \mathbf{X}_2 \subseteq \mathbf{V}$ , we will use  $G_{\overline{\mathbf{X}_1} \underline{\mathbf{X}_2}}$ as a subgraph of G in which all incoming edges to the node in  $X_1$  and all outgoing edges to the node in  $X_2$  are cut. For a discrete (e.g., binary) random vector X and its realized value  $\mathbf{x} \in \mathfrak{D}_{\mathbf{X}}$  where  $\mathfrak{D}_{\mathbf{X}}$  is the domain of  $\mathbf{X},$  we use  $\mathbbm{1}_{\mathbf{x}}(\mathbf{X})$  to represent the indicator function such that  $\mathbb{1}_{\mathbf{x}}(\mathbf{X}) = 1$  if  $\mathbf{X} = \mathbf{x}$ ;  $\mathbb{1}_{\mathbf{x}}(\mathbf{X}) = 0$  otherwise. For a random vector  $\mathbf{X}$ , we use  $P(\mathbf{X})$  to denote its distribution and  $p(\mathbf{x})$  as a corresponding density function at  $\mathbf{X} = \mathbf{x}$ . For a function f,  $\mathbb{E}_P[f(\mathbf{X})] \coloneqq \int_{\mathfrak{S}_{\mathbf{X}}} f(\mathbf{x}) p(\mathbf{x}) d[\mathbf{x}]$  where  $\mathfrak{S}_{\mathbf{X}}$  is the support for **X**. For a sample set  $D \coloneqq \{\mathbf{V}_{(i)}\}_{i=1}^{n}$ where  $\mathbf{V}_{(i)}$  denotes the *i*th samples, we use  $\mathbb{E}_D \left[ f(\mathbf{V}) \right] \coloneqq$  $(1/n)\sum_{i=1}^{n} f(\mathbf{V}_{(i)})$ . We use  $\|f\|_P \coloneqq \sqrt{\mathbb{E}_P[(f(\mathbf{X}))^2]}$ . If a function  $\hat{f}$  is a consistent estimator of f having a rate  $r_n$ , we will use  $\hat{f} - f = o_P(r_n)$ . We will say  $\hat{f}$  is  $L_2$ -consistent if  $\|\hat{f} - f\|_P = o_P(1)$ . We will use  $\hat{f} - f = O_P(1)$  if  $\hat{f} - f$ is bounded in probability. Also,  $\hat{f} - f$  is said to be bounded in probability at rate  $r_n$  if  $\hat{f} - f = O_P(r_n)$ . Throughout the paper, we assume that samples D are independent.

Structural Causal Models. We use Structural Causal Models (SCMs) as our framework (Pearl, 2000; Bareinboim et al., 2022). An SCM  $\mathcal{M}$  is a quadruple  $\mathcal{M} =$  $\langle \mathbf{U}, \mathbf{V}, P(\mathbf{U}), F \rangle$ . U is a set of exogenous (latent) variables following a joint distribution  $P(\mathbf{U})$ . V is a set of endogenous (observable) variables whose values are determined by functions  $F = \{f_{V_i}\}_{V_i \in \mathbf{V}}$  such that  $V_i \leftarrow f_{V_i}(pa_i, u_i)$ where  $PA_i \subseteq V$  and  $U_i \subseteq U$ . Each SCM  $\mathcal{M}$  induces a distribution  $P(\mathbf{V})$  and a causal graph  $G = G(\mathcal{M})$  over  $\mathbf{V}$ in which there exists a directed edge from every variable in  $PA_i$  to  $V_i$  and dashed-bidirected arrows encode common latent variables (e.g., see Fig. 1a). Performing an intervention fixing  $\mathbf{X} = \mathbf{x}$  is represented through the do-operator,  $do(\mathbf{X} = \mathbf{x})$ , which encodes the operation of replacing the original equations of X (i.e.,  $f_X(pa_x, u_x)$ ) by the constant  $x \in \mathfrak{D}_X$  for all  $X \in \mathbf{X}$  and induces an interventional distribution  $P(\mathbf{V}|do(\mathbf{x}))$ . We will sometimes employ the shorthand notation  $P_{\mathbf{x}}(\mathbf{y})$  to represent  $P(\mathbf{y}|do(\mathbf{x}))$ . We will use  $P_{\text{rand}(\mathbf{X})}(\mathbf{Y}) \coloneqq \{P_{\mathbf{x}}(\mathbf{Y})\}_{\mathbf{x}\in\mathfrak{S}_{\mathbf{X}}}$ . For a sample set  $D \coloneqq {\mathbf{V}_{(i)}}_{i=1}^{n}, D$  is said to follow  $\hat{P}_{rand(\mathbf{X})}(\mathbf{V})$  if each subsamples  $D_{\mathbf{x}} \coloneqq {\mathbf{V}_{(i)}}_{\mathbf{V}_{(i)} \in D, \mathbf{X}_{(i)} = \mathbf{x}}$  follows  $P_{\mathbf{x}}(\mathbf{V})$ .

## 3. Combining Two Experiments

In this section, we address the challenge of estimating the combined effects by leveraging the results of two distinct experiments. In Section 3.1, we delve into the estimation

of treatment-treatment interactions (TTI) based on the outcomes of two separate marginal experiments. Then, in Section 3.2, we extend our investigation to accommodate scenarios where the treatments in the source and target experiments may not align perfectly.

#### **3.1. Treatment-Treatment Interaction**

Our goal is to estimate joint effects by combining two *randomized controlled experiments*, as formally defined below.

**Task TTI** (**Treatment-Treatment Interaction (TTI**)). The task of estimating the treatment-treatment interactions (TTI) from two marginal experiments composes of

• **Input**: Two sets of samples,  $D_1$  and  $D_2$ , which follow interventional distributions  $P_{\text{rand}(X_1)}(C_1, X_1, W)$  and  $P_{\text{rand}(X_2)}(C_1, W, X_1, X_2, Y)$ , respectively.  $C_1$  is a covariate for the experiment randomizing  $X_1$  (i.e.,  $\text{rand}(X_1)$ ), and W and Y represent the outcomes of the experiments randomizing  $X_1$  (rand $(X_1)$ ) and  $X_2$  (rand $(X_2)$ ), respectively.

• Query: Estimation of  $\mathbb{E}[Y|do(x_1, x_2)]$ .

## 3.1.1. ADJUSTMENT CRITERION FOR TTI (AC-TTI)

A sufficient graphical criterion for identifying the treatmenttreatment interaction is the following:

**Definition 1 (Adjustment criterion for Treatment-Treatment Interaction (AC-TTI)).** A set  $\{C_1, W\}$  is said to satisfy the *adjustment criterion for treatment-treatment interaction (AC-TTI)* w.r.t  $\{(X_1, X_2), Y\}$  in G if

1.  $(\{C_1, W\} \perp X_2 | X_1)_{G_{\overline{X_1, X_2}}}$ , i.e., there are no direct paths from  $X_2$  to  $\{C_1, W\}$  in  $G_{\overline{X_1, X_2}}$ ; and

2.  $(Y \perp X_1 | C_1, W, X_2)_{G_{\underline{X_1 X_2}}}$ , i.e., the back-door paths from  $X_1$  to Y are blocked by  $\{C_1, W\}$  in  $G_{\overline{X_2}}$ .

We make the following positivity assumption:

Assumption 1 (Positivity Assumption for AC-TTI).  $P_{x_1}(C_1, W), P_{x_2}(C_1, W), P_{x_2}(X_1|C_1, W)$  are strictly positive distributions  $\forall x_1, x_2 \in \mathfrak{D}_{X_1, X_2}$ .

Under AC-TTI and Assumption 1, the joint treatment effect  $\mathbb{E}[Y|do(x_1, x_2)]$  is identifiable and given as follows:

**Theorem 1** (Identification through AC-TTI). Suppose AC-TTI in Def. 1 and Assumption 1 hold. Then,  $\mathbb{E}[Y|do(x_1, x_2)]$  is identifiable from  $P_{rand(X_1)}(C_1, W)$  and  $P_{rand(X_2)}(C_1, W, X_1, Y)$  and the expression is:

$$\mathbb{E}\left[Y|do(x_1, x_2)\right] = \mathbb{E}_{P_{x_1}}\left[\mathbb{E}_{P_{x_2}}\left[Y|C_1, W, x_1\right]\right].$$
 (1)

For example, in Fig. 1a,  $\{C_1, W\}$  satisfies AC-TTI w.r.t.  $\{(X_1, X_2), Y\}$ . Therefore, with Assumption. 1,  $\mathbb{E}[Y|do(x_1, x_2)]$  is identifiable from  $P_{\text{rand}(X_1)}(C_1, W)$  and  $P_{\text{rand}(X_2)}(C_1, W, X_1, Y)$  as in Eq. (1).



Figure 1: Example causal graphs for Section 3. Nodes representing the treatment and outcome are marked in blue and red respectively.

#### 3.1.2. ESTIMATORS FOR AC-TTI

We define the nuisance functionals for estimating the AC-TTI functional in Eq. (1) as follows:

**Definition 2** (Nuisance for AC-TTI). Nuisance functions for the AC-TTI functional in Eq. (1) are defined as follows: For a fixed  $x_1, x_2 \in \mathfrak{D}_{X_1, X_2}$  where  $x_1, x_2$  are specified in Eq. (1),  $\pi_0 \coloneqq \pi_0(C_1, X_1, W) \coloneqq \frac{P_{x_1}(W|C_1)}{P_{x_2}(W, X_1|C_1)}$ . Also,  $\mu_0 \coloneqq \mu_0(C_1, X_1, W) \coloneqq \mathbb{E}_{P_{x_2}}[Y|X_1, W, C_1]$ . We will use  $\pi \coloneqq \pi(C_1, X_1, W) > 0$  and  $\mu \coloneqq \mu(C_1, X_1, W)$  to denote arbitrary<sup>1</sup> finite functions.

Now, we construct regression-based ('REG'), probability weighting ('PW'), and double/debiased machine learning ('DML') (Chernozhukov et al., 2018) based estimators:

**Definition 3** (AC-TTI estimators). Let  $D_1$  and  $D_2$  denote two separate samples following the distributions  $P_{\text{rand}(X_1)}(C_1, W)$  and  $P_{\text{rand}(X_2)}(C_1, W, X_1, Y)$ , respectively. For fixed  $x_1, x_2 \in \mathfrak{D}_{X_1, X_2}$ , we define  $D_{x_1}$  and  $D_{x_2}$  as subsamples of  $D_1$  and  $D_2$  such that  $X_1 = x_1$  and  $X_2 = x_2$ . Let  $\mu$  and  $\pi$  denote the nuisances as defined in Definition 2. We now introduce the {REG, PW, DML} estimators for the AC-TTI-functional specified in Equation (1) as follows:

$$T^{reg} \coloneqq \mathbb{E}_{D_{x_1}} \left[ \mu(W, C_1, x_1) \right],$$
  

$$T^{pw} \coloneqq \mathbb{E}_{D_{x_2}} \left[ \pi(W, C_1, X_1) \mathbb{1}_{x_1}(X_1) Y \right],$$
  

$$T^{dml} \coloneqq \mathbb{E}_{D_{x_1}} \left[ \pi \mathbb{1}_{x_1}(X_1) \{Y - \mu\} \right] + \mathbb{E}_{D_1} \left[ \mu(W, C_1, x_1) \right) \right]$$

We assume that samples used for training the nuisance functions and evaluating the nuisances are independent:

**Assumption 2** (**Sample-splitting**). Samples for training nuisances and evaluating the estimators equipped with the trained nuisance are separate and independent<sup>2</sup>.

<sup>&</sup>lt;sup>1</sup>Throughout the paper,  $\mu, \pi$  are understood as estimated nuisances for  $\mu_0, \pi_0$ .

<sup>&</sup>lt;sup>2</sup>This assumption is satisfied by applying cross-fitting algo-

We assume that nuisances can be estimated  $L_2$  consistently. In practice, this assumption can be easily satisfied by employing flexible machine learning models.

Assumption 3 ( $L_2$  consistency of nuisances). Estimated nuisances are  $L_2$  consistent; i.e.,  $\forall i \in \{1,2\}, \forall x_i \in \mathfrak{D}_{X_i}$ ,

$$\|\mu(W, C_1, x_1) - \mu_0(W, C_1, x_1)\|_{P_{x_i}} = o_{P_{x_i}}(1), \|\pi(W, C_1, X_1) - \pi_0(W, C_1, X_1)\|_{P_{x_2}} = o_{P_{x_2}}(1).$$

We also assume that the baseline covariates have the same distribution over all marginal experiments. Specifically,

Assumption 4 (Shared Covariates). The distributions of the covariates  $C_1$  are the same; i.e., for all  $x_1, x_2 \in \mathfrak{D}_{X_1,X_2}$ ,  $P_{x_1}(C_1) = P_{x_2}(C_1)$ .

In practice, this assumption may not hold due to "covariate shift". We discuss relaxing Assumption 4 in Appendix F.1.

Then, the errors for each estimator are given as follows:

**Theorem 2 (Error analysis for AC-TTI estimators).** Under Assumptions (1,2,3,4) and AC-TTI in Def. 1, the error of the estimators in Def. 3, denoted  $\epsilon^{est} \coloneqq T^{est} - \mathbb{E}[Y|do(x_1, x_2)]$  for est  $\in \{reg, pw, dml\}$  are:

$$\begin{aligned} \epsilon^{reg} &= R_1 + O_{P_{x_1}} \left( \| \mu - \mu_0 \| \right), \\ \epsilon^{pw} &= R_2 + O_{P_{x_2}} \left( \| \pi - \pi_0 \| \right), \\ \epsilon^{dml} &= R_1 + R_2 + O_{P_{x_2}} \left( \| \pi - \pi_0 \| \| \mu - \mu_0 \| \right) \end{aligned}$$

where  $R_i$  is a random variable such that  $\sqrt{n_i}R_i$  converges in distribution to a zero-mean normal random variable, where  $n_i := |D_{x_i}|$  for  $i \in \{1, 2\}$ .

We highlight that the DML estimator  $T^{dml}$  exhibits robustness property since  $\epsilon^{dml}$  is bounded in probability at  $n^{-1/2}$ rate (for  $n = \min\{n_1, n_2\}$ ) whenever  $\|\pi - \pi_0\|_{P_{x_2}} = O_{P_{x_2}}(n^{-1/4})$  and  $\|\mu - \mu_0\|_{P_{x_2}} = O_{P_{x_2}}(n^{-1/4})$ . Furthermore, the DML estimator displays the following doubly robustness property:

**Corollary 2** (Doubly robustness of the DML estimators (Corollary of Thm. 2)). Suppose Assumptions (1,2,3,4) and AC-TTI in Def. 1 hold. Suppose either  $\pi = \pi_0$  or  $\mu = \mu_0$ . Then,  $T^{dml}$  is an unbiased estimator of  $\mathbb{E}[Y|do(x_1, x_2)]$ .

## 3.2. Combining Two Arbitrary Experiments

In this section, we extend Task TTI to cases where the effect with two or more treatments (i.e.,  $|\mathbf{X}| \ge 2$ ) can be identified from two arbitrary experiments. For example, let's consider a scenario extending Example TTI where

we are interested in studying the effect of three factors: the antihypertensive drug  $(Z_1)$ , the anti-diabetic drug  $(Z_2)$ , and the individual's diet habits  $(X_0)$ , on the occurrence of cardiovascular disease (as depicted in Fig. 1b) when we are given two marginal experiments randomizing  $Z_1$  and  $Z_2$ respectively. This extended task is referred to as *generalized treatment-treatment interactions* (gTTI) and is defined as follows:

**Task gTTI (Generalized TTI).** The task of *generalized TTI* composes of

• Input: Two sets of samples  $D_1$ ,  $D_2$  following distributions  $P_{\text{rand}(Z_1)}(\mathbf{V})$  and  $P_{\text{rand}(Z_2)}(\mathbf{V})$ , respectively.

• Query: Estimation of  $\mathbb{E}[Y|do(\mathbf{x})]$ .

We note that Task gTTI generalizes Task TTI in the sense that it does not require that X is identical to  $Z_1 \cup Z_2$ .

3.2.1. Adjustment Criterion for GTTI (AC-GTTI)

A graphical criterion for identifying  $\mathbb{E}[Y|do(\mathbf{x})]$  from two distributions  $P_{\text{rand}(Z_1)}(\mathbf{V})$  and  $P_{\text{rand}(Z_2)}(\mathbf{V})$  is the following:

**Definition 4 (Adjustment criterion for combining two experiments (AC-gTTI)).** A set of variables **A** is said to satisfy *adjustment criterion for generalized TTI (AC-gTTI)* w.r.t (**X**, Y) in G if

1.  $Z_1 \subseteq \mathbf{X}$  and  $(\mathbf{A} \perp \mathbf{X} \setminus Z_1 | Z_1)_{G_{\overline{\mathbf{X}}}}$ , i.e., there are no direct paths from  $\mathbf{X} \setminus Z_1$  to  $\mathbf{A}$  in  $G_{\overline{\mathbf{X}}}$ ; and

2.  $Z_2 \subseteq \mathbf{X}$  and  $(Y \perp \mathbf{X} \setminus Z_2 | \mathbf{A}, Z_2)_{G_{\mathbf{X} \setminus Z_2} \overline{Z_2}}$ , i.e., the back-door paths from  $\mathbf{X} \setminus Z_2$  to Y are blocked by  $\mathbf{A}$  in  $G_{\overline{Z_2}}$ .

We make the following positivity assumption:

Assumption 5 (Positivity Assumption for AC-gTTI).  $P_{z_1}(\mathbf{A}), P_{z_2}(\mathbf{A}), P_{z_2}(\mathbf{X} \setminus Z_2 | \mathbf{A})$  are strictly positive distributions  $\forall z_1, z_2 \in \mathfrak{D}_{Z_1, Z_2}$ .

Under AC-gTTI, the joint treatment effect  $\mathbb{E}[Y|do(\mathbf{x})]$  is identifiable and given as follows:

**Theorem 3 (Identification through AC-gTTI).** Suppose *AC-gTTI in Def.* 4 and Assumption 5 hold. Then, the query  $\mathbb{E}[Y|do(\mathbf{x})]$  is identifiable from  $P_{rand(Z_1)}(\mathbf{A})$  and  $P_{rand(Z_2)}(\mathbf{A}, \mathbf{X}, Y)$  and given as follows:

$$\mathbb{E}\left[Y|do(\mathbf{x})\right] = \mathbb{E}_{P_{z_1}}\left[\mathbb{E}_{P_{z_2}}\left[Y|\mathbf{A}, \mathbf{x} \setminus \mathbf{z}_2\right]\right].$$
 (2)

For example, in Fig. 1b,  $\mathbf{A} \coloneqq \{W, C\}$  satisfies AC-gTTI criterion w.r.t.  $\{\mathbf{X} = (X_0, Z_1, Z_2), Y\}$ . Therefore, under positivity,  $\mathbb{E}[Y|do(\mathbf{x})]$  is expressible as in Eq. (2).

,

rithms (e.g., (Klaassen, 1987; Robins & Ritov, 1997; Zheng & van der Laan, 2011; Chernozhukov et al., 2018)), which split the samples and using one for training nuisances and another for evaluating the trained nuisances.

#### 3.2.2. ESTIMATORS FOR AC-GTTI

We define the nuisance functionals for the AC-gTTI functional in Eq. (2) as follows:

**Definition 5** (Nuisances for AC-gTTI). Nuisance functions for estimating AC-gTTI functional in Eq. (2) are defined as follows: For a fixed  $z_1, z_2 \in \mathfrak{D}_{Z_1,Z_2}$  where  $z_1, z_2$ are specified in Eq. (2),  $\pi_0 \coloneqq \pi_0(\mathbf{A}, \mathbf{X}) \coloneqq \frac{P_{z_1}(\mathbf{A})}{P_{z_2}(\mathbf{A}, \mathbf{X} \setminus Z_2)}$ , and  $\mu_0 \coloneqq \mu_0(\mathbf{A}, \mathbf{X}) \coloneqq \mathbb{E}_{P_{z_2}}[Y|\mathbf{X} \setminus Z_2, \mathbf{A}]$ . We will use  $\pi \coloneqq \pi(\mathbf{A}, \mathbf{X}) > 0$  and  $\mu \coloneqq \mu(\mathbf{A}, \mathbf{X})$  to denote estimated nuisances.

Now, we construct regression-based ('REG'), probability weighting ('PW'), and double/debiased machine learning ('DML') estimators:

**Definition 6** (AC-gTTI estimators). Let  $D_1$ ,  $D_2$  denote two sample sets following distributions  $P_{\text{rand}(Z_1)}(\mathbf{A})$  and  $P_{\text{rand}(Z_2)}(\mathbf{A}, \mathbf{X}, Y)$ , respectively. For a fixed  $z_1, z_2 \in \mathfrak{D}_{Z_1, Z_2}$ , we define  $D_{z_1}$  and  $D_{z_2}$  as subsamples of  $D_1$  and  $D_2$  such that  $Z_1 = z_1$  and  $Z_2 = z_2$ . Let  $\mu, \pi$  denote nuisances defined in Def. 5. Then, {REG, PW, DML} estimators for the AC-gTTI functional are defined as follows:

$$\begin{split} T^{reg} &\coloneqq \mathbb{E}_{D_{z_1}} \left[ \mu(\mathbf{A}, \mathbf{x}) \right], \\ T^{pw} &\coloneqq \mathbb{E}_{D_{z_2}} \left[ \pi(\mathbf{A}, \mathbf{X}) \mathbb{1}_{\mathbf{x}}(\mathbf{X}) Y \right], \\ T^{dml} &\coloneqq \mathbb{E}_{D_{z_2}} \left[ \pi \mathbb{1}_{\mathbf{x}}(\mathbf{X}) \{Y - \mu\} \right] + \mathbb{E}_{D_{z_1}} \left[ \mu(\mathbf{A}, \mathbf{x}) \right] \end{split}$$

We assume that nuisances can be estimated  $L_2$  consistently.

Assumption 6 ( $L_2$  consistency of nuisances). Estimated nuisances are  $L_2$  consistent; i.e.,  $\forall i \in \{1,2\}, \forall z_i \in \mathfrak{D}_{Z_i}$ ,

$$\|\mu(\mathbf{A}, \mathbf{x}) - \mu_0(\mathbf{A}, \mathbf{x})\|_{P_{z_i}} = o_{P_{z_i}}(1),$$
  
$$\|\pi(\mathbf{A}, \mathbf{X}) - \pi_0(\mathbf{A}, \mathbf{X})\|_{P_{z_2}} = o_{P_{z_2}}(1).$$

Then, the error of each estimator is given as follows:

**Theorem 4 (Error analysis for AC-gTTI estimators).** Under Assumptions (2,5,6) and AC-gTTI in Def. 4, the errors of the estimators in Def. 6, denoted  $\epsilon^{est} := T^{est} - \mathbb{E}[Y|do(\mathbf{x})]$ for  $est \in \{reg, pw, dml\}$ , are:

$$\begin{aligned} \epsilon^{reg} &= R_1 + O_{P_{z_1}} \left( \| \mu - \mu_0 \| \right), \\ \epsilon^{pw} &= R_2 + O_{P_{z_2}} \left( \| \pi - \pi_0 \| \right), \\ \epsilon^{dml} &= R_1 + R_2 + O_{P_{z_2}} \left( \| \pi - \pi_0 \| \| \mu - \mu_0 \| \right), \end{aligned}$$

where  $R_i$  is a random variable such that  $\sqrt{n_i}R_i$  converges in distribution to a zero-mean normal random variable, where  $n_i := |D_{z_i}|$ .

We highlight that the DML estimator  $T^{\text{dml}}$  exhibits robustness property since  $\epsilon^{\text{dml}}$  is bounded in probability at  $n^{-1/2}$  (for  $n = \min\{n_1, n_2\}$ ) rate whenever  $\|\pi - \pi_0\|_{P_{z_2}} = O_{P_{z_2}}(n^{-1/4})$  and  $\|\mu - \mu_0\|_{P_{z_2}} = O_{P_{z_2}}(n^{-1/4})$ . Furthermore, the DML estimator displays the following doubly robustness property:

**Corollary 4** (Doubly robustness of the DML estimators (Corollary of Thm. 4)). Suppose Assumptions (2,5,6) and AC-gTTI in Def. 4 hold. Suppose either  $\pi = \pi_0$  or  $\mu = \mu_0$ . Then,  $T^{dml}$  is an unbiased estimator of  $\mathbb{E}[Y|do(\mathbf{x})]$ .

## **4.** Combining Multiple ( $\geq$ 2) Experiments

In this section, we address the estimation of joint effects by leveraging multiple (more than two) experiments. Specifically, in Sec. 4.1, we focus on estimating multiple treatment interactions (MTI) using multiple marginal experiments. In Sec. 4.2, we extend this setting to estimate multiple treatment effects from multiple experiments that are not necessarily over each element in **X**.

#### 4.1. Multiple Treatment Interaction

We first introduce the formal version of the task.

**Task MTI** (**Multiple-Treatment Interaction (MTI**)). Estimating multiple treatment interaction (MTI) composes of

• Input: Multiple sets of samples  $\{D_i\}_{i=1}^m$  drawn from a sequence of interventional distributions  $\{P_{\text{rand}(X_i)}(\mathbf{C}^{(i)}, \mathbf{W}^{(i)}, \mathbf{X}^{(i)})\}_{i=1}^m$ .  $\{C_i, X_i, W_i\}$  for  $i = 1, \cdots, m$  is the *i*th triplet corresponding to the covariate, the treatment, and the outcome.

• Query: Estimation of  $\mathbb{E}[Y|do(\mathbf{x})]$  where  $\mathbf{x} = \{x_i\}_{i=1}^m$  is a realization of the ordered set  $\mathbf{X} := \{X_1, \cdots, X_m\}$ , and  $Y := W_m$ .

### 4.1.1. Adjustment Criterion for MTI (AC-MTI)

A sufficient graphical criterion for identifying the multiple treatment interaction is the following:

**Definition 7** (Adjustment criterion for Multiple **Treatment Interaction (AC-MTI)).** An ordered set  $\{C_1, W_1, C_2, W_2, \dots, C_{m-1}, W_{m-1}\}$  satisfies *adjustment criterion for multiple treatment interaction (AC-MTI)* w.r.t.  $\{\mathbf{X}, Y\}$  for  $\mathbf{X} = \{X_i\}_{i=1}^m$  in *G* if, for  $i = 1, 2, \dots, m$ ,

1.  $\{X_i\}_{i>i}$  is non-ancestor of  $\{\mathbf{X}^{(i)}, \mathbf{W}^{(i)}, \mathbf{C}^{(i)}\}$ ; and

2.  $(Y \perp X_i | \mathbf{C}^{(i-1)}, \mathbf{X}^{(i-1)}, \mathbf{W}^{(i)}, \mathbf{X}^{>i})_{G_{\underline{X_i}, \mathbf{X}^{>i}}}$ , i.e., the back-door paths from  $X_i$  to Y are blocked by

 $\mathbf{C}^{(i-1)}, \mathbf{X}^{(i-1)}, \mathbf{W}^{(i)}, \mathbf{X}^{>i}$  in the graph  $G_{\overline{\mathbf{X}^{>i}}}$ .

We make the following positivity assumption:

Assumption 7 (Positivity Assumption for AC-MTI).  $\{P_{x_i}(W_i, C_i | \mathbf{W}^{(i-1)}, \mathbf{C}^{(i-1)}, \mathbf{X}^{(i-1)})\}_{i=1}^m$ 

 $P_{x_{i+1}}(X_i | \mathbf{W}^{(i)}, \mathbf{C}^{(i)}, \mathbf{X}^{(i-1)})$  for  $i = 1, \dots, m-1$  are strictly positive  $\forall \mathbf{x} \in \mathfrak{D}_{\mathbf{X}}$ .

Under AC-MTI, the joint treatment effects  $\mathbb{E}[Y|do(\mathbf{x})]$  is identifiable and given as follows:

**Theorem 5** (Identification through AC-MTI). Suppose AC-MTI in Def. 7 and Assumption 7 hold. Then,  $\mathbb{E}[Y(\mathbf{x})]$ is identifiable from  $\{P_{rand(X_i)}(\mathbf{C}^{(i)}, \mathbf{W}^{(i)}, \mathbf{X}^{(i-1)})\}_{i=1}^{m}$  as follows: Let  $\mu_0^m \coloneqq \mathbb{E}_{P_{x_m}}[Y|\mathbf{W}^{(m-1)}, \mathbf{C}^{(m-1)}, \mathbf{X}^{(m-1)}]$ , and for  $i = m - 1, \dots, 2$ ,

$$\boldsymbol{\mu}_{0}^{i}\coloneqq \mathbb{E}_{P_{x_{i}}}\left[\overline{\boldsymbol{\mu}}_{0}^{i+1}|\mathbf{W}^{(i-1)},\mathbf{C}^{(i-1)},\mathbf{X}^{(i-1)}\right]$$

where  $\overline{\mu}_{0}^{i+1} \coloneqq \mu_{0}^{i+1}(\mathbf{W}^{(i)}, \mathbf{C}^{(i)}, x_{i}, \mathbf{X}^{(i-1)})$ . Then,  $\mathbb{E}[Y(\mathbf{x})] = \mathbb{E}_{P_{x_{1}}} \left[ \mu_{0}^{2}(W_{1}, C_{1}, x_{1}) \right].$  (3)

For example, in Fig. 2a,  $\{C_1, W_1, C_2, W_2\}$  satisfies AC-MTI w.r.t.  $\{(X_1, X_2), Y\}$  in Def. 7. Therefore, with the positivity assumption in Assumption. 7,  $\mathbb{E}[Y|do(\mathbf{x})]$  is identifiable from  $\{P_{\text{rand}(X_i)}(\mathbf{C}^{(i)}, \mathbf{W}^{(i)}, \mathbf{X}^{(i-1)})\}_{i=1}^m$  as in Eq. (3).

## 4.1.2. ESTIMATORS FOR AC-MTI

We define the nuisance functionals for estimating the AC-MTI functional in Eq. (3) as follows:

**Definition 8** (Nuisances for AC-MTI). Nuisance functions for AC-MTI are defined as follows: For a fixed  $\mathbf{x} \coloneqq \{x_1, \dots, x_m\} \in \mathfrak{D}_{\mathbf{X}}, \text{ let } \{\mu^i\}_{i=2}^m \text{ and } \{\overline{\mu}^i\}_{i=2}^m$  be the nuisances defined in Thm. 5. For  $i = 1, \dots, m-1, \pi_0^i \coloneqq \frac{P_{x_i}(W_i|C_i, \mathbf{C}^{(i-1)}, \mathbf{X}^{(i-1)}, \mathbf{W}^{(i-1)})}{P_{x_m}(W_i, X_i|C_i, \mathbf{C}^{(i-1)}, \mathbf{X}^{(i-1)}, \mathbf{W}^{(i-1)})}$ , and  $\pi_0^{(i)} \coloneqq \prod_{j=1}^i \pi_0^j(\mathbf{W}^{(j)}, \mathbf{C}^{(j)}, \mathbf{X}^{(j)})$ . We will use  $\pi^i(\mathbf{W}^{(i)}, \mathbf{C}^{(i)}, \mathbf{X}^{(i)}) > 0$  and  $\mu^i(\mathbf{W}^{(i-1)}, \mathbf{C}^{(i-1)}, \mathbf{X}^{(i-1)})$  to denote estimated nuisances.

Now, we construct regression-based ('REG'), probability weighting ('PW'), and double/debiased machine learning ('DML') estimators:

**Definition 9 (AC-MTI estimators).** Let  $D_i$  denote samples following  $P_{\text{rand}(X_i)}(\mathbf{C}^{(i)}, \mathbf{W}^{(i)}, \mathbf{X}^{(i)})$  for  $i = 1, 2, \cdots, m$ . For a fixed  $x_i \in \mathfrak{D}_{X_i}$ , let  $D_{x_i}$  denote the subsamples of  $D_i$  such that  $X_i = x_i$ . Let  $A_i \coloneqq \{W_i, C_i\}$  and  $V_i \coloneqq$  $\{A_i, X_i\}$ . Let  $\mu^{m+1} \coloneqq Y$ . Let  $\mathbb{1}_{\mathbf{x}}^{i-1} \coloneqq \mathbb{1}_{\mathbf{x}^{(i-1)}}(\mathbf{X}^{(i-1)})$ for  $i = 2, \cdots, m$ . Then {REG, PW, DML} estimators are defined as follows:

$$T^{reg} := \mathbb{E}_{D_{x_1}} \left[ \mu^2(W_1, C_1, x_1)) \right],$$
  

$$T^{pw} := \mathbb{E}_{D_{x_m}} \left[ \pi^{(m-1)} \mathbb{1}_{\mathbf{x}}(\mathbf{X}) Y \right],$$
  

$$T^{dml} := \sum_{i=2}^m \mathbb{E}_{D_{x_i}} \left[ \pi^{(i-1)} \mathbb{1}_{\mathbf{x}}^{i-1} \{ \overline{\mu}^{i+1} - \mu^i \} \right] + \mathbb{E}_{D_{x_1}} \left[ \overline{\mu}^2 \right]$$

We assume that nuisances can be estimated  $L_2$  consistently.



Figure 2: Example causal graphs for Section 4.

**Assumption 8** ( $L_2$  consistency of nuisances). *Estimated* nuisances are  $L_2$ -consistent; specifically,

$$\begin{split} \|\mu^{i+1} - \mu_0^{i+1}\|_{P_{x_i}} &= o_{P_{x_i}}(1), \; \forall i \in \{1, 2, \cdots, m-1\} \\ \|\mu^i - \mu_0^i\|_{P_{x_i}} &= o_{P_{x_i}}(1), \; \forall i \in \{2, \cdots, m\} \\ \|\pi^i - \pi^i\|_{P_{x_{i+1}}} &= o_{P_{x_{i+1}}}(1), \; \forall i \in \{1, \cdots, m-1\}. \end{split}$$

We assume that the baseline covariates  $C_i$  in the *i*th experiment follow the same distribution as in the *j*th experiment for j > i:

Assumption 9 (Shared Covariates). For any fixed  $i, j \in \{1, 2, \dots, m-1\}$  s.t. j > i and any fixed  $x_i, x_j \in \mathfrak{D}_{X_i, X_j}$ , the baseline covariates  $C_i$ 's distribution satisfies the following:  $P_{x_i}(C_i|\mathbf{C}^{(i-1)}, \mathbf{X}^{(j-1)}, \mathbf{W}^{(j-1)}) = P_{x_j}(C_i|\mathbf{C}^{(i-1)}, \mathbf{X}^{(i-1)}, \mathbf{W}^{(i-1)}).$ 

We discuss the relaxation of Assumption 9 in Appendix F.2.

**Theorem 6 (Error analysis of AC-MTI estimators).** Under Assumptions (2,7,8,9) and AC-MTI in Def. 7, the errors of the estimators in Def. 9, denoted  $\epsilon^{est} := T^{est} - \mathbb{E}[Y|do(\mathbf{x})]$  for  $est \in \{reg, pw, dml\}$ , are:

$$\begin{split} \epsilon^{reg} &= R_1 + O_{P_{x_1}} \left( \| \mu^1 - \mu_0^1 \| \right), \\ \epsilon^{pw} &= R_m + O_{P_{x_m}} (\| \pi^{(m-1)} - \pi_0^{(m-1)} \|), \\ \epsilon^{dml} &= \sum_{i=1}^m R_i + \sum_{i=2}^m O_{P_{x_i}} \left( \| \mu^i - \mu_0^i \| \| \pi^{i-1} - \pi_0^{i-1} \| \right). \end{split}$$

where  $R_i$  is a random variable such that  $\sqrt{n_i}R_i$  converges in distribution to a zero-mean normal random variable, where  $n_i := |D_{x_i}|$  for  $i \in \{1, \dots, m\}$ .

We highlight that the DML estimator  $T^{dml}$  in Def. 9 exhibits robustness property since the error  $\epsilon^{dml}$  is bounded in probability at a rate  $O_{P_{x_i}}(n^{-1/2})$  (for  $n = \min\{n_1, n_2, \cdots, n_m\}$ ) rate whenever  $\|\mu^i - \mu_0^i\|_{P_{x_i}} = O_{P_{x_i}}(n^{-1/4})$  and  $\|\pi^{i-1} - \pi_0^{i-1}\|_{P_{x_i}} = O_{P_{x_i}}(n^{-1/4})$ . Furthermore, the DML estimator displays the following multiply robustness property:

**Corollary 6** (Multiply robustness of the DML estimators (Corollary of Thm. 6)). *Suppose Asumptions* (2,7,8,9) *and* 

AC-MTI in Def. 7 hold. For  $i = 2, \dots, m-1$ , suppose either  $\pi^{i-1} = \pi_0^{i-1}$  or  $\mu^i = \mu_0^i$ . Then,  $T^{dml}$  in Def. 9 is an unbiased estimator of  $\mathbb{E}[Y|do(\mathbf{x})]$ .

#### 4.2. Combining Multiple Arbitrary Experiments

In this section, we generalize Task MTI to the case where multiple treatment effects  $\mathbb{E}[Y|do(\mathbf{x})]$  can be identified from arbitrary sets of experiments. We label this task as *'generalized multiple-teratment-interaction (gMTI)*':

**Task gMTI (Generalized MTI).** The task of *generaliazed MTI* composes of

• Input: Multiple sets of samples  $\{D_i\}_{i=1}^m$  following distributions  $\{P_{\text{rand}(Z_i)}(\mathbf{V})\}_{i=1}^m$ .

• **Query**: Estimation of  $\mathbb{E}[Y|do(\mathbf{x})]$ .

We note that Task gMTI is a generalization of Task MTI since X is not necessarily identical to  $\bigcup_i Z_i$ .

#### 4.2.1. Adjustment Criterion for GMTI

A graphical criterion for identifying the effect  $\mathbb{E}[Y|do(\mathbf{x})]$  is the following:

**Definition 10** (Adjustment criterion for gMTI (AC-gMTI)). Let  $\mathbf{Z} := \{Z_1, \dots, Z_m\} \subseteq \mathbf{X}$  denote the subset of treatments. Let  $\{\ell_i\}_{i=1}^m \subseteq \{1, 2, \dots, |\mathbf{X}|\}$ denote the index of  $\mathbf{Z}$ ; i.e.,  $\mathbf{Z} = \{X_{\ell_1}, \dots, X_{\ell_m}\}$ . Let  $\overline{X}_1 := \{X_j\}_{j \leq \ell_1}, \overline{X}_{m+1} := \{X_j\}_{j > \ell_m}$ , and  $\overline{X}_i := \{X_j\}_{\ell_{i-1} < j \leq \ell_i}$  for  $i = 2, 3, \dots, m$ . An ordered set  $\mathbf{A} := \{\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_m\}$  satisfies adjustment criterion for combining multiple experiments (AC-gMTI) w.r.t. ( $\mathbf{X}, Y$ ) in G if, for  $i = 1, 2, \dots, m-1$ ,

1. 
$$(\mathbf{A}_{i} \perp \mathbf{\overline{X}}^{>i-1} \setminus Z_{i} | \mathbf{\overline{X}}^{(i-1)}, \mathbf{A}^{(i-1)}, Z_{i})_{G_{\mathbf{\overline{X}}^{>i-1}}};$$
  
2.  $(Y \perp \mathbf{\overline{X}}_{i} | \mathbf{A}^{(i)}, \mathbf{\overline{X}}^{(i-1)}, \mathbf{\overline{X}}^{>i})_{G_{\underline{\overline{X}}_{i}, \mathbf{\overline{X}}^{>i}}};$  and  
3.  $(Y \perp \mathbf{\overline{X}}^{\geq m} \setminus Z_{m} | \mathbf{A}^{(m-1)}, \mathbf{\overline{X}}^{(m-1)}, Z_{m})_{G_{\overline{Zm}, \mathbf{\overline{X}}^{\geq m} \setminus Z_{m}}}$ 

We make the following positivity assumption:

Assumption 10 (Positivity Assumption for AC-gMTI).  $P_{z_m}(\overline{X}_m \setminus Z_m, \overline{X}_{m+1} | \mathbf{A}^{(m-1)}, \overline{\mathbf{X}}^{(m-1)})$  and  $\{P_{z_i}(\mathbf{A}_i | \mathbf{A}^{(i-1)}, \overline{\mathbf{X}}^{(i-1)}), P_{z_{i+1}}(\mathbf{A}_i | \mathbf{A}^{(i-1)}, \overline{\mathbf{X}}^{(i-1)})\}_{i=1}^{m-1}, \{P^{i+1}(\overline{X}_i | \mathbf{A}^{(i)}, \overline{\mathbf{X}}^{(i-1)})\}_{i=1}^{m-1}$  are strictly positive distributions  $\forall i \in \{1, \cdots, m\}, \forall z_i \in \mathfrak{D}_{Z_i}.$ 

Under AC-gMTI in Def. 10,  $\mathbb{E}[Y|do(\mathbf{x})]$  is identifiable and given as follows:

**Theorem 7** (Identification through AC-gMTI). Suppose AC-gMTI in Def. 10 and Assumption 10 hold. Then,

 $\mathbb{E}[Y|do(\mathbf{x})]$  is identifiable from  $\{P_{rand(Z_i)}(\mathbf{A}^{(i)}, \overline{\mathbf{X}}^{(i)})\}_{i=1}^{m}$ and given as follows. Denote

$$\mu_0^m \coloneqq \mathbb{E}_{P_{z_m}} \left[ Y | \mathbf{A}^{(m-1)}, \mathbf{X} \backslash Z_m \right]$$
$$\overline{\mu}_0^m \coloneqq \mathbb{E}_{P_{z_m}} \left[ Y | \mathbf{A}^{(m-1)}, \overline{\mathbf{x}}_{m-1:m+1}, \overline{\mathbf{X}}^{(m-2)} \right]$$
$$\mu_0^{m-1} \coloneqq \mathbb{E}_{P_{z_{m-1}}} \left[ \overline{\mu}_0^m | \mathbf{A}^{(m-2)}, \overline{\mathbf{X}}^{(m-2)} \right],$$

where  $\overline{X}_{m-1:m+1} \coloneqq \{\overline{X}_{m-1}, \overline{X}_m, \overline{X}_{m+1}\}$ . For  $i = m - 2, \cdots, 2$ ,

$$\mu_0^i \coloneqq \mathbb{E}_{P_{z_i}}\left[\mu^{i+1}(\mathbf{A}^{(i)}, \overline{x}_i, \overline{\mathbf{X}}^{(i-1)}) \middle| \mathbf{A}^{(i-1)}, \overline{\mathbf{X}}^{(i-1)}\right],$$

and 
$$\overline{\mu}_{0}^{i+1} \coloneqq \mu_{0}^{i+1}(\mathbf{A}^{(i)}, \overline{x}_{i}, \overline{\mathbf{X}}^{(i-1)})$$
. Then,  
$$\mathbb{E}\left[Y(\mathbf{x})\right] = \mathbb{E}_{P_{z_{1}}}\left[\overline{\mu}_{0}^{2}\right]. \tag{4}$$

For example, in Fig. 2b,  $\{\mathbf{A}_1, \mathbf{A}_2\}$  where  $\mathbf{A}_1 := \{C_1, W_1\}$ and  $\mathbf{A}_2 := \{C_2, W_2\}$  satisfies AC-gMTI criterion in Def. 10 w.r.t.  $\{\mathbf{X}, Y\}$  where  $\mathbf{X} := \{X_0, Z_1, Z_2, Z_3\}$ . Therefore, with the positivity in Assumption 10,  $\mathbb{E}[Y|do(\mathbf{x})]$  is identifiable from  $\{P_{z_i}(\mathbf{A}^{(i)}, \overline{\mathbf{X}}^{(i)})\}_{i=1}^m$  as in Eq. (4).

### 4.2.2. ESTIMATORS FOR AC-GMTI

We define the nuisance functionals for estimating the ACgMTI functional in Eq. (4) as follows:

**Definition 11** (Nuisances for AC-gMTI). Nuisance functions for AC-gMTI are defined as follows: For a fixed  $\mathbf{z} := \{z_1, \dots, z_m\} \in \mathfrak{D}_{\mathbf{Z}}$ , let  $\{\mu_0^i\}_{i=2}^m$  be the nuisances defined in Thm. 7. For  $i = 1, \dots, m-2, \ \pi_0^i := \frac{P_{z_i}(A_i|\mathbf{A}^{(i-1)}, \overline{\mathbf{X}}^{(i-1)})}{P_{z_m}(A_i, \overline{X}_i|\mathbf{A}^{(i-1)}, \overline{\mathbf{X}}^{(i-1)})},$  and  $\pi_0^{(i)} := \prod_{j=1}^i \pi_0^j (\mathbf{A}^{(j)}, \overline{\mathbf{X}}^{(j)})$ . Also,  $\mathbf{x}_1^{(m-1)} := \frac{P_{z_m}(A_{m-1}|\mathbf{A}^{(m-2)}, \overline{\mathbf{X}}^{(m-2)})}{P_{z_m}(A_{m-1}, \overline{X}_{m-1:m+1}|\mathbf{A}^{(m-2)}, \overline{\mathbf{X}}^{(m-2)})},$  and  $p_{i_0^{(m-1)}} := \pi_0^{(m-2)} \times \pi_0^{m-1},$  where  $\overline{X}_{m-1:m+1} := \{\overline{X}_{m-1}, \overline{X}_m, \overline{X}_{m+1}\}$ . For all  $i = 1, 2, \dots, m-1$ , we will use  $\pi^i(\mathbf{W}^{(i)} \mathbf{C}^{(i)} \mathbf{X}^{(i)}) > 0$  and  $\mu^i$  and  $\overline{\mu^i}$  to denote

will use  $\pi^i(\mathbf{W}^{(i)}, \mathbf{C}^{(i)}, \mathbf{X}^{(i)}) > 0$  and  $\mu^i$  and  $\overline{\mu}^i$  to denote estimated nuisances.

Now, we construct regression-based ('REG'), probability weighting ('PW'), and double/debiased machine learning ('DML') estimators:

**Definition 12 (AC-gMTI estimators).** Let  $D_i$  denote samples following  $P_{\text{rand}(Z_i)}(\mathbf{V})$  for  $i = 1, 2, \dots, m$ . For a fixed  $z_i \in \mathfrak{D}_{Z_i}$ , let  $D_{z_i}$  denote the subsamples of  $D_i$  such that  $Z_i = z_i$ . Let  $\mu^{m+1} \coloneqq Y$ . Let  $\mathbb{1}_{\mathbf{x}}^{i-1} \coloneqq \mathbb{1}_{\overline{\mathbf{x}}^{(i-1)}}(\overline{\mathbf{X}}^{(i-1)})$ .

	Case 1	Case 2	Case 3	Case 4
Fig. 1a (TTI)	0.5 Model 0.4 Model 0.4 Model 0.5 Model 0.6 Model 0.6 Model 0.7 Model		$\begin{array}{c} 0.5 \\ 0.4 \\ 0.3 \\ \hline \hline \\ 0.3 \\ \hline \\ \hline \\ 0.1 \\ \hline \\ 1 \\ -1 \\ -1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ $	$\begin{array}{c} 0.5 \\ 0.4 \\ 0.3 \\ \hline + 1 \\ - 1 \\ + 1 \\ - 1$
Fig. 1b (gTTI)	$\begin{array}{c} 0.5 \\ 0.4 \\ 0.3 \\ 0.1 \\$	$\begin{array}{c} 0.5 \\ 0.4 \\ 0.3 \\ 0.2 \\ 0.1 \\ 0.2 \\ 0.1 \\ 0.2 \\$	0.5 + - + + + + + + + + + + + + + + + + +	$\begin{array}{c} 0.5 \\ 0.4 \\ 0.3 \\ 0.2 \\ 0.1 \\ \begin{array}{c} \end{array} \\ \end{array} \\ \begin{array}{c} \end{array} \\ \end{array} \\ \begin{array}{c} \end{array} \\ \begin{array}{c} \end{array} \\ \end{array} \\ \begin{array}{c} \end{array} \\ \begin{array}{c} \end{array} \\ \begin{array}{c} \end{array} \\ \end{array} \\ \end{array} \\ \end{array} \\ \end{array} \\ \begin{array}{c} \end{array} \\ \end{array} $
Fig. 2a (MTI)	$\begin{array}{c} 0.5 \\ 0.4 \\ 0.3 \\ 0.4 \\ 0.3 \\ 0.4 \\ 0.5 \\ 0.4 \\ 0.5 \\ 0.4 \\ 0.5 \\ 0.4 \\ 0.5 \\ 0.4 \\ 0.5 \\ 0.4 \\ 0.5 \\ 0.4 \\ 0.5 \\ 0.4 \\ 0.5 \\ 0.4 \\ 0.5 \\ 0.4 \\ 0.5 \\ 0.4 \\ 0.5 \\ 0.4 \\ 0.5 \\ 0.4 \\ 0.5 \\ 0.4 \\ 0.5 \\ 0.4 \\ 0.5 \\ 0.4 \\ 0.5 \\ 0.4 \\ 0.5 \\ 0.4 \\ 0.5 \\ 0.4 \\$	$\begin{array}{c} 0.5 \\ 0.4 \\ 0.3 \\ 0.2 \\ 0.1 \\ 0.1 \\ 0.1 \\ 0.2 \\ 0.2 \\ 0.2 \\ 0.2 \\ 0.3 \\ 0.2 \\ 0.3 \\ 0.2 \\ 0.3 \\ 0.4 \\$	$\begin{array}{c} 0.5 \\ 0.4 \\ 0.3 \\ 0.4 \\ 0.3 \\ 0.4 \\ 0.5 \\ 0.4 \\ 0.5 \\$	$\begin{array}{c} 0.5 \\ 0.4 \\ 0.3 \\ 0.2 \\ 0.1 \\ \hline 1 \\ 1 \\$
Fig. 2b (gMTI)		$\begin{array}{c} 0.5 \\ 0.4 \\ 0.3 \\ 0.2 \\ 0.1 \\ 0.2 \\ 0.1 \\ 0.2 \\$		$ \begin{array}{c} 0.5 \\ 0.4 \\ 0.3 \\ 0.1 $

Figure 3: AAE Plots for Figs. (1a, 1b, 2a, 2b) for Cases  $\{1,2,3,4\}$  depicted in the Experimental Setup section. The *x*-axis and *y*-axis are the number of samples and AAE, respectively. Plots can be zoomed in.

Then {REG, PW, DML} estimators are defined as:

$$T^{reg} \coloneqq \mathbb{E}_{D_{z_1}} \left[ \mu^2(A_1, \overline{x}_1) \right],$$
  

$$T^{pw} \coloneqq \mathbb{E}_{D_{z_m}} \left[ \pi^{(m-1)}(\mathbf{A}^{(m-1)}, \mathbf{X}) \mathbb{1}_{\mathbf{x}}(\mathbf{X}) Y \right],$$
  

$$T^{dml} \coloneqq \sum_{i=2}^m \mathbb{E}_{D_{z_i}} \left[ \pi^{(i-1)} \mathbb{1}_{\mathbf{x}}^{i-1} \{ \overline{\mu}^{i+1} - \mu^i \} \right] + \mathbb{E}_{D_{z_1}} \left[ \overline{\mu}^2 \right]$$

We assume that nuisances can be estimated  $L_2$  consistently.

Assumption 11 ( $L_2$  consistency of nuisances). Estimated nuisances  $\{\mu^i\}_{i=2}^m$  and  $\{\pi^i\}_{i=1}^{m-1}$  are  $L_2$  consistent; specifically,

$$\begin{split} \|\mu^{i+1} - \mu_0^{i+1}\|_{P_{z_i}} &= o_{P_{z_i}}(1), \; \forall i \in \{1, 2, \cdots, m-1\} \\ \|\mu^i - \mu_0^i\|_{P_{z_i}} &= o_{P_{z_i}}(1), \; \forall i \in \{2, \cdots, m\} \\ \|\pi^i - \pi^i\|_{P_{z_{i+1}}} &= o_{P_{z_{i+1}}}(1), \; \forall i \in \{1, \cdots, m-1\}. \end{split}$$

**Theorem 8 (Error analysis of the AC-gMTI estimators).** Under Assumptions (2,10,11) and AC-gMTI in Def. 10, the errors of the estimators in Def. 12, denoted  $\epsilon^{est} := T^{est} - \mathbb{E}[Y|do(\mathbf{x})]$  for  $est \in \{reg, pw, dml\}$ , are:

$$\begin{aligned} \epsilon^{reg} &= R_1 + O_{P_{z_1}}(\|\mu^1 - \mu_0^1\|), \\ \epsilon^{pw} &= R_m + O_{P_{z_m}}(\|\pi^{(m-1)} - \pi_0^{(m-1)}\|), \\ \epsilon^{dml} &= \sum_{i=1}^m R_i + \sum_{i=2}^m O_{P_{z_i}}(\|\mu^i - \mu_0^i\|\|\pi^{i-1} - \pi_0^{i-1}\|), \end{aligned}$$

where  $R_i$  is a variable such that  $\sqrt{n_i}R_i$  converges in distribution to a zero-mean normal random variable, where  $n_i \coloneqq |D_i|$  for  $i \in \{1, \dots, m\}$ .

We highlight that the DML estimator  $T^{dml}$  in Def. 12 exhibits robustness property since the error  $\epsilon^{dml}$  is bounded in probability at a rate  $O_{Pz_i}(n^{-1/2})$  (for  $n = \min\{n_1, n_2, \cdots, n_m\}$ ) rate whenever  $\|\mu^i - \mu_0^i\|_{Pz_i} = O_{Pz_i}(n^{-1/4})$  and  $\|\pi^{i-1} - \pi_0^{i-1}\|_{Pz_i} = O_{Pz_i}(n^{-1/4})$ . Furthermore, the DML estimator displays the following multiply robustness property:

**Corollary 8** (Multiply robustness of the DML estimators (Corollary of Thm. 8)). Suppose Assumptions (2,10,11) and AC-gMTI in Def. 10 hold. For  $i = 2, \dots, m-1$ , suppose either  $\pi^{i-1} = \pi_0^{i-1}$  or  $\mu^i = \mu_0^i$ . Then,  $T^{dml}$  in Def. 12 is an unbiased estimator of  $\mathbb{E}[Y|do(\mathbf{x})]$ .

## **5.** Experiments

In this section, we demonstrate the proposed estimators in Defs. (3,6) for combining two experiments and those in Defs. (9,12) for combining multiple experiments. We first compared the estimators on synthetic data. The synthetic data analysis provided evidence of fast convergence and doubly robustness behaviors of the proposed estimators. We then applied the estimators on a real-world dataset Project STAR (Krueger & Whitmore, 2001; Schanzenbach, 2006).

We will use the notations  $T^{\text{est}}(\mathbf{x})$  for  $\text{est} \in \{\text{reg}, \text{pw}, \text{dml}\}$ 

to represent the estimators {REG, PW, DML} for the joint treatment effects  $\mathbb{E}[Y|do(\mathbf{x})]$ . To assess the quality of each estimator est  $\in$  {reg, pw, dml}, we compute the *average absolute error* (AAE) as AAE<sup>est</sup> =  $\frac{1}{|\mathfrak{D}_{\mathbf{x}}|} \sum_{\mathbf{x} \in \mathfrak{D}_{\mathbf{X}}} |T^{\text{est}}(\mathbf{x}) - \mathbb{E}[Y|do(\mathbf{x})]|$ , where  $|\mathfrak{D}_{\mathbf{X}}|$  is the cardinality of  $\mathfrak{D}_{\mathbf{X}}$ . The box plots for AAE values for each estimator are referred to as "AAE-plots". We used XGBoost (Chen & Guestrin, 2016) to estimate nuisances.

#### **5.1. Sythetic Dataset Analysis**

**Experimental Setup.** We ran 100 simulations for each  $N = \{2000, 4000, 6000, 8000, 10000\}$  where N is the sample size. We measure the AAE<sup>est</sup> for each four scenarios: (**Case 1**) there were no noises nor misspecification in estimating nuisances; (**Case 2**); The 'converging noise'  $\epsilon$ , decaying at a  $N^{-\alpha}$  rate (i.e.,  $\epsilon \sim \text{Normal}(N^{-\alpha}, N^{-2\alpha}))$  for  $\alpha = 1/4$ , was added in estimating nuisances; (**Case 3**) Nuisances  $\{\mu^i\}_{i=2}^m$  are wrongly estimated; (**Case 4**)  $\{\pi^i\}_{i=1}^{m-1}$  are wrongly estimated. Case 2 is a scenario to highlight fast convergence property of the DML estimator implied in Thms. (2,4,6,8) where  $T^{\text{dml}}$  converge at  $n^{-1/2}$ -rate) when other estimators  $T^{\text{reg}}, T^{\text{pw}}$  converge at  $n^{-1/4}$ -rate. Cases {3,4} are designed to exhibit doubly robustness property of  $T^{\text{dml}}$  formalized in Corolaries (2,4,6,8). Details of the experiments are provided in Appendix E.

**Experimental Results.** The AAE plots for all cases are presented in Fig. 3. All {DML, REG, IPW} estimators converges in Case 1 as the sample size grows. In Case 2 where the estimated nuisances are controlled to converge at  $n^{-1/4}$  rate, the DML estimators  $T^{\text{dml}}$  outperform the other two estimators by achieving a fast convergence. This result corroborates the robustness property in Thms. (2,4,6,8). In Cases (3,4) where the estimated nuisances for  $\{\mu^i\}_{i=2}^m$  or  $\{\pi^i\}_{i=1}^{m-1}$  are wrongly specified, the DML estimator  $T^{\text{dml}}$  converges while other estimators fail to converge. This result corrobrates the doubly robustness property in Coros. (2,4,6,8).

### 5.2. Project STAR Dataset

This section provides an overview of the analysis conducted on the Project STAR dataset. The detailed procedures and results can be found in Appendix D. Project STAR investigated the impact of teacher/student ratios on academic achievement for students in kindergarten through third grade. The dataset, denoted as D, includes variables such as class size for kindergarten (X<sub>1</sub>), academic outcome in kindergarten (W), class size for third grade (X<sub>2</sub>), academic outcome in third grade (Y), and pre-treatment variables (C). Project STAR is a longitudinal experimental study, with samples for the variables {C, X<sub>1</sub>, W} following a distribution  $P_{\text{rand}(X_1)}(C, X_1, W)$  and samples for the variables {C, X<sub>1</sub>, W, X<sub>2</sub>, Y} following a distribution  $P_{\text{rand}(X_1, X_2)}(C, X_1, W, X_2, Y)$ . We assume that the SCM



Figure 4: A graph and the AAE-plot for Project STAR.

 $\mathcal{M}$  generating D induces a causal graph in Figure 4a.

**Experimental Setup.** To simulate Task TTI, we generate two datasets  $D_1$  and  $D_2$  from the original dataset D.  $D_1$ is a random subsample of D with only  $\{C, X_1, W\}$  and follows  $P_{\text{rand}(X_1)}(C, X_1, W)$ .  $D_2$  is constructed by resampling from D in a way that the confounding bias between  $X_1$  and W presents, following  $P_{\text{rand}(X_2)}(C, X_1, W, X_2, Y)$ . We conducted 100 simulations by generating new instances of  $D_1$  and  $D_2$  to create the AAE plot. Estimators were constructed solely from  $D_1$  and  $D_2$ , with D used exclusively to construct the ground-truth estimate. In this empirical study, we aim to study  $\mathbb{E}[Y|do(x_1, x_2)]$ .

**Experimental Results.** We evaluated the AAE<sup>est</sup> of estimators  $T^{\text{est}}$  for est  $\in$  {reg, pw, dml}. The AAE plot is in Fig. 4b. Our findings indicate that the proposed estimators, especially  $T^{\text{reg}}$  and  $T^{\text{dml}}$ , consistently provided reliable estimates for the ground-truth quantity. Finally, additional simulations in Appendix D demonstrated the fast convergence and doubly robustness properties of  $T^{\text{dml}}$ .

## 6. Conclusions

We proposed a set of identification conditions for estimating joint causal effects  $\mathbb{E}[Y|do(\mathbf{x})]$  by combining multiple marginal experiments (Thms (1,3,5,7)). Next, we developed corresponding estimators (Defs (3,6,9,12)) that are robust to model misspecification and slow convergence in learning nuisance (Thms (2,4,6,8) and Coros (2,4,6,8)) for Tasks (TTI,gTTI,MTI,gMTI). Our experimental results corroborate theories. We hope this work can help data scientists to estimate joint treatment effects from multiple experiments in a more principled and efficient manner.

## Acknowledgement

We thank Iván Díaz and the reviewers for their valuable feedback and assistance in improving this paper. Elias Bareinboim and Yonghan Jung were supported in part by funding from the Alfred P. Sloan Foundation, NSF, ONR, AFOSR, DoE, Amazon, and JP Morgan. Jin Tian was partially supported by NSF grant IIS-2231797.

## References

- Ajjan, R. A. and Grant, P. J. Cardiovascular disease prevention in patients with type 2 diabetes: The role of oral anti-diabetic agents. *Diabetes and Vascular Disease Research*, 3(3):147–158, 2006.
- Athey, S., Chetty, R., Imbens, G. W., and Kang, H. The surrogate index: Combining short-term proxies to estimate long-term treatment effects more rapidly and precisely. Technical report, National Bureau of Economic Research, 2019.
- Athey, S., Chetty, R., and Imbens, G. Combining experimental and observational data to estimate treatment effects on long term outcomes. *arXiv preprint arXiv:2006.09676*, 2020.
- Bareinboim, E. and Pearl, J. Causal inference by surrogate experiments: z-identifiability. In *In Proceedings of the* 28th Conference on Uncertainty in Artificial Intelligence, pp. 113–120. AUAI Press, 2012a.
- Bareinboim, E. and Pearl, J. Transportability of causal effects: Completeness results. In *Proceedings of the 26th AAAI Conference on Artificial Intelligence*, pp. 698–704, 2012b.
- Bareinboim, E. and Pearl, J. Causal inference and the datafusion problem. *Proceedings of the National Academy of Sciences*, 113(27):7345–7352, 2016.
- Bareinboim, E., Correa, J. D., Ibeling, D., and Icard, T. On pearl's hierarchy and the foundations of causal inference. In *Probabilistic and causal inference: the works of judea pearl*, pp. 507–556. 2022.
- Bhattacharya, R., Nabi, R., and Shpitser, I. Semiparametric inference for causal effects in graphical models with hidden variables. *Journal of Machine Learning Research*, 23:1–76, 2022.
- Chen, T. and Guestrin, C. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 785–794, 2016.
- Chernozhukov, V., Chetverikov, D., Demirer, M., Duflo, E., Hansen, C., Newey, W., and Robins, J. Double/debiased machine learning for treatment and structural parameters: Double/debiased machine learning. *The Econometrics Journal*, 21(1), 2018.
- Colnet, B., Mayer, I., Chen, G., Dieng, A., Li, R., Varoquaux, G., Vert, J.-P., Josse, J., and Yang, S. Causal inference methods for combining randomized trials and observational studies: a review. arXiv preprint arXiv:2011.08047, 2020.

- Correa, J., Lee, S., and Bareinboim, E. Nested counterfactual identification from arbitrary surrogate experiments. *Advances in Neural Information Processing Systems*, 34, 2021.
- Crump, R. K., Hotz, V. J., Imbens, G. W., and Mitnik, O. A. Dealing with limited overlap in estimation of average treatment effects. *Biometrika*, 96(1):187–199, 2009.
- Dahabreh, I. J., Robertson, S. E., Tchetgen, E. J., Stuart, E. A., and Hernán, M. A. Generalizing causal inferences from individuals in randomized trials to all trial-eligible individuals. *Biometrics*, 75(2):685–694, 2019.
- Degtiar, I. and Rose, S. A review of generalizability and transportability. *arXiv preprint arXiv:2102.11904*, 2021.
- Egami, N. and Imai, K. Causal interaction in factorial experiments: Application to conjoint analysis. *Journal of the American Statistical Association*, 2018.
- Ferrannini, E. and Cushman, W. C. Diabetes and hypertension: the bad companions. *The Lancet*, 380(9841): 601–610, 2012.
- Gazzard, B., Clark, R., Borirakchanyavat, V., and Williams, R. A controlled trial of heparin therapy in the coagulation defect of paracetamol-induced hepatic necrosis. *Gut*, 15 (2):89–93, 1974.
- Gentzel, A. M., Pruthi, P., and Jensen, D. How and why to use experimental data to evaluate methods for observational causal inference. In *International Conference on Machine Learning*, pp. 3660–3671. PMLR, 2021.
- Hansson, L., Zanchetti, A., Carruthers, S. G., Dahlöf, B., Elmfeldt, D., Julius, S., Ménard, J., Rahn, K. H., Wedel, H., Westerling, S., et al. Effects of intensive bloodpressure lowering and low-dose aspirin in patients with hypertension: principal results of the hypertension optimal treatment (hot) randomised trial. *The Lancet*, 351 (9118):1755–1762, 1998.
- Hansson, L., Lindholm, L. H., Ekbom, T., Dahlöf, B., Lanke, J., Scherstén, B., Wester, P., Hedner, T., de Faire, U., Group, S.-H.-. S., et al. Randomised trial of old and new antihypertensive drugs in elderly patients: cardiovascular mortality and morbidity the swedish trial in old patients with hypertension-2 study. *The Lancet*, 354(9192):1751– 1756, 1999.
- Hill, J. L. Bayesian nonparametric modeling for causal inference. *Journal of Computational and Graphical Statistics*, 20(1):217–240, 2011.
- Imbens, G., Kallus, N., Mao, X., and Wang, Y. Long-term causal inference under persistent confounding via data combination. arXiv preprint arXiv:2202.07234, 2022.

- Jung, Y., Tian, J., and Bareinboim, E. Learning causal effects via weighted empirical risk minimization. *Advances in Neural Information Processing Systems*, 33, 2020.
- Jung, Y., Tian, J., and Bareinboim, E. Estimating identifiable causal effects on markov equivalence class through double machine learning. In *Proceedings of the 38th International Conference on Machine Learning*, 2021a.
- Jung, Y., Tian, J., and Bareinboim, E. Estimating identifiable causal effects through double machine learning. In *Proceedings of the 35th AAAI Conference on Artificial Intelligence*, 2021b.
- Jung, Y., Kasiviswanathan, S., Tian, J., Janzing, D., Blöbaum, P., and Bareinboim, E. On measuring causal contributions via do-interventions. In *International Conference on Machine Learning*, pp. 10476–10501. PMLR, 2022.
- Kennedy, E. H., Balakrishnan, S., G'Sell, M., et al. Sharp instruments for classifying compliers and generalizing causal effects. *Annals of Statistics*, 48(4):2008–2030, 2020.
- Klaassen, C. A. Consistent estimation of the influence function of locally asymptotically linear estimators. *The Annals of Statistics*, pp. 1548–1562, 1987.
- Krueger, A. B. and Whitmore, D. M. The effect of attending a small class in the early grades on college-test taking and middle school test results: Evidence from project star. *The Economic Journal*, 111(468):1–28, 2001.
- Lee, S. and Bareinboim, E. Causal effect identifiability under partial-observability. In *Proceedings of the 37th International Conference on Machine Learning*, 2020.
- Lee, S., Correa, J. D., and Bareinboim, E. General identifiability with arbitrary surrogate experiments. In *Proceed*ings of the 35th Conference on Uncertainty in Artificial Intelligence. AUAI Press, 2019.
- Lee, S., Correa, J., and Bareinboim, E. General transportability–synthesizing observations and experiments from heterogeneous domains. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 10210–10217, 2020.
- Lesko, S. M. and Mitchell, A. A. An assessment of the safety of pediatric ibuprofen: a practitioner-based randomized clinical trial. *Jama*, 273(12):929–933, 1995.
- Louizos, C., Shalit, U., Mooij, J. M., Sontag, D., Zemel, R., and Welling, M. Causal effect inference with deep latent-variable models. *Advances in neural information processing systems*, 30, 2017.

- Moore, N., Pollack, C., and Butkerait, P. Adverse drug reactions and drug–drug interactions with over-the-counter nsaids. *Therapeutics and clinical risk management*, 11: 1061, 2015.
- Parbhoo, S., Bauer, S., and Schwab, P. Ncore: Neural counterfactual representation learning for combinations of treatments. arXiv preprint arXiv:2103.11175, 2021.
- Pearl, J. Causality: Models, Reasoning, and Inference. Cambridge University Press, New York, 2000. 2nd edition, 2009.
- Robins, J. M. and Ritov, Y. Toward a curse of dimensionality appropriate (coda) asymptotic theory for semi-parametric models. *Statistics in medicine*, 16(3):285–319, 1997.
- Rotnitzky, A., Robins, J., and Babino, L. On the multiply robust estimation of the mean of the g-functional. *arXiv preprint arXiv:1705.08582*, 2017.
- Saengkyongam, S. and Silva, R. Learning joint nonlinear effects from single-variable interventions in the presence of hidden confounders. In *Conference on Uncertainty in Artificial Intelligence*, pp. 300–309. PMLR, 2020.
- Schanzenbach, D. W. What have researchers learned from project star? *Brookings papers on education policy*, (9): 205–228, 2006.
- Shi, X., Pan, Z., and Miao, W. Data integration in causal inference. Wiley Interdisciplinary Reviews: Computational Statistics, pp. e1581, 2022.
- Stock, J. H., Watson, M. W., et al. Introduction to econometrics, volume 104. Addison Wesley Boston, 2003.
- Tam, V., Patel, N., Turcotte, M., Bossé, Y., Paré, G., and Meyre, D. Benefits and limitations of genome-wide association studies. *Nature Reviews Genetics*, 20(8):467–484, 2019.
- VanderWeele, T. J. and Knol, M. J. A tutorial on interaction. *Epidemiologic methods*, 3(1):33–72, 2014.
- Zhang, J. and Bareinboim, E. Near-optimal reinforcement learning in dynamic treatment regimes. Advances in Neural Information Processing Systems, 32, 2019.
- Zheng, W. and van der Laan, M. J. Cross-validated targeted minimum-loss-based estimation. In *Targeted Learning*, pp. 459–474. Springer, 2011.

# Supplement to "Estimating Joint Treatment Effects by Combining Multiple Experiments"

## Contents

1	<b>Intr</b> 1.1	oduction Related Work	1 2
2	Prel	liminaries	2
3	Con	nbining Two Experiments	2
-	3.1	Treatment-Treatment Interaction	3
	5.1	3.1.1 Adjustment Criterion for TTI (AC-TTI)	3
		3.1.2 Estimators for AC-TTI	3
	3.2	Combining Two Arbitrary Experiments	4
		3.2.1 Adjustment Criterion for gTTI (AC-gTTI)	4
		3.2.2 Estimators for AC-gTTI	5
4	Con	nbining Multiple ( $\geq$ 2) Experiments	5
	4.1	Multiple Treatment Interaction	5
		4.1.1 Adjustment Criterion for MTI (AC-MTI)	5
		4.1.2 Estimators for AC-MTI	6
	4.2	Combining Multiple Arbitrary Experiments	7
		4.2.1 Adjustment Criterion for gMTI	7
		4.2.2 Estimators for AC-gMTI	7
5	Exp	eriments	8
	5.1	Sythetic Dataset Analysis	9
	5.2	Project STAR Dataset	9
6	Con	clusions	9
A	Prel	liminaries	14
	A.1	The Axioms of Structural Counterfactuals	14
B	Ider	ntification based on Potential Outcome Framework	14
	<b>B</b> .1	Treatment-Treatment Interaction based on Potential Outcome Framework	14
	B.2	Combining Two Experiments based on Potential Outcome Framework	16
	B.3	Multiple Treatment Interaction based on Potential Outcome Framework	18
	<b>B.</b> 4	Combining Multiple Experiments based on Potential Outcome Framework	22

С	Proc	ofs	27
	<b>C</b> .1	Preliminaries	27
	<b>C</b> .2	Proof of Theorem 1	28
	C.3	Proof of Theorem 2 and Corollary 2	29
	<b>C.</b> 4	Proof of Theorem 3	33
	C.5	Proof of Theorem 4 and Corollary 4	33
	C.6	Proof of Theorem 5	37
	C.7	Proof of Theorem 6 and Corollary 6	38
	C.8	Proof of Theorem 7	47
	C.9	Proof of Theorem 8 and Corollary 8	50
D	Proj	ject STAR: Estimating Joint Effects of Class Sizes to Academic Outcomes	59
D E	Proj Deta	ject STAR: Estimating Joint Effects of Class Sizes to Academic Outcomes ails of Experiments	59 62
D E	Proj Deta E.1	ject STAR: Estimating Joint Effects of Class Sizes to Academic Outcomes ails of Experiments Designs of Simulations	<b>59</b> <b>62</b> 62
D E	Proj Deta E.1	ject STAR: Estimating Joint Effects of Class Sizes to Academic Outcomes         ails of Experiments         Designs of Simulations         E.1.1         Task TTI	<ul> <li>59</li> <li>62</li> <li>62</li> <li>62</li> </ul>
D E	Proj Deta E.1	ject STAR: Estimating Joint Effects of Class Sizes to Academic Outcomes         ails of Experiments         Designs of Simulations         E.1.1         Task TTI         E.1.2         Task gTTI	<b>59</b> <b>62</b> 62 62 64
D	Proj Deta E.1	ject STAR: Estimating Joint Effects of Class Sizes to Academic Outcomes     ails of Experiments   Designs of Simulations   E.1.1   Task TTI   E.1.2   Task gTTI   E.1.3   Task MTI	<b>59</b> <b>62</b> 62 62 64 65
D E	Proj Deta E.1	ject STAR: Estimating Joint Effects of Class Sizes to Academic Outcomes         ails of Experiments         Designs of Simulations         E.1.1       Task TTI         E.1.2       Task gTTI         E.1.3       Task MTI         E.1.4       Task gMTI	<b>59</b> <b>62</b> 62 62 64 65 67
D	Proj Deta E.1	ject STAR: Estimating Joint Effects of Class Sizes to Academic Outcomes     ails of Experiments   Designs of Simulations   E.1.1   Task TTI   E.1.2   Task gTTI   E.1.3   Task MTI   E.1.4   Task gMTI   E.1.5   Data Generation for Project STAR	<b>59</b> <b>62</b> 62 64 65 67 69
D E F	Proj Deta E.1	ject STAR: Estimating Joint Effects of Class Sizes to Academic Outcomes     ails of Experiments   Designs of Simulations   E.1.1   Task TTI   E.1.2   Task gTTI   E.1.3   Task MTI   E.1.4   Task gMTI   E.1.5   Data Generation for Project STAR   cussion on Relaxation of Shared Covariates Assumptions	<ul> <li>59</li> <li>62</li> <li>62</li> <li>64</li> <li>65</li> <li>67</li> <li>69</li> <li>74</li> </ul>
D E F	Proj Deta E.1 Disc F.1	ject STAR: Estimating Joint Effects of Class Sizes to Academic Outcomes         ails of Experiments         Designs of Simulations         E.1.1         Task TTI         E.1.2         Task gTTI         E.1.3         Task MTI         E.1.4         Task gMTI         E.1.5         Data Generation for Project STAR         Cussion on Relaxation of Shared Covariates Assumptions         On Assumption 4.	<ul> <li>59</li> <li>62</li> <li>62</li> <li>64</li> <li>65</li> <li>67</li> <li>69</li> <li>74</li> <li>74</li> </ul>

## **A. Preliminaries**

In this section, we present the preliminary concepts and notation used in this paper. Let W and X be sets of variables that are subsets of V, which are induced from the structural causal model (SCM)  $\mathcal{M}$ . Given a realization x of X, we denote W(x) as the counterfactual W under the hypothetical scenario where X takes the value x. In other words, W(x) represents a random vector generated by the submodel  $\mathcal{M}_x$ .

## A.1. The Axioms of Structural Counterfactuals

**Definition A.1** (The Axioms of Structural Counterfactuals (Pearl, 2000, Chapter 7.3.1)). For any three sets of endogenous variables  $\mathbf{X}, \mathbf{Y}, \mathbf{W}$  in a causal model and  $\mathbf{x}, \mathbf{w} \in \mathfrak{D}_{\mathbf{X}, \mathbf{W}}$ , the following holds:

- Composition:  $\mathbf{W}(\mathbf{x}) = \mathbf{w} \implies \mathbf{Y}(\mathbf{x}, \mathbf{w}) = \mathbf{Y}(\mathbf{x}).$
- Effectiveness:  $\mathbf{X}(\mathbf{w}, \mathbf{x}) = \mathbf{x}$ .
- Reversibility:  $\mathbf{Y}(\mathbf{x}, \mathbf{w}) = \mathbf{y}$  and  $\mathbf{W}(\mathbf{x}, \mathbf{y}) = \mathbf{w} \implies \mathbf{Y}(\mathbf{x}) = \mathbf{y}$ .

**Theorem A.1 (Soundness and Completeness of the Axioms** (Pearl, 2000, Theorems {7.3.3, 7.3.6})). *The Axioms of structural counterfactuals in Def. A.1 are sound and complete for all causal models.* 

**Remark 1** ((Pearl, 2000, page 230)). In the recursive (acyclic) system, Reversibility is followed from Composition. Therefore, Composition and Effectiveness are sound and complete.

**Definition A.2** (Potential Response, Counterfactuals (Pearl, 2000, Def. 7.1.4)). Let  $(\mathbf{X}, \mathbf{Y}) \subseteq \mathbf{V}$  generated by the SCM  $\mathcal{M}$ . The counterfactual of  $\mathbf{Y}$  at  $\mathbf{x}$ , denoted  $\mathbf{Y}(\mathbf{x})$ , is the variable  $\mathbf{Y}$  induced by the submodel  $\mathcal{M}_{\mathbf{x}}$ .

In this section, we will use  $P_{\text{rand}(\mathbf{X})}(\mathbf{V}) \coloneqq \{P(\mathbf{V}(\mathbf{x}))\}_{\mathbf{x}\in\mathfrak{D}_{\mathbf{X}}}$  to denote a collection of counterfactual distributions  $P(\mathbf{V}(\mathbf{x}))$  over all possible realizations  $\mathbf{x}\in\mathfrak{D}_{\mathbf{X}}$ . We will denote the density of P as p.

## **B. Identification based on Potential Outcome Framework**

In this section, we introduce more results on identifying causal effects. For  $\mathbf{x} \in \mathfrak{D}_{\mathbf{X}}$ , we use  $(\mathbf{W} \setminus \mathbf{X})(\mathbf{x})$  to denote the counterfactual of  $\mathbf{W} \setminus \mathbf{X}$  at  $\mathbf{x}$ . We use  $(\mathbf{w} \setminus \mathbf{x})(\mathbf{x})$  to denote its realization.

## **B.1. Treatment-Treatment Interaction based on Potential Outcome Framework**

We present a sufficient identification criterion for estimating treatment-treatment interactions using potential outcome frameworks based on two marginal experiments.

**Definition B.1 (Adjustment criterion for treatment-treatment-interaction – Potential Outcome (AC-TTI-PO)).** A set of variables  $\{C_1, W\}$  is said to satisfy the adjustment criterion for treatment-treatment interaction (AC-TTI) w.r.t. discrete treatments  $(X_1, X_2)$  and the outcome Y from two sets of distributions  $P_{\text{rand}(X_1)}(C_1, X_1, W)$  and  $P_{\text{rand}(X_2)}(C_1, X_1, W, X_2, Y)$  if

- 1.  $W(x_1, x_2) = W(x_1), C_1(x_1, x_2) = C_1(x_1)$ ; i.e., the outcome W and the covariate  $C_1$  is invariant of the second intervention  $X_2 = x_2$ .
- 2.  $Y(x_1, x_2) \perp X_1(x_2) | W_1(x_1, x_2), C_1(x_1, x_2);$  i.e., the first intervention  $X_1 = x_1$  is non-informative to the joint experimental outcome  $Y(x_1, x_2)$  given covariates  $C_1(x_1, x_2)$  and the first outcome  $W(x_1, x_2)$ .

The treatment-treatment interaction can be identified as follow:

**Theorem B.1 (Identification through AC2-TTI-PO).** Suppose the condition AC-TTI-PO in Def. B.1 holds. For any fixed  $x_1, x_2 \in \mathfrak{D}_{X_1,X_2}$ , define  $P^1 \coloneqq P_{x_1} \in P_{rand(X_1)}$  and  $P^2 \coloneqq P_{x_2} \in P_{rand(X_2)}$ . Assume the following positivity condition

holds for  $\forall x_1, x_2, w, c_1 \in \mathfrak{D}_{X_1, X_2, W, C_1}$ :

$$\frac{p^2(w,c_1)}{p^2(w,c_1)}p^2(X_1 = x_1|w,c_1) > 0.$$
(B.1)

Then, the query  $\mathbb{E}[Y(x_1, x_2)]$  is identifiable from two distributions  $P^1, P^2$  and given as follow:

$$\mathbb{E}[Y(x_1, x_2)] = \mathbb{E}_{P^1} \left[ \mathbb{E}_{P^2} \left[ Y | W, C_1, X_1 = x_1 \right] \right].$$
(B.2)

Proof of Theorem **B.1**.

$$\begin{split} \mathbb{E}\left[Y(x_1, x_2)\right] &= \mathbb{E}\left[\mathbb{E}\left[Y(x_1, x_2)|W(x_1, x_2), C_1(x_1, x_2)\right]\right] \\ &= \mathbb{E}\left[\mathbb{E}\left[Y(x_1, x_2)|W(x_1, x_2), C_1(x_1, x_2), X_1(x_2) = x_1\right]\right] \\ &= \mathbb{E}\left[\mathbb{E}\left[Y(x_1, x_2)|W(x_1), C_1(x_1), X_1(x_2) = x_1\right]\right] \\ &= \mathbb{E}\left[\mathbb{E}\left[Y(x_2)|W(x_1), C_1(x_1), X_1(x_2) = x_1\right]\right] \\ &= \mathbb{E}_{P^1}\left[\mathbb{E}\left[Y(x_2)|W, C_1, X_1(x_2) = x_1\right]\right], \end{split}$$

where

•  $\stackrel{1}{=}$  holds by the given condition that  $Y(x_1, x_2) \perp X_1(x_2) | W(x_1, x_2), C_1(x_1, x_2)$  and the positivity  $P(X_1(x_2) = x_1 | W(x_1, x_2), C_1(x_1, x_2)) > 0$  for any  $x_1, x_2$ . To witness, it suffices to show that  $p_{X_1(x_2)|W(x_1, x_2), C_1(x_1, x_2)}(x_1 | w, c_1) > 0$ .

$$\begin{split} p_{X_{1}(x_{2})|W(x_{1},x_{2}),C_{1}(x_{1},x_{2})}(x_{1}|w,c_{1}) \\ &= \frac{p_{X_{1}(x_{2}),W(x_{1},x_{2}),C_{1}(x_{1},x_{2})}(x_{1},w,c_{1})}{p_{W(x_{1},x_{2}),C_{1}(x_{1},x_{2})}(w,c_{1})} \\ \frac{1a}{2} \frac{p_{X_{1}(x_{2}),W(x_{1},x_{2}),C_{1}(x_{1},x_{2})}(x_{1},w,c_{1})}{p_{W(x_{1}),C_{1}(x_{1})}(w,c_{1})} \\ \frac{1b}{2} \frac{p_{X_{1}(x_{2}),W(x_{2}),C_{1}(x_{2})}(x_{1},w,c_{1})}{p_{W(x_{1}),C_{1}(x_{1})}(w,c_{1})} \\ &= \frac{p_{X_{1}(x_{2}),W(x_{2}),C_{1}(x_{2})}(x_{1},w,c_{1})}{p_{W(x_{1}),C_{1}(x_{1})}(w,c_{1})} \frac{p_{W(x_{2}),C_{1}(x_{2})}(w,c_{1})}{p_{W(x_{2}),C_{1}(x_{2})}(w,c_{1})} \\ &= \frac{p_{W(x_{2}),C_{1}(x_{2})}(w,c_{1})}{p_{W(x_{1}),C_{1}(x_{1})}(w,c_{1})} p_{X_{1}(x_{2})|W(x_{2}),C_{1}(x_{2})}(x_{1}|w,c_{1}) \\ &= \frac{p_{W,C_{1}}^{2}(w,c_{1})}{p_{W,C_{1}}^{1}(w,c_{1})} p_{X_{1}|W,C_{1}}^{2}(x_{1}|w,c_{1}) \\ &= \frac{b_{W,C_{1}}^{2}(w,c_{1})}{p_{W,C_{1}}^{1}(w,c_{1})} p_{X_{1}|W,C_{1}}^{2}(x_{1}|w,c_{1}) \end{split}$$

where

- $\stackrel{1a}{=}$  holds by the first condition of the AC-TTI-PO in Def. B.1, stating that  $W(x_1, x_2) = W(x_1)$  and  $C_1(x_1, x_2) = C_1(x_1)$ .
- $\stackrel{1b}{=}$  holds by the Composition axiom in Def. A.1. Specifically,  $X_1(x_2) = x_1$  implies  $W(x_1, x_2) = W(x_2)$  and  $C_1(x_1, x_2) = C_1(x_2)$ .
- $\stackrel{1c}{>}$  holds by the given assumption.
- $\stackrel{2}{=}$  holds since  $W(x_1, x_2) = W(x_1)$  and  $C_1(x_1, x_2) = C_1(x_1)$  by the first condition of the AC-TTI-PO in Def. B.1.

- $\stackrel{3}{=}$  holds by the Composition axiom in Def. A.1. Specifically,  $X_1(x_2) = x_1$  implies  $Y(x_1, x_2) = Y(x_2)$ .
- $\stackrel{4}{=}$  by the definition of  $P^1$ .

We note that  $\mathbb{E}[Y(x_2)|W(x_1), C_1(x_1), X_1(x_2) = x_1]$  is estimable from  $P^2(Y|W, C_1, X_1)$  since

$$P^{2}(Y|W, C_{1}, X_{1} = x_{1}) \stackrel{5}{=} P(Y(x_{2})|W(x_{2}), C_{1}(x_{2}), X_{1}(x_{2}) = x_{1})$$

$$\stackrel{6}{=} P(Y(x_{2})|W(x_{1}, x_{2}), C_{1}(x_{1}, x_{2}), X_{1}(x_{2}) = x_{1})$$

$$\stackrel{7}{=} P(Y(x_{2})|W(x_{1}), C_{1}(x_{1}), X_{1}(x_{2}) = x_{1}),$$

where

- $\stackrel{5}{=}$  holds by the definition of  $P^2$ .
- $\stackrel{6}{=}$  holds since  $W(x_1, x_2) = W(x_2)$  and  $C_1(x_1, x_2) = C_1(x_2)$  when  $X_1(x_2) = x_1$  by Composition axiom in Def. A.1.
- $\stackrel{7}{=}$  holds by the first condition of the AC-TTI-PO in Def. B.1, stating that  $W(x_1, x_2) = W(x_1)$  and  $C_1(x_1, x_2) = C_1(x_1)$ .

Therefore,

$$\mathbb{E}_{P^1}\left[\mathbb{E}\left[Y(x_2)|W, C_1, X_1(x_2) = x_1\right]\right] = \mathbb{E}_{P^1}\left[\mathbb{E}_{P^2}\left[Y|W, C_1, X_1 = x_1\right]\right].$$

#### **B.2.** Combining Two Experiments based on Potential Outcome Framework

We provide an adjustment criterion based on potential outcome frameworks for combining two experiments as follow:

**Definition B.2** (Adjustment criterion for combining two experiments – Potential Outcome (AC2-PO)). A set of variables **A** is said to satisfy the adjustment criterion (AC2) w.r.t discrete treatments **X** and the outcome Y from two sets of distributions  $P_{\text{rand}(Z_1)}(\mathbf{A})$  and  $P_{\text{rand}(Z_2)}(\mathbf{A}, \mathbf{X}, Y)$  if

- 1.  $\mathbf{A}(\mathbf{x}) = \mathbf{A}(z_1)$ ; i.e.,  $Z_1 \subseteq \mathbf{X}$  and  $\mathbf{X} \setminus Z_1$  is causally irrelevant to  $\mathbf{A}$  given  $Z_1$ ;
- 2.  $Z_2 \subseteq \mathbf{X}$  and  $Y(\mathbf{x}) \perp (\mathbf{X} \setminus Z_2)(z_2) | \mathbf{A}(\mathbf{x})$ .

Under Def. B.2, the causal effect is identified as follows:

Theorem B.2 (Identification through AC2-PO). Suppose the condition AC2-PO in Def. B.2 holds. Let

$$P^{1}(\mathbf{A}) \coloneqq P(\mathbf{A}(z_{1}))$$

$$P^{2}(Y, \mathbf{X} \setminus Z_{2}, \mathbf{A}) \coloneqq P(Y(z_{2}), (\mathbf{X} \setminus Z_{2})(z_{2}), \mathbf{A}(z_{2})),$$

and  $p^1, p^2$  are densities for distributions  $P^1, P^2$ . Assume the following positivity condition:

$$\frac{p^2(\mathbf{a})}{p^1(\mathbf{a})}p^2(\mathbf{X}\backslash Z_2 = \mathbf{x}\backslash z_2|\mathbf{a}), \ \forall \mathbf{x}, \mathbf{a} \in \mathcal{X} \times \mathcal{A}.$$
(B.3)

Then, the query  $\mathbb{E}[Y(\mathbf{x})]$  is identifiable from two distributions  $P^1, P^2$  and given as follow:

$$\mathbb{E}\left[Y(\mathbf{x})\right] = \mathbb{E}_{P^1}\left[\mathbb{E}_{P^2}\left[Y|\mathbf{A}, \mathbf{x} \setminus z_2\right]\right].$$
(B.4)

Proof of Theorem **B.2**.

$$\begin{split} \mathbb{E}\left[Y(\mathbf{x})\right] &= \mathbb{E}\left[\mathbb{E}\left[Y(\mathbf{x})|\mathbf{A}(\mathbf{x})\right]\right] \\ &\stackrel{1}{=} \mathbb{E}\left[\mathbb{E}\left[Y(\mathbf{x})|\mathbf{A}(\mathbf{x}), (\mathbf{X}\backslash Z_2)(z_2) = \mathbf{x}\backslash z_2\right]\right] \\ &\stackrel{2}{=} \mathbb{E}\left[\mathbb{E}\left[Y(\mathbf{x})|\mathbf{A}(z_1), (\mathbf{X}\backslash Z_2)(z_2) = \mathbf{x}\backslash z_2\right]\right] \\ &\stackrel{3}{=} \mathbb{E}\left[\mathbb{E}\left[Y(z_2)|\mathbf{A}(z_1), (\mathbf{X}\backslash Z_2)(z_2) = \mathbf{x}\backslash z_2\right]\right] \\ &\stackrel{4}{=} \mathbb{E}_{P_1}\left[\mathbb{E}\left[Y(z_2)|\mathbf{A}, (\mathbf{X}\backslash Z_2)(z_2) = \mathbf{x}\backslash z_2\right]\right], \end{split}$$

where

•  $\stackrel{1}{=}$  holds since  $Y(\mathbf{x}) \perp (\mathbf{X} \setminus Z_2)(z_2) | \mathbf{A}(\mathbf{x})$ , and the positivity  $P((\mathbf{X} \setminus Z_2)(z_2) | \mathbf{A}(\mathbf{x})) > 0$  for any  $\mathbf{x}$ . To witness, it suffices to show that  $p_{(\mathbf{X} \setminus Z_2)(z_2) | \mathbf{A}(\mathbf{x})}(\mathbf{x} \setminus z_2 | \mathbf{a}) > 0$ .

$$p_{(\mathbf{X}\backslash Z_{2})(z_{2})|\mathbf{A}(\mathbf{x})}(\mathbf{x}\backslash z_{2}|\mathbf{a}) = \frac{p_{\mathbf{X}\backslash Z_{2},\mathbf{A}(\mathbf{x})}(\mathbf{x}\backslash z_{2},\mathbf{a})}{p_{\mathbf{A}(\mathbf{x})}(\mathbf{a})}$$

$$\stackrel{1a}{=} \frac{p_{(\mathbf{X}\backslash Z_{2})(z_{2}),\mathbf{A}(\mathbf{x})}(\mathbf{x}\backslash z_{2},\mathbf{a})}{p_{\mathbf{A}(z_{1})}(\mathbf{a})}$$

$$\stackrel{1b}{=} \frac{p_{(\mathbf{X}\backslash Z_{2})(z_{2}),\mathbf{A}(z_{2})}(\mathbf{x}\backslash z_{2},\mathbf{a})}{p_{\mathbf{A}(z_{1})}(\mathbf{a})}$$

$$= \frac{p_{(\mathbf{X}\backslash Z_{2})(z_{2}),\mathbf{A}(z_{2})}(\mathbf{x}\backslash z_{2},\mathbf{a})}{p_{\mathbf{A}(z_{1})}(\mathbf{a})} \frac{p_{\mathbf{A}(z_{2})}(\mathbf{a})}{p_{\mathbf{A}(z_{2})}(\mathbf{a})}$$

$$= \frac{p_{\mathbf{A}(z_{2})}(\mathbf{a})}{p_{\mathbf{A}(z_{1})}(\mathbf{a})} p_{(\mathbf{X}\backslash Z_{2})(z_{2})|\mathbf{A}(z_{2})}(\mathbf{x}\backslash z_{2}|\mathbf{a})$$

$$= \frac{p^{2}(\mathbf{a})}{p^{1}(\mathbf{a})} p^{2}(\mathbf{x}\backslash z_{2}|\mathbf{a})$$

$$\stackrel{1c}{>} 0,$$

where

-  $\stackrel{1a}{=}$  holds by the first condition of the AC2-PO in Def. B.2, stating that  $\mathbf{A}(\mathbf{x}) = \mathbf{A}(z_1)$ .

-  $\stackrel{1b}{=}$  holds by the Composition axiom in Def. A.1. Specifically,  $(\mathbf{X} \setminus Z_2)(z_2) = \mathbf{x} \setminus z_2$  implies  $\mathbf{A}(\mathbf{x}) = \mathbf{A}(z_2)$ .

- $>^{1c}$  holds by the given assumption.
- $\stackrel{2}{=}$  holds since  $\mathbf{A}(\mathbf{x}) = \mathbf{A}(z_2)$ .
- $\stackrel{3}{=}$  holds by the Composition axiom in Def. A.1. Specifically,  $(\mathbf{X} \setminus Z_2)(z_2) = \mathbf{x} \setminus \mathbf{z}_2$  implies  $Y(\mathbf{x}) = Y(z_2)$ .
- $\stackrel{4}{=}$  holds by the definition of  $P^1$ .

We note that  $\mathbb{E}[Y(z_2)|\mathbf{A}(z_1), (\mathbf{X}\setminus Z_2)(z_2) = \mathbf{x}\setminus z_2]$  is estimable from  $P_2(Y|\mathbf{A}, \mathbf{X}\setminus Z_2)$  since

$$P_{2}(Y|\mathbf{A}, \mathbf{X} \setminus Z_{2} = \mathbf{x} \setminus z_{2}) \stackrel{\scriptscriptstyle{5}}{=} P(Y(z_{2})|\mathbf{A}(z_{2}), (\mathbf{X} \setminus Z_{2})(z_{2}) = \mathbf{x} \setminus z_{2})$$
$$\stackrel{\scriptscriptstyle{6}}{=} P(Y(z_{2})|\mathbf{A}(\mathbf{x}), (\mathbf{X} \setminus Z_{2})(z_{2}) = \mathbf{x} \setminus z_{2})$$
$$\stackrel{\scriptscriptstyle{7}}{=} P(Y(z_{2})|\mathbf{A}(z_{1}), (\mathbf{X} \setminus Z_{2})(z_{2}) = \mathbf{x} \setminus z_{2}),$$

where

•  $\stackrel{5}{=}$  holds by the definition.

- $\stackrel{6}{=}$  holds since  $\mathbf{A}(\mathbf{x})$  when  $(\mathbf{X} \setminus Z_2)(z_2) = \mathbf{x} \setminus z_2$  by Composition axiom in Def. A.1.
- $\stackrel{7}{=}$  holds since  $\mathbf{A}(z_1) = \mathbf{A}(\mathbf{x})$ .

Therefore,

$$\mathbb{E}_{P^1}\left[\mathbb{E}\left[Y(x_2)|\mathbf{A}, (\mathbf{X}\backslash Z_2)(z_2)=\mathbf{x}\backslash z_2\right]\right]=\mathbb{E}_{P^1}\left[\mathbb{E}_{P^2}\left[Y|\mathbf{A}, \mathbf{x}\backslash z_2\right]\right].$$

## **B.3. Multiple Treatment Interaction based on Potential Outcome Framework**

We first provide a sufficient identification criterion based on potential outcome frameworks for estimating multiple treatment interaction from multiple marginal experiments.

**Definition B.3 (Adjustment Criterion for MTI – Potential Outcome (AC-MTI-PO)).** A set of variables  $\{\mathbf{C}, \mathbf{W}\}$  is said to satisfy the adjustment criterion for multiple treatment interaction (AC-MTI-PO) w.r.t. discrete treatments  $\mathbf{x} = \{x_i\}_{i=1}^m$  and the outcome Y from multiple distributions  $\{P_{x_i}(\mathbf{V})\}_{i=1}^m$  if

- 1.  $W_i(\mathbf{x}) = W_i(\mathbf{x}^{(i)}), C_i(\mathbf{x}) = C_i(\mathbf{x}^{(i)})$  for  $i = 1, 2, \dots, m-1$ ; i.e., the *i*th joint outcome  $W_i(\mathbf{x})$  and the covariate  $C_i(\mathbf{x})$  are invariant to the next interventions  $X_{i+1}, \dots, X_m$ .
- 2. For all  $i = 1, 2, \dots, m-1$ ,  $X_i = X_i(\mathbf{x}^{(i+1:k)}), \forall k \in \{i+1, \dots, m\}$ ; i.e.,  $X_i$  is invariant to any intervention  $X_k = x_k$  for k > i.
- 3.  $Y(\mathbf{x}) \perp X_i | \mathbf{C}^{(i)}(\mathbf{x}), \mathbf{X}^{(i-1)}, \mathbf{W}^{(i)}(\mathbf{x})$  for  $i = 1, 2, \dots, m-1$ ; i.e., the *i*th intervention  $X_i = x_i$  is non-informative to the joint outcome  $Y(\mathbf{x})$  given the *i*th outcome  $W_i(\mathbf{x})$ , covariate  $C_i(\mathbf{x})$  and previous observations  $\mathbf{X}^{(i-1)}, \mathbf{W}^{(i-1)}, \mathbf{C}^{(i-1)}$ .

Under Def. B.3, the causal effect  $\mathbb{E}\left[Y(\mathbf{x})|\mathbf{A}^{(i-1)}(\mathbf{x}), \mathbf{X}^{(i-1)}\right]$  for  $i = 1, 2, \dots, m-1$  can be expressed in a recursive form as follow:

**Lemma B.1.** Suppose the condition AC-MTI-PO in Def. B.3 holds. Let  $A_i := \{W_i, C_i\}$ . Assume the following positivity condition holds: For  $i = 1, 2, \dots, m-1$ 

$$p_{X_i|\mathbf{A}^{(i)}(\mathbf{x}),\mathbf{X}^{(i-1)}}(x_i|\mathbf{a}^{(i)},\mathbf{x}^{(i-1)}) > 0, \ \forall \mathbf{a}, \mathbf{x} \in \mathfrak{D}_{\mathbf{A},\mathbf{X}}.$$
(B.5)

Then,

$$\mathbb{E}\left[Y(\mathbf{x})\middle|\mathbf{A}^{(i-1)}(\mathbf{x}),\mathbf{X}^{(i-1)}\right] = \mathbb{E}\left[\mathbb{E}\left[Y(\mathbf{x})|\mathbf{A}^{(i)}(\mathbf{x}),X_i=x_i,\mathbf{X}^{(i-1)}\right]\middle|\mathbf{A}^{(i-1)}(\mathbf{x}),\mathbf{X}^{(i-1)}\right]\right]$$

where  $\mathbf{A}_0 \coloneqq \emptyset$  and  $X_0 \coloneqq \emptyset$ .

Proof of Lemma B.1.

$$\mathbb{E}\left[Y(\mathbf{x})\middle|\mathbf{A}^{(i-1)}(\mathbf{x}), \mathbf{X}^{(i-1)}\right] \stackrel{1}{=} \mathbb{E}\left[\mathbb{E}\left[Y(\mathbf{x})\middle|\mathbf{A}^{(i)}(\mathbf{x}), \mathbf{X}^{(i-1)}\right]\middle|\mathbf{A}^{(i-1)}(\mathbf{x}), \mathbf{X}^{(i-1)}\right] \\ \stackrel{2}{=} \mathbb{E}\left[\mathbb{E}\left[Y(\mathbf{x})\middle|\mathbf{A}^{(i)}(\mathbf{x}), X_{i} = x_{i}, \mathbf{X}^{(i-1)}\right]\middle|\mathbf{A}^{(i-1)}(\mathbf{x}), \mathbf{X}^{(i-1)}\right],$$

where

•  $\stackrel{\perp}{=}$  holds by marginalizing over  $\mathbf{A}_i(\mathbf{x})$ .

•  $\stackrel{2}{=}$  holds by the condition in the AC-MTI-PO; i.e.,  $Y(\mathbf{x}) \perp X_i | \mathbf{A}^{(i)}(\mathbf{x}), \mathbf{X}^{(i-1)}$  for  $i = 1, 2, \dots, m-1$ , and the given positivity condition in Eq. (B.5).

**Corollary B.1** (Corollary of Lemma B.1). Suppose the condition AC-MTI-PO in Def. B.3 holds. Let  $A_i := \{W_i, C_i\}$ . Assume the positivity condition given in Eq. (B.5). Let

$$\nu_0^i(\mathbf{A}^{(i-1)}(\mathbf{x}), \mathbf{X}^{(i-1)}) \coloneqq \mathbb{E}\left[Y(\mathbf{x}) | \mathbf{A}^{(i-1)}(\mathbf{x}), \mathbf{X}^{(i-1)}\right], \text{ for } i = m, m-1, \cdots, 2.$$

Then,

$$\nu_0^i(\mathbf{A}^{(i-1)}(\mathbf{x}), \mathbf{X}^{(i-1)}) = \mathbb{E}\left[\nu_0^{i+1}(\mathbf{A}^{(i)}(\mathbf{x}), x_i, \mathbf{X}^{(i-1)}) \middle| \mathbf{A}^{(i-1)}(\mathbf{x}), \mathbf{X}^{(i-1)}\right] \text{ for } i = m - 1, \cdots, 2,$$

where  $\mathbf{A}_0 \coloneqq \emptyset$  and  $X_0 \coloneqq \emptyset$ . Furthermore,

$$\mathbb{E}[Y(\mathbf{x})] = \mathbb{E}\left[\nu_0^2(\mathbf{A}^{(1)}(\mathbf{x}), x_1)\right].$$

Proof of Corollary B.1. We first note that the equations

$$\nu_0^i(\mathbf{A}^{(i-1)}(\mathbf{x}), \mathbf{X}^{(i-1)}) \coloneqq \mathbb{E}\left[Y(\mathbf{x}) | \mathbf{A}^{(i-1)}(\mathbf{x}), \mathbf{X}^{(i-1)}\right], \text{ for } i = m, m-1, \cdots, 2,$$

is immediately followed by Lemma B.1. Therefore, it only suffices to show the following:

$$\mathbb{E}[Y(\mathbf{x})] = \mathbb{E}\left[\nu_0^2(\mathbf{A}^{(1)}(\mathbf{x}), x_1)\right].$$

To witness,

$$\mathbb{E}\left[Y(\mathbf{x})\right] = \mathbb{E}\left[\mathbb{E}\left[Y(\mathbf{x})|\mathbf{A}^{(1)}(\mathbf{x})\right]\right]$$
$$= \mathbb{E}\left[\mathbb{E}\left[Y(\mathbf{x})|\mathbf{A}^{(1)}(\mathbf{x}), X_1 = x_1\right]\right]$$
$$= \mathbb{E}\left[\nu_0^2(\mathbf{A}^{(1)}(\mathbf{x}), x_1)\right],$$

where the second equation holds by the condition in the AC-MTI-PO; i.e.,  $Y(\mathbf{x}) \perp X_i | \mathbf{A}^{(i)}(\mathbf{x}), \mathbf{X}^{(i-1)}$  for  $i = 1, 2, \dots, m-1$ , and the given positivity condition in Eq. (B.5).

**Lemma B.2.** Suppose the condition AC-MTI-PO in Def. B.3 holds. Let  $A_i := \{W_i, C_i\}$ . Assume the positivity condition given in Eq. (B.5). For  $i = m, \dots, 1$ , and

$$\nu_0^i(\mathbf{A}^{(i-1)}(\mathbf{x}), \mathbf{X}^{(i-1)}) \coloneqq \mathbb{E}\left[Y(\mathbf{x}) | \mathbf{A}^{(i-1)}(\mathbf{x}), \mathbf{X}^{(i-1)}\right].$$

For  $i = 1, 2, \dots, m$ , let  $P^i$  denote a distribution defined as follow:

$$P^{i}(\mathbf{W}^{(i)}, \mathbf{C}^{(i)}, \mathbf{X}^{(i)}) \coloneqq P(\mathbf{W}^{(i)}(x_{i}), \mathbf{C}^{(i)}(x_{i}), \mathbf{X}^{(i)}(x_{i})),$$

and  $p^1, \dots, p^{m-1}$  are densities for distributions  $P^1, \dots, P^m$ . Let

$$\mu_0^m(\mathbf{A}^{(m-1)}, \mathbf{X}^{(m-1)}) \coloneqq \mathbb{E}_{P^m}\left[Y | \mathbf{A}^{(m-1)}, \mathbf{X}^{(m-1)}\right],$$

and for  $i = m - 1, \dots, 1$ ,

$$\mu_0^i(\mathbf{A}^{(i-1)}, \mathbf{X}^{(i-1)}) \coloneqq \mathbb{E}_{P^i} \left[ \mu_0^{i+1}(\mathbf{A}^{(i)}, x_i, \mathbf{X}^{(i-1)}) \middle| \mathbf{A}^{(i-1)}, \mathbf{X}^{(i-1)} \right].$$

Then, for  $i = m, \cdots, 1$ ,

$$\mu_0^i(\mathbf{A}^{(i-1)}(\mathbf{x}), \mathbf{x}^{(i-1)}) = \nu_0^i(\mathbf{A}^{(i-1)}(\mathbf{x}), \mathbf{x}^{(i-1)}).$$

## **Proof of Lemma B.2.** We first show that

$$\mu_0^m(\mathbf{A}^{(m-1)}, \mathbf{x}^{(m-1)}) = \nu_0^m(\mathbf{A}^{(m-1)}(\mathbf{x}), \mathbf{x}^{(m-1)}).$$

To witness,

$$\nu_{0}^{m}(\mathbf{A}^{(m-1)}(\mathbf{x}), \mathbf{x}^{(m-1)}) \coloneqq \mathbb{E}\left[Y(\mathbf{x}) | \mathbf{A}^{(m-1)}(\mathbf{x}), \mathbf{X}^{(m-1)} = \mathbf{x}^{(m-1)}\right]$$
  
$$\stackrel{1}{=} \mathbb{E}\left[Y(\mathbf{x}) | \mathbf{A}^{(m-1)}(\mathbf{x}), \mathbf{X}^{(m-1)}(x_{m}) = \mathbf{x}^{(m-1)}\right]$$
  
$$\stackrel{2}{=} \mathbb{E}\left[Y(x_{m}) | \mathbf{A}^{(m-1)}(\mathbf{x}), \mathbf{X}^{(m-1)}(x_{m}) = \mathbf{x}^{(m-1)}\right]$$
  
$$\stackrel{3}{=} \mathbb{E}\left[Y(x_{m}) | \mathbf{A}^{(m-1)}(x_{m}), \mathbf{X}^{(m-1)}(x_{m}) = \mathbf{x}^{(m-1)}\right]$$
  
$$= \mathbb{E}_{P^{m}}\left[Y | \mathbf{A}^{(m-1)}, \mathbf{X}^{(m-1)} = \mathbf{x}^{(m-1)}\right]$$
  
$$=: \mu_{0}^{m}(\mathbf{A}^{(m-1)}, \mathbf{x}^{(m-1)}),$$

where

- $\stackrel{1}{=}$  holds since  $\mathbf{X}^{(m-1)}(x_m) = \mathbf{X}^{(m-1)}$  by the condition of the AC-MTI-PO in Def. B.3 stating that treatment variables are invariant to the next interventions.
- $\stackrel{2}{=}$  holds since

$$\mathbf{X}^{(m-1)}(x_m) = \mathbf{x}^{(m-1)} \implies Y(\mathbf{x}) = Y(\mathbf{x}^{(m-1)}, x_m) = Y(x_m),$$

by Composition axiom in Axiom A.1.

•  $\stackrel{3}{=}$  holds since

$$\mathbf{X}^{(m-1)}(x_m) = \mathbf{x}^{(m-1)} \implies \mathbf{A}^{(m-1)}(\mathbf{x}) = \mathbf{A}^{(m-1)}(\mathbf{x}^{(m-1)}, x_m) = \mathbf{A}^{(m-1)}(x_m),$$

by Composition axiom in Axiom A.1.

We now make an induction hypothesis as follow: For any given  $i \in \{2, \dots, m-1\}$  suppose the following holds:

$$\mu_0^{i+1}(\mathbf{A}^{(i)}, \mathbf{x}^{(i)}) = \nu_0^{i+1}(\mathbf{A}^{(i)}(\mathbf{x}), \mathbf{x}^{(i)}).$$

Then,

$$\begin{split} \nu_0^i(\mathbf{A}^{(i-1)}(\mathbf{x}), \mathbf{x}^{(i-1)}) &\stackrel{4}{=} \mathbb{E} \left[ \nu_0^{i+1}(\mathbf{A}^{(i)}(\mathbf{x}), x_i, \mathbf{X}^{(i-1)}) \big| \mathbf{A}^{(i-1)}(\mathbf{x}), \mathbf{X}^{(i-1)} = \mathbf{x}^{(i-1)} \right] \\ &= \mathbb{E} \left[ \nu_0^{i+1}(\mathbf{A}^{(i)}(\mathbf{x}), \mathbf{x}^{(i)}) \big| \mathbf{A}^{(i-1)}(\mathbf{x}), \mathbf{x}^{(i-1)} \right] \\ &\stackrel{5}{=} \mathbb{E} \left[ \mu_0^{i+1}(\mathbf{A}^{(i)}, \mathbf{x}^{(i)}) \big| \mathbf{A}^{(i-1)}(\mathbf{x}), \mathbf{x}^{(i-1)} \right] \\ &\stackrel{6}{=} \mathbb{E} \left[ \mu_0^{i+1}(\mathbf{A}^{(i)}, \mathbf{x}^{(i)}) \big| \mathbf{A}^{(i-1)}(x_i), \mathbf{x}^{(i-1)}(x_i) \right] \\ &= \mathbb{E}_{P^i} \left[ \mu_0^{i+1}(\mathbf{A}^{(i)}, \mathbf{x}^{(i)}) \big| \mathbf{A}^{(i-1)}, \mathbf{x}^{(i-1)} \right] \\ &= \mu_0^i(\mathbf{A}^{(i-1)}, \mathbf{x}^{(i-1)}), \end{split}$$

where

- $\stackrel{4}{=}$  implies from Corollary B.1.
- $\stackrel{5}{=}$  because of the induction hypothesis.
- $\stackrel{6}{=}$  holds because  $\mathbf{X}^{(i-1)} = \mathbf{X}^{(i-1)}(x_i)$  by the Composition Axiom in Def. A.1, and

$$\mathbf{X}^{(i-1)}(x_i) = \mathbf{x}^{(i-1)} \implies \mathbf{A}^{(i-1)}(\mathbf{x}) = \mathbf{A}^{(i-1)}(\mathbf{x}^{(i-1)}) = \mathbf{A}^{(i-1)}(\mathbf{X}^{(i-1)}(x_i) = \mathbf{x}^{(i-1)}, x_i) = \mathbf{A}^{(i-1)}(x_i)$$

by applying the Composition Axiom in Def. A.1.

This proves that, for  $i = m, m - 1, \cdots, 1$ ,

$$\mu_0^i(\mathbf{A}^{(i-1)}, \mathbf{x}^{(i-1)}) = \nu_0^i(\mathbf{A}^{(i-1)}(\mathbf{x}), \mathbf{x}^{(i-1)}).$$

**Theorem B.3** (Identification through AC-MTI-PO). Suppose the condition AC-MTI-PO in Def. B.3 holds. For  $i = 1, 2, \dots, m$ , let  $P^i$  denote a distribution defined as follow:

$$P^{i}(\mathbf{W}^{(i)}, \mathbf{C}^{(i)}, \mathbf{X}^{(i)}) \coloneqq P(\mathbf{W}^{(i)}(x_{i}), \mathbf{C}^{(i)}(x_{i})), \mathbf{X}^{(i)}(x_{i})),$$

and  $p^1, \dots, p^{m-1}$  are densities for distributions  $P^1, \dots, P^m$ . Assume the following positivity condition holds: For all  $i = 1, 2, \dots, m-1$ 

$$\frac{p^{i+1}(w_i, c_i | \mathbf{w}^{(i-1)}, \mathbf{c}^{(i-1)}, \mathbf{x}^{(i-1)})}{p^i(w_i, c_i | \mathbf{w}^{(i-1)}, \mathbf{c}^{(i-1)}, \mathbf{x}^{(i-1)})} p^{i+1}(X_i = x_i | \mathbf{w}^{(i)}, \mathbf{c}^{(i)}, \mathbf{x}^{(i-1)}) > 0; \ \forall \mathbf{w}, \mathbf{c}, \mathbf{x} \in \mathfrak{D}_{\mathbf{W}, \mathbf{C}, \mathbf{X}}.$$
(B.6)

Then, the query  $\mathbb{E}[Y(\mathbf{x})]$  is identifiable from distributions  $P^1, \dots, P^m$ , and given as follow: Let

$$\mu_0^m(\mathbf{A}^{(m-1)}, \mathbf{X}^{(m-1)}) \coloneqq \mathbb{E}_{P^m}\left[Y | \mathbf{A}^{(m-1)}, \mathbf{X}^{(m-1)}\right],$$

and for  $i = m - 1, m - 2, \cdots, 2$ ,

$$\mu_0^i(\mathbf{A}^{(i-1)}, \mathbf{X}^{(i-1)}) \coloneqq \mathbb{E}_{P^i} \left[ \mu_0^{i+1}(\mathbf{A}^{(i)}, x_i, \mathbf{X}^{(i-1)}) | \mathbf{A}^{(i-1)}, \mathbf{X}^{(i-1)} \right].$$

Then,

$$\mathbb{E}[Y(\mathbf{x})] = \mathbb{E}_{P^1}\left[\mu_0^2(\mathbf{A}^{(1)}, \mathbf{x}^{(1)})\right].$$

*Proof of Theorem B.3*. Suppose the positivity condition in Eq. (B.5) is equivalent to the condition in Eq. (B.6). Then, Theorem B.3 is implied by Lemmas (B.1, B.2) and Corollary B.1.

The equivalence between Eq. (B.5) and Eq. (B.6) are the following. We will use  $A_i := \{W_i, C_i\}$ . Then,

$$\begin{split} p_{X_{i}|\mathbf{A}^{(i)}(\mathbf{x}),\mathbf{X}^{(i-1)}(x_{i}|\mathbf{a}^{(i)},\mathbf{x}^{(i-1)})} \\ &= \frac{p_{X_{i},A_{i}(\mathbf{x})|\mathbf{A}^{(i-1)}(\mathbf{x}),\mathbf{X}^{(i-1)}(x_{i},a_{i}|\mathbf{a}^{(i-1)},\mathbf{x}^{(i-1)})}{p_{A_{i}(\mathbf{x})|\mathbf{A}^{(i-1)}(\mathbf{x}),\mathbf{X}^{(i-1)}(a_{i}|\mathbf{a}^{(i-1)},\mathbf{x}^{(i-1)})} \\ &= \frac{p_{X_{i}(x_{i+1}),A_{i}(x_{i+1})|\mathbf{A}^{(i-1)}(x_{i+1}),\mathbf{X}^{(i-1)}(x_{i+1})(x_{i},a_{i}|\mathbf{a}^{(i-1)},\mathbf{x}^{(i-1)})}{p_{A_{i}(x_{i})|\mathbf{A}^{(i-1)}(x_{i}),\mathbf{X}^{(i-1)}(x_{i}),\mathbf{x}^{(i-1)}(x_{i+1})(x_{i},a_{i}|\mathbf{a}^{(i-1)},\mathbf{x}^{(i-1)})} \\ &= \frac{p_{X_{i}(x_{i+1}),A_{i}(x_{i+1})|\mathbf{A}^{(i-1)}(x_{i+1}),\mathbf{X}^{(i-1)}(x_{i+1})(x_{i},a_{i}|\mathbf{a}^{(i-1)},\mathbf{x}^{(i-1)})}{p_{A_{i}(x_{i})|\mathbf{A}^{(i-1)}(x_{i}),\mathbf{X}^{(i-1)}(x_{i})(a_{i}|\mathbf{a}^{(i-1)},\mathbf{x}^{(i-1)})}{p_{A_{i}(x_{i})|\mathbf{A}^{(i-1)}(x_{i+1}),\mathbf{X}^{(i-1)}(x_{i+1})(a_{i}|\mathbf{a}^{(i-1)},\mathbf{x}^{(i-1)})}} \\ &= \frac{p_{A_{i}(x_{i+1})|\mathbf{A}^{(i-1)}(x_{i+1}),\mathbf{X}^{(i-1)}(x_{i+1})(a_{i}|\mathbf{a}^{(i-1)},\mathbf{x}^{(i-1)})}{p_{A_{i}(x_{i})|\mathbf{A}^{(i-1)}(x_{i}),\mathbf{X}^{(i-1)}(x_{i})}p_{A_{i}(x_{i+1})|\mathbf{A}^{(i-1)}(x_{i+1}),\mathbf{X}^{(i-1)}(x_{i+1})(a_{i}|\mathbf{a}^{(i-1)},\mathbf{x}^{(i-1)})}} \\ &= \frac{p_{A_{i}(x_{i+1})|\mathbf{A}^{(i-1)}(x_{i}),\mathbf{X}^{(i-1)}(x_{i})(a_{i}|\mathbf{a}^{(i-1)},\mathbf{x}^{(i-1)})}{p_{A_{i}(x_{i})|\mathbf{A}^{(i-1)}(x_{i}),\mathbf{X}^{(i-1)}(x_{i})}p_{A_{i}(x_{i+1})|\mathbf{A}^{(i-1)}(x_{i+1})(x_{i+1})(a_{i}|\mathbf{a}^{(i-1)},\mathbf{x}^{(i-1)})}}{p_{A_{i}(x_{i})|\mathbf{A}^{(i-1)}(x_{i}),\mathbf{X}^{(i-1)}(x_{i})}p_{A_{i}(x_{i+1})|\mathbf{A}^{(i)}(x_{i+1}),\mathbf{X}^{(i-1)}(x_{i+1})(x_{i}|\mathbf{a}^{(i)},\mathbf{x}^{(i-1)})}} \\ &= \frac{p^{i+1}(w_{i},c_{i}|\mathbf{w}^{(i-1)},\mathbf{c}^{(i-1)},\mathbf{x}^{(i-1)})}{p^{i}(w_{i},c_{i}|\mathbf{tw}^{(i-1)},\mathbf{c}^{(i-1)},\mathbf{x}^{(i-1)})}p^{i+1}(X_{i}=x_{i}|\mathbf{w}^{(i)},\mathbf{c}^{(i)},\mathbf{x}^{(i-1)}). \end{split}$$

## **B.4.** Combining Multiple Experiments based on Potential Outcome Framework

We provide an adjustment criterion based on potential outcome frameworks for combining two experiment as follow:

**Definition B.4 (Adjustment criterion for combining multiple experiments – Potential Outcome (AC-gMTI-PO)).** Let  $\mathbf{X} \coloneqq \{X_1, \dots, X_{m_x}\}$  and Y denote an ordered treatments and outcome variables. Let  $\mathbf{Z} \coloneqq \{Z_1, \dots, Z_m\} \subseteq \mathbf{X}$  denote the subset of treatments. Let  $\{\ell_i\}_{i=1}^m \subseteq \{1, 2, \dots, m_x\}$  denote the index of  $\mathbf{Z}$ ; i.e.,  $\mathbf{Z} = \{X_{\ell_1}, \dots, X_{\ell_m}\}$ . Let  $\overline{X}_1 \coloneqq \{X_j\}_{j \leq \ell_1}, \overline{X}_{m+1} \coloneqq \{X_j\}_{j > \ell_m}$ , and  $\overline{X}_i \coloneqq \{X_j\}_{\ell_{i-1} < j \leq \ell_i}$  for  $i = 2, 3, \dots, m$ . A set of topologically ordered variables  $\mathbf{A} \coloneqq \{\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_m\}$  is said to satisfy the adjustment criterion for combining multiple experiments (AC-gMTI-PO) w.r.t. treatments  $\mathbf{x} = \{x_i\}_{i=1}^{m_x}$  and the outcome Y from multiple distributions  $\{P_{z_i}(\mathbf{A}^{(i)}, \overline{\mathbf{X}}^{(i)})\}_{i=1}^m$  if

1.  $\mathbf{A}_i(\mathbf{x}) = \mathbf{A}_i(\mathbf{z}^{(i)})$  for  $i = 1, 2, \dots, m-1$  and  $\overline{X}_i(z_j) = \overline{X}_i$  for all  $i.j \in \{1, 2, \dots, m\}$  where  $i \leq j$ .

2. 
$$Y(\mathbf{x}) \perp \perp \overline{X}_i | \mathbf{A}^{(i)}(\mathbf{x}), \overline{\mathbf{X}}^{(i-1)} \text{ for } i = 1, 2, \cdots, m-1$$

3. 
$$Y(\mathbf{x}) \perp \{\overline{X}_m \setminus Z_m, \overline{X}_{m+1}(z_m)\} | \mathbf{A}^{(m-1)}(\mathbf{x}), \overline{\mathbf{X}}^{(m-1)}$$

We first note that the causal effect  $\mathbb{E}[Y(\mathbf{x})]$  can be represented in a recursive form as follow:

**Lemma B.3.** Suppose the condition AC-gMTI-PO in Def. B.4 holds. Assume the following positivity condition holds: For  $\forall \overline{x}_m, \overline{x}_{m+1}, \mathbf{a}^{(m-1)}, \overline{\mathbf{x}}^{(m-1)} \in \mathfrak{D}_{\overline{X}_m, \overline{X}_{m+1}, \mathbf{A}^{(m-1)}, \overline{\mathbf{x}}^{(m-1)}}$ ,

$$p_{\overline{X}_m \setminus Z_m, \overline{X}_{m+1}(z_m) | \mathbf{A}^{(m-1)}(\mathbf{x}), \overline{\mathbf{X}}^{(m-1)}(\overline{x}_m \setminus z_m, \overline{x}_{m+1} | \mathbf{a}^{(m-1)}, \overline{\mathbf{x}}^{(m-1)}) > 0,$$
(B.7)

and for  $i = 1, 2, \dots, m - 1$ 

$$p_{\overline{X}_{i}|\mathbf{A}^{(i)}(\mathbf{x}),\overline{\mathbf{X}}^{(i-1)}}(\overline{x}_{i}|\mathbf{a}^{(i)},\overline{\mathbf{x}}^{(i-1)}) > 0, \ \forall \mathbf{a}^{(i)},\mathbf{x}^{(i)} \in \mathcal{A}^{(i)} \times \mathcal{X}^{(i)}.$$
(B.8)

Let

$$\nu_0^m(\mathbf{A}^{(m-1)}(\mathbf{x}), \mathbf{x} \backslash z_m) \coloneqq \mathbb{E}\left[Y(\mathbf{x}) | \mathbf{A}^{(m-1)}(\mathbf{x}), (\mathbf{X} \backslash Z_m)(z_m) = \mathbf{x} \backslash z_m\right]$$
$$\nu_0^{m-1}(\mathbf{A}^{(m-2)}(\mathbf{x}), \overline{\mathbf{x}}^{(m-2)}) \coloneqq \mathbb{E}\left[\nu_0^m(\mathbf{A}^{(m-1)}(\mathbf{x}), \mathbf{x} \backslash z_m) \middle| \mathbf{A}^{(m-2)}(\mathbf{x}), \overline{\mathbf{x}}^{(m-2)}\right],$$

and for  $i = m - 2, \dots, 2$ ,

$$\nu_0^i(\mathbf{A}^{(i-1)}(\mathbf{x}), \overline{\mathbf{x}}^{(i-1)}) \coloneqq \mathbb{E}\left[\nu_0^{i+1}(\mathbf{A}^{(i)}(\mathbf{x}), \overline{\mathbf{x}}^{(i)}) \middle| \mathbf{A}^{(i-1)}(\mathbf{x}), \overline{\mathbf{x}}^{(i-1)}\right].$$

Then,

$$\mathbb{E}\left[Y(\mathbf{x})\right] = \mathbb{E}\left[\nu^2(A_1(\mathbf{x}), \overline{x}_1)\right].$$

Proof of Lemma B.3. Let

$$\eta_0^{i+1}(\mathbf{A}^{(i)}(\mathbf{x}), \overline{\mathbf{x}}^{(i)}) \coloneqq \mathbb{E}\left[Y(\mathbf{x}) | \mathbf{A}^{(i)}(\mathbf{x}), \overline{\mathbf{x}}^{(i)}\right], \text{ for } i = 1, 2, \cdots, m-1.$$

Then, the causal effect can be written as

$$\mathbb{E}\left[Y(\mathbf{x})\right] = \mathbb{E}\left[\eta_0^2(A_1(\mathbf{x}), \overline{x}_1)\right],$$

since

$$\mathbb{E}\left[Y(\mathbf{x})\right] = \mathbb{E}\left[\mathbb{E}\left[Y(\mathbf{x})|A_1(\mathbf{x})\right]\right] = \mathbb{E}\left[\mathbb{E}\left[Y(\mathbf{x})|\overline{x}_1, A_1(\mathbf{x})\right]\right] = \mathbb{E}\left[\eta^2(A_1(\mathbf{x}), \overline{x}_1)\right],$$

where the second equation holds since  $Y(\mathbf{x}) \perp \overline{X}_i | \mathbf{A}^{(i)}(\mathbf{x}), \overline{\mathbf{X}}^{(i-1)}$  for  $i = 1, 2, \dots, m-1$  by Def. B.4 and the positivity condition in Eq. (B.8).

We will show the following:

$$\eta_0^m(\mathbf{A}^{(m-1)}(\mathbf{x}), \overline{\mathbf{x}}^{(m-1)}) = \nu_0^m(\mathbf{A}^{(m-1)}(\mathbf{x}), \mathbf{x} \setminus z_m)$$
  
$$\eta_0^{i+1}(\mathbf{A}^{(i)}(\mathbf{x}), \overline{\mathbf{x}}^{(i)}) = \nu_0^{i+1}(\mathbf{A}^{(i)}(\mathbf{x}), \overline{\mathbf{x}}^{(i)}) \text{ for } i = m - 2, \cdots, 1$$

First equation can be witnessed by

$$\eta_0^m(\mathbf{A}^{(m-1)}(\mathbf{x}), \overline{\mathbf{x}}^{(m-1)}) = \mathbb{E}\left[Y(\mathbf{x})|\mathbf{A}^{(m-1)}(\mathbf{x}), \overline{\mathbf{x}}^{(m-1)}\right]$$

$$\stackrel{2}{=} \mathbb{E}\left[Y(\mathbf{x})|\mathbf{A}^{(m-1)}(\mathbf{x}), \overline{\mathbf{x}}^{(m-1)}, \{\overline{X}_m \setminus Z_m = \overline{x}_m \setminus z_m, \overline{X}_{m+1}(z_m) = \overline{x}_{m+1}\}\right]$$

$$\stackrel{3}{=} \mathbb{E}\left[Y(\mathbf{x})|\mathbf{A}^{(m-1)}(\mathbf{x}), (\mathbf{X} \setminus Z_m)(z_m) = \mathbf{x} \setminus z_m\right],$$

$$= \nu_0^m(\mathbf{A}^{(m-1)}(\mathbf{x}), \mathbf{x} \setminus z_m),$$

where

• 
$$\stackrel{2}{=}$$
 holds by  $Y(\mathbf{x}) \perp \{\overline{X}_m \setminus Z_m, \overline{X}_{m+1}(z_m)\} | \mathbf{A}^{(m-1)}(\mathbf{x}), \overline{\mathbf{X}}^{(m-1)}$  in Def. B.4 and the positivity condition in Eq. (B.7).

•  $\stackrel{3}{=}$  holds since  $\overline{X}_i(z_j) = \overline{X}_i$  for all i, j, as given in Def. B.4.

The second equation can be witnessed as follow:

$$\begin{split} \eta_0^{m-1}(\mathbf{A}^{(m-2)}(\mathbf{x}), \overline{\mathbf{x}}^{(m-2)}) &\coloneqq \mathbb{E}\left[Y(\mathbf{x}) \middle| \mathbf{A}^{(m-2)}(\mathbf{x}), \overline{\mathbf{x}}^{(m-2)}\right] \\ &= \mathbb{E}\left[\mathbb{E}\left[Y(\mathbf{x}) \middle| \mathbf{A}^{(m-1)}(\mathbf{x}), \overline{\mathbf{x}}^{(m-2)}\right] \middle| \mathbf{A}^{(m-2)}(\mathbf{x}), \overline{\mathbf{X}}^{(m-2)}\right] \\ &\stackrel{4}{=} \mathbb{E}\left[\mathbb{E}\left[Y(\mathbf{x}) \middle| \mathbf{A}^{(m-1)}(\mathbf{x}), \overline{\mathbf{x}}^{(m-1)}\right] \middle| \mathbf{A}^{(m-2)}(\mathbf{x}), \overline{\mathbf{x}}^{(m-2)}\right] \\ &= \mathbb{E}\left[\eta_0^m(\mathbf{A}^{(m-1)}(\mathbf{x}), \overline{\mathbf{x}}^{(m-1)}) \middle| \mathbf{A}^{(m-2)}(\mathbf{x}), \overline{\mathbf{x}}^{(m-2)}\right] \\ &\stackrel{5}{=} \mathbb{E}\left[\nu^m(\mathbf{A}^{(m-1)}(\mathbf{x}), \overline{\mathbf{x}}^{(m-1)}) \middle| \mathbf{A}^{(m-2)}(\mathbf{x}), \overline{\mathbf{x}}^{(m-2)}\right] \\ &= \nu_0^{m-1}(\mathbf{A}^{(m-2)}(\mathbf{x}), \overline{\mathbf{x}}^{(m-2)}), \end{split}$$

where

- $\stackrel{4}{=}$  holds since  $Y(\mathbf{x}) \perp \overline{X}_i | \mathbf{A}^{(i)}(\mathbf{x}), \overline{\mathbf{X}}^{(i-1)}$  for  $i = 1, 2, \cdots, m-1$  by Def. B.4 and the positivity condition in Eq. (B.8).
- $\stackrel{5}{=}$  holds since  $\eta^m(\mathbf{A}^{(m-1)}(\mathbf{x}), \overline{\mathbf{x}}^{(m-1)}) = \nu_0^m(\mathbf{A}^{(m-1)}(\mathbf{x}), \overline{\mathbf{x}}^{(m-1)}).$

Now, we make an induction hypothesis as follow: For any  $i + 1 \in \{m - 1 \dots, 3\}$ , suppose the following holds:

$$\eta_0^{i+1}(\mathbf{A}^{(i)}(\mathbf{x}), \overline{\mathbf{x}}^{(i)}) = \nu_0^{i+1}(\mathbf{A}^{(i)}(\mathbf{x}), \overline{\mathbf{x}}^{(i)}).$$

Then,

$$\begin{split} \eta_0^i(\mathbf{A}^{(i-1)}(\mathbf{x}), \overline{\mathbf{x}}^{(i-1)}) &\coloneqq \mathbb{E}\left[Y(\mathbf{x})|\mathbf{A}^{(i-1)}(\mathbf{x}), \overline{\mathbf{x}}^{(i-1)}\right] \\ &= \mathbb{E}\left[\mathbb{E}\left[Y(\mathbf{x})|\mathbf{A}^{(i)}(\mathbf{x}), \overline{\mathbf{x}}^{(i-1)}\right] \left|\mathbf{A}^{(i-1)}(\mathbf{x}), \overline{\mathbf{x}}^{(i-1)}\right]\right] \\ &\stackrel{7}{=} \mathbb{E}\left[\mathbb{E}\left[Y(\mathbf{x})|\mathbf{A}^{(i)}(\mathbf{x}), \overline{\mathbf{x}}^{(i)}\right] \left|\mathbf{A}^{(i-1)}(\mathbf{x}), \overline{\mathbf{x}}^{(i-1)}\right]\right] \\ &= \mathbb{E}\left[\eta_0^{i+1}(\mathbf{A}^{(i)}(\mathbf{x}), \overline{\mathbf{x}}^{(i)}) \left|\mathbf{A}^{(i-1)}(\mathbf{x}), \overline{\mathbf{x}}^{(i-1)}\right]\right] \\ &\stackrel{8}{=} \mathbb{E}\left[\nu_0^{i+1}(\mathbf{A}^{(i)}(\mathbf{x}), \overline{\mathbf{x}}^{(i)}) \left|\mathbf{A}^{(i-1)}(\mathbf{x}), \overline{\mathbf{x}}^{(i-1)}\right]\right] \\ &= \nu_0^i(\mathbf{A}^{(i-1)}(\mathbf{x}), \overline{\mathbf{X}}^{(i-1)}, \overline{x}_i \backslash z_i), \end{split}$$

where

•  $\stackrel{7}{=}$  holds since  $Y(\mathbf{x}) \perp \overline{X}_i | \mathbf{A}^{(i)}(\mathbf{x}), \overline{\mathbf{X}}^{(i-1)}$  for  $i = 1, 2, \cdots, m-1$  by Def. B.4 and the positivity condition in Eq. (B.8).

•  $\stackrel{8}{=}$  by induction hypothesis.

Therefore,  $\eta^i = \nu^i$  for all  $i = 1, 2, \cdots, m$ . This completes the proof.

**Theorem B.4** (Identification through AC-gMTI-PO). Suppose the condition AC-gMTI-PO in Def. B.4 holds. For  $i = 1, 2, \dots, m$ , let  $P^i$  denote a distribution defined as follow:

$$P^{i}(\mathbf{A}^{(i)}, \overline{\mathbf{X}}^{(i)} \setminus Z_{i}) \coloneqq P(\mathbf{A}^{(i)}(z_{i}), (\overline{\mathbf{X}}^{(i)} \setminus Z_{i})(z_{i})), \text{ for } i = 1, 2, \cdots, m-1$$
$$P^{m}(\mathbf{A}^{(m)}, \mathbf{X} \setminus Z_{m}, Y) \coloneqq P(\mathbf{A}^{(m)}(z_{m}), (\mathbf{X} \setminus Z)(z_{m}), Y(z_{m})).$$

Assume the following positivity condition holds: For  $\forall \overline{x}_m, \overline{x}_{m+1}, \mathbf{a}^{(m-1)}, \overline{\mathbf{x}}^{(m-1)} \in \mathfrak{D}_{\overline{X}_m, \overline{X}_{m+1}, \mathbf{A}^{(m-1)}, \overline{\mathbf{X}}^{(m-1)}}$ 

$$p_{\overline{X}_m \setminus Z_m, \overline{X}_{m+1} | \mathbf{A}^{(m-1)}, \overline{\mathbf{X}}^{(m-1)}}^{m}(\overline{x}_m \setminus z_m, \overline{x}_{m+1} | \mathbf{a}^{(m-1)}, \overline{\mathbf{x}}^{(m-1)}) > 0,$$
(B.9)

and for  $i = 1, 2, \cdots, m - 1$ ,

$$\frac{p^{i+1}(a_i|\mathbf{a}^{(i-1)}, \overline{\mathbf{x}}^{(i-1)})}{p^i(a_i|\mathbf{a}^{(i-1)}, \overline{\mathbf{x}}^{(i-1)})} p^{i+1}(\overline{x}_i|\mathbf{a}^{(i)}, \overline{\mathbf{x}}^{(i-1)}) > 0, \ \forall \mathbf{a}^{(i-1)}, \overline{\mathbf{x}}^{(i)} \in \overline{\mathcal{X}}^{(i)} \times \mathcal{A}^{(i-1)}.$$
(B.10)

Then, the query  $\mathbb{E}[Y(\mathbf{x})]$  is identifiable from distributions  $P^1, \dots, P^m$ , and given as follow: Let

$$\mu_0^m(\mathbf{A}^{(m-1)}, \mathbf{x} \backslash z_m) \coloneqq \mathbb{E}_{P^m} \left[ Y | \mathbf{A}^{(m-1)}, \mathbf{x} \backslash z_m \right]$$
$$\mu_0^{m-1}(\mathbf{A}^{(m-2)}, \overline{\mathbf{x}}^{(m-2)}) \coloneqq \mathbb{E}_{P^{m-1}} \left[ \mu_0^m(\mathbf{A}^{(m-1)}, \mathbf{x} \backslash z_m) \middle| \mathbf{A}^{(m-2)}, \overline{\mathbf{x}}^{(m-2)} \right],$$

and for  $i = m - 2, \cdots, 2$ 

$$\mu_0^i(\mathbf{A}^{(i-1)}, \overline{\mathbf{x}}^{(i-1)}) \coloneqq \mathbb{E}_{P^i}\left[\mu^{i+1}(\mathbf{A}^{(i)}, \overline{\mathbf{x}}^{(i)}) \middle| \mathbf{A}^{(i-1)}, \overline{\mathbf{x}}^{(i-1)}\right].$$

Then,

$$\mathbb{E}\left[Y(\mathbf{x})\right] = \mathbb{E}_{P^1}\left[\mu^2(A_1, \overline{x}_1)\right].$$

*Proof of Theorem B.4.* We first show that the positivity conditions in Eqs (B.7, B.8) match with Eqs. (B.9, B.10). Eq. (B.7) = Eq. (B.9) holds since

$$\begin{aligned} \mathsf{Eq.} \ (\mathbf{B.7}) &= p_{(\overline{X}_m \setminus Z_m(z_m), \overline{X}_{m+1}(z_m) | \mathbf{A}^{(m-1)}(\mathbf{x}), \overline{\mathbf{X}}^{(m-1)}}(\overline{x}_m \setminus z_m, \overline{x}_{m+1} | \mathbf{a}^{(m-1)}, \overline{\mathbf{x}}^{(m-1)}) \\ &= p_{(\overline{X}_m \setminus Z_m(z_m), \overline{X}_{m+1}(z_m) | \mathbf{A}^{(m-1)}(\overline{\mathbf{x}}^{(m-1)}), \overline{\mathbf{X}}^{(m-1)}(\overline{x}_m \setminus z_m, \overline{x}_{m+1} | \mathbf{a}^{(m-1)}, \overline{\mathbf{x}}^{(m-1)}) \\ &= p_{(\overline{X}_m \setminus Z_m(z_m), \overline{X}_{m+1}(z_m) | \mathbf{A}^{(m-1)}, \overline{\mathbf{X}}^{(m-1)}(\overline{x}_m \setminus z_m, \overline{x}_{m+1} | \mathbf{a}^{(m-1)}, \overline{\mathbf{x}}^{(m-1)}) \\ &= p_{(\overline{X}_m \setminus Z_m(z_m), \overline{X}_{m+1}(z_m) | \mathbf{A}^{(m-1)}(z_m), \overline{\mathbf{X}}^{(m-1)}(z_m)}(\overline{x}_m \setminus z_m, \overline{x}_{m+1} | \mathbf{a}^{(m-1)}, \overline{\mathbf{x}}^{(m-1)}) \\ &= \mathsf{Eq.} \ (\mathbf{B.9}). \end{aligned}$$

Eq. (B.8) = Eq. (B.10) holds since

$$\begin{split} & \mathsf{Eq.} \, (\mathbf{B.8}) \\ &= \frac{p_{\overline{X}_i,A_i(\mathbf{x})|\mathbf{A}^{(i-1)}(\mathbf{x}),\overline{\mathbf{X}}^{(i-1)}(\overline{x}_i,a_i|\mathbf{a}^{(i-1)},\overline{\mathbf{x}}^{(i-1)})}{p_{A_i(\mathbf{x})|\mathbf{A}^{(i-1)}(\mathbf{x}),\overline{\mathbf{X}}^{(i-1)}(\overline{x}_i,a_i|\mathbf{a}^{(i-1)},\overline{\mathbf{x}}^{(i-1)})}{p_{A_i(\mathbf{z}^{(i)})|\mathbf{A}^{(i-1)}(\overline{\mathbf{z}^{(i-1)}}),\overline{\mathbf{X}}^{(i-1)}(\overline{\mathbf{x}}_i,a_i|\mathbf{a}^{(i-1)},\overline{\mathbf{x}}^{(i-1)})}{p_{A_i(\mathbf{z}^{(i)})|\mathbf{A}^{(i-1)}(\overline{\mathbf{x}}^{(i-1)},\overline{\mathbf{x}}^{(i-1)}(\overline{\mathbf{x}}^{(i-1)},\overline{\mathbf{x}}^{(i-1)})}{p_{A_i(\mathbf{z}^{(i)})|\mathbf{A}^{(i-1)},\overline{\mathbf{x}}^{(i-1)}(a_i|\mathbf{a}^{(i-1)},\overline{\mathbf{x}}^{(i-1)})}{p_{A_i(\mathbf{z}^{(i)})|\mathbf{A}^{(i-1)},\overline{\mathbf{x}}^{(i-1)}(a_i|\mathbf{a}^{(i-1)},\overline{\mathbf{x}}^{(i-1)})}}{p_{A_i(\mathbf{z}^{(i)})|\mathbf{A}^{(i-1)}(z_i,\overline{\mathbf{x}}^{(i-1)})(z_i,\overline{\mathbf{x}}^{(i-1)})}{p_{A_i(\mathbf{z}^{(i)})|\mathbf{A}^{(i-1)}(z_i,\overline{\mathbf{x}}^{(i-1)})(z_i,\overline{\mathbf{x}}^{(i-1)})}{p_{A_i(\mathbf{z}^{(i)})|\mathbf{A}^{(i-1)}(z_i,\overline{\mathbf{x}}^{(i-1)})(z_i,\overline{\mathbf{x}}^{(i-1)})(z_{i+1})}(\overline{x}^{(i,a}|\mathbf{a}^{(i-1)},\overline{\mathbf{x}}^{(i-1)})}{p_{A_i(z_i)|\mathbf{A}^{(i-1)}(z_i,\overline{\mathbf{x}}^{(i-1)})(z_i,\overline{\mathbf{x}}^{(i-1)})(z_{i+1})}(\overline{x}^{(i,a}|\mathbf{a}^{(i-1)},\overline{\mathbf{x}}^{(i-1)})})} \\ &= \frac{p_{\overline{X}_i(z_{i+1}),A_i(z_{i+1})|\mathbf{A}^{(i-1)}(z_{i+1}),\overline{\mathbf{x}}^{(i-1)}(z_{i+1})}(a_i|\mathbf{a}^{(i-1)},\overline{\mathbf{x}}^{(i-1)})}{p_{A_i(z_i)|\mathbf{A}^{(i-1)}(z_i,\overline{\mathbf{x}}^{(i-1)})(z_{i+1})}(a_i|\mathbf{a}^{(i-1)},\overline{\mathbf{x}}^{(i-1)})}}{p_{A_i(z_i)|\mathbf{A}^{(i-1)}(z_i,\overline{\mathbf{x}}^{(i-1)})(z_i)}(a_i|\mathbf{a}^{(i-1)},\overline{\mathbf{x}}^{(i-1)})}p_{\overline{X}_i(z_{i+1})|\mathbf{A}^{(i)}(z_{i+1}),\overline{\mathbf{x}}^{(i-1)}(z_{i+1})}(a_i|\mathbf{a}^{(i-1)},\overline{\mathbf{x}}^{(i-1)})}} \\ &= \frac{p_{A_i(z_{i+1})|\mathbf{A}^{(i-1)}(z_i,\overline{\mathbf{x}}^{(i-1)})(z_i)}(a_i|\mathbf{a}^{(i-1)},\overline{\mathbf{x}}^{(i-1)})}{p_{A_i(z_i)|\mathbf{A}^{(i-1)}(z_i,\overline{\mathbf{x}}^{(i-1)})}(z_i)}(a_i|\mathbf{a}^{(i-1)},\overline{\mathbf{x}}^{(i-1)})}}{p_{\overline{X}_i(z_{i+1})|\mathbf{A}^{(i)}(z_{i+1}),\overline{\mathbf{X}}^{(i-1)})}(z_i)}(\overline{x}_i|\mathbf{a}^{(i)},\overline{\mathbf{x}}^{(i-1)})})} \\ &= \frac{p_{i+1}(a_i|\mathbf{a}^{(i-1)},\overline{\mathbf{x}}^{(i-1)})}{p^i(a_i|\mathbf{a}^{(i-1)},\overline{\mathbf{x}}^{(i-1)})}}p^{i+1}(\overline{x}_i|\mathbf{a}^{(i)},\overline{\mathbf{x}}^{(i-1)})}) =: \mathsf{Eq.} \,(\mathsf{B.10}). \end{split}$$

Now, it suffices to show that

$$\nu_0^m(\mathbf{A}^{(m-1)}(\mathbf{x}), \mathbf{x} \backslash z_m) = \mu_0^m(\mathbf{A}^{(m-1)}(\mathbf{x}), \mathbf{x} \backslash z_m)$$
$$\nu_0^{i+1}(\mathbf{A}^{(i)}(\mathbf{x}), \overline{\mathbf{x}}^{(i)}) = \mu_0^{i+1}(\mathbf{A}^{(i)}, \overline{\mathbf{x}}^{(i)}), \text{ for } i = m - 2, \cdots, 1,$$

where  $\nu^i$  for  $i = m, \cdots, 2$  are defined in Lemma B.3.

First,

$$\nu_0^m(\mathbf{A}^{(m-1)}(\mathbf{x}), \mathbf{x} \backslash z_m) \coloneqq \mathbb{E}\left[Y(\mathbf{x}) | \mathbf{A}^{(m-1)}(\mathbf{x}), (\mathbf{X} \backslash Z_m)(z_m) = \mathbf{x} \backslash z_m\right]$$
$$= \mathbb{E}\left[Y(z_m) | \mathbf{A}^{(m-1)}(z_m), (\mathbf{X} \backslash Z_m)(z_m) = \mathbf{x} \backslash z_m\right]$$
$$= \mathbb{E}_{P^m}\left[Y | \mathbf{A}^{(m-1)}, \mathbf{x} \backslash z_m\right]$$
$$=: \mu^m(\mathbf{A}^{(m-1)}, \mathbf{x} \backslash z_m).$$

Also,

$$\begin{split} \nu_{0}^{m-1}(\mathbf{A}^{(m-2)}(\mathbf{x}), \overline{\mathbf{x}}^{(m-2)}) &\coloneqq \mathbb{E} \left[ \nu_{0}^{m}(\mathbf{A}^{(m-1)}(\mathbf{x}), \mathbf{x} \setminus z_{m}) \Big| \mathbf{A}^{(m-2)}(\mathbf{x}), \overline{\mathbf{x}}^{(m-2)} \right], \\ &= \mathbb{E} \left[ \mu_{0}^{m}(\mathbf{A}^{(m-1)}, \mathbf{x} \setminus z_{m}) \Big| \mathbf{A}^{(m-2)}(\mathbf{x}), \overline{\mathbf{x}}^{(m-2)} \right] \\ &= \mathbb{E} \left[ \mu_{0}^{m}(\mathbf{A}^{(m-1)}(\mathbf{z}^{(m-1)}), \mathbf{x} \setminus z_{m}) \Big| \mathbf{A}^{(m-2)}(\mathbf{x}), \overline{\mathbf{x}}^{(m-2)} \right] \\ &= \mathbb{E} \left[ \mu_{0}^{m}(\mathbf{A}^{(m-1)}(z_{m-1}), \mathbf{x} \setminus z_{m}) \Big| \mathbf{A}^{(m-2)}(z_{m-1}), \overline{\mathbf{X}}^{(m-2)}(z_{m-1}) = \overline{\mathbf{x}}^{(m-1)} \right] \\ &= \mathbb{E}_{P^{m-1}} \left[ \mu_{0}^{m}(\mathbf{A}^{(m-1)}, \mathbf{x} \setminus z_{m}) \Big| \mathbf{A}^{(m-2)}, \overline{\mathbf{x}}^{(m-2)} \right] \\ &=: \mu_{0}^{m-1}(\mathbf{A}^{(m-2)}, \overline{\mathbf{x}}^{(m-2)}). \end{split}$$

Finally, for  $i + 1 \in \{m - 1, \dots, 3\}$ , suppose

$$\nu_0^{i+1}(\mathbf{A}^{(i)}(\mathbf{x}), \overline{\mathbf{x}}^{(i)}) = \mu_0^{i+1}(\mathbf{A}^{(i)}, \overline{\mathbf{x}}^{(i)}).$$

Then,

$$\begin{split} \nu_{0}^{i}(\mathbf{A}^{(i-1)}(\mathbf{x}), \overline{\mathbf{x}}^{(i-1)}) &= \mathbb{E}\left[\nu_{0}^{i+1}(\mathbf{A}^{(i)}(\mathbf{x}), \overline{\mathbf{x}}^{(i)}) \middle| \mathbf{A}^{(i-1)}(\mathbf{x}), \overline{\mathbf{x}}^{(i-1)} \right] \\ &= \mathbb{E}\left[\mu_{0}^{i+1}(\mathbf{A}^{(i)}, \overline{\mathbf{x}}^{(i)}) \middle| \mathbf{A}^{(i-1)}(\mathbf{x}), \overline{\mathbf{x}}^{(i-1)} \right] \\ &= \mathbb{E}\left[\mu_{0}^{i+1}(\mathbf{A}^{(i)}(z_{i}), \overline{\mathbf{x}}^{(i)}) \middle| \mathbf{A}^{(i-1)}(z_{i}), \overline{\mathbf{X}}^{(i-1)}(z_{i}) = \overline{\mathbf{x}}^{(i-1)} \right] \\ &= \mathbb{E}_{P^{i}}\left[\mu_{0}^{i+1}(\mathbf{A}^{(i)}, \overline{\mathbf{x}}^{(i)}) \middle| \mathbf{A}^{(i-1)}, \overline{\mathbf{x}}^{(i-1)} \right] \\ &= \mu_{0}^{i}(\mathbf{A}^{(i-1)}, \overline{\mathbf{x}}^{(i-1)}). \end{split}$$

## 

## C. Proofs

## **C.1.** Preliminaries

**Lemma C.1 (Continuous Mapping Theorem for**  $L_2(P)$ ). Let  $X_n, X$  denote a random sequence defined on a metric space S. Suppose a function  $g: S \to S'$  (where S' is another metric space) is continuous almost everywhere. Suppose g is bounded. Then,

$$X_n \stackrel{L_2(P)}{\to} X \implies g(X_n) \stackrel{L_2(P)}{\to} g(X)$$

**Proof of Lemma C.1.** We first note that  $X_n \xrightarrow{L_2(P)} X$  implies  $X_n \xrightarrow{p} X$ . Then, by continuous mapping theorem,  $g(X_n) \xrightarrow{p} g(X)$ . Then,

$$\lim_{n \to \infty} \|g(X_n) - g(X)\|^2 = \lim_{n \to \infty} \int_{\mathcal{X}} |g(X_n) - g(X)|^2 \ d[P] \stackrel{*}{=} \int_{\mathcal{X}} \lim_{n \to \infty} |g(X_n) - g(X)|^2 \ d[P] = 0,$$

where the equation  $\stackrel{*}{=}$  holds by dominated convergence theorem in  $L_2(P)$  space, which is applicable since  $g(X_n), g(X)$  are bounded functions (from the given condition) and  $X_n \stackrel{p}{\to} X$ .

**Lemma C.2** (Asymptotic Unbiasedness implies Consistency). Suppose an estimator  $T_N$  is asymptotically unbiased to  $\mu$ ; *i.e.*,  $\mathbb{E}_P[T_N - \mu] \to 0$  as  $N \to \infty$ . Suppose an estimator has vanishing variance; *i.e.*,  $var(T_N) \to 0$  as  $N \to \infty$ . Then,  $T_N$  is a consistent estimator of  $\mu$ .

Proof of Lemma C.2. By Markov inequality,

$$P(|T_N - \mu| > \epsilon) = P((T_N - \mu)^2 > \epsilon^2) \le \mathbb{E}_P\left[(T_N - \mu)^2\right] / \epsilon^2.$$

Also, for  $\mu_N \coloneqq \mathbb{E}_P[T_N]$ ,

$$\mathbb{E}_{P}\left[ (T_{N} - \mu)^{2} \right] \leq 2\mathbb{E}_{P}\left[ (T_{N} - \mu_{N})^{2} \right] + 2(\mu_{N} - \mu)^{2}$$
  
=  $2\mathbb{V}_{P}\left[ T_{N} \right] + 2(\mu_{N} - \mu)^{2}$   
 $\rightarrow 0.$ 

where  $\operatorname{var}(T_N) + (\mu_N - \mu) \to 0$  by the given assumptions that  $\operatorname{var}(T_N) \to 0$  and  $\mathbb{E}_P[T_N - \mu] = \mu_N - \mu \to 0$  as  $N \to \infty$ .

**Lemma C.3** (Decomposition (Kennedy et al., 2020, Lemma 2)). Let  $f_{\eta} \equiv f(\mathbf{V}; \eta)$  denote a finite and continuous functional and  $\eta$  denote its nuisances. For some samples  $D \sim P$ , let  $T \equiv \mathbb{E}_D[f_{\eta}]$ . Let  $\theta_0 \equiv \mathbb{E}_P[f_{\eta_0}]$  for some  $\eta_0$ . Let  $\mathbb{E}_{D-P}[f_{\eta}] \equiv \mathbb{E}_D[f_{\eta}] - \mathbb{E}_P[f_{\eta}]$ . Then, the following decomposition holds:

$$\mathbb{E}_{D}[f_{\eta}] - \theta_{0} = \mathbb{E}_{D-P}[f_{\eta_{0}}] + \mathbb{E}_{D-P}[f_{\eta} - f_{\eta_{0}}] + \mathbb{E}_{P}[f_{\eta} - f_{\eta_{0}}].$$
(C.1)

Suppose further that

- 1. Samples used for estimating  $\eta$  are independent and separate from D; and
- 2.  $\|\eta \eta_0\| = o_P(1)$ .

Then, Eq. (C.1) reduces to

$$\mathbb{E}_D\left[f_\eta\right] - \theta_0 = R + \mathbb{E}_P\left[f_\eta - f_{\eta_0}\right],\tag{C.2}$$

where R is a random variable such that  $\sqrt{nR}$  converges in distribution to a mean-zero normal random variable, where  $n \equiv |D|$ .

Proof of Lemma C.3. We first prove the equality in Eq. (C.1).

$$\mathbb{E}_{D} [f_{\eta}] - \theta_{0} = \mathbb{E}_{D} [f_{\eta}] - \mathbb{E}_{P} [f_{\eta_{0}}]$$

$$= \mathbb{E}_{D-P} [f_{\eta}] + \mathbb{E}_{P} [f_{\eta} - f_{\eta_{0}}]$$

$$= \underbrace{\mathbb{E}_{D-P} [f_{\eta_{0}}]}_{\equiv A} + \underbrace{\mathbb{E}_{D-P} [f_{\eta} - f_{\eta_{0}}]}_{\equiv B} + \mathbb{E}_{P} [f_{\eta} - f_{\eta_{0}}].$$

We now prove Eq. (C.2).

- A converges in distribution to the zero-mean normal distribution at  $\sqrt{n}$  rate by the central limit theorem.
- We note that a given condition  $\|\eta \eta_0\| = o_P(1)$  implies  $\|f_\eta f_{\eta_0}\| = o_P(1)$  by continuous mapping theorem for  $L_2(P)$  in Lemma C.1. In particular, Lemma C.1 is applicable since  $f_\eta$ ,  $f_{\eta_0}$  is a bounded and continuous function, and  $\|\eta \eta_0\| = o_P(1)$ . Then, B converges to zero at  $o_P(1/\sqrt{N})$  rate by (Kennedy et al., 2020, Lemma 2).

Finally, define  $R \equiv A + B$ . Then, the proof completes by applying Slutsky's theorem.

#### 

#### C.2. Proof of Theorem 1

**Definition 1** (Adjustment criterion for Treatment-Treatment Interaction (AC-TTI)). A set  $\{C_1, W\}$  is said to satisfy the *adjustment criterion for treatment-treatment interaction (AC-TTI)* w.r.t  $\{(X_1, X_2), Y\}$  in G if

- 1.  $(\{C_1, W\} \perp X_2 | X_1)_{G_{\overline{X_1, X_2}}}$ , i.e., there are no direct paths from  $X_2$  to  $\{C_1, W\}$  in  $G_{\overline{X_1, X_2}}$ ; and
- 2.  $(Y \perp X_1 | C_1, W, X_2)_{G_{X_1 \overline{X_2}}}$ , i.e., the back-door paths from  $X_1$  to Y are blocked by  $\{C_1, W\}$  in  $G_{\overline{X_2}}$ .

Assumption 1 (Positivity Assumption for AC-TTI).  $P_{x_1}(C_1, W), P_{x_2}(C_1, W), P_{x_2}(X_1|C_1, W)$  are strictly positive distributions  $\forall x_1, x_2 \in \mathfrak{D}_{X_1, X_2}$ .

**Theorem 1** (Identification through AC-TTI). Suppose AC-TTI in Def. 1 and Assumption 1 hold. Then,  $\mathbb{E}[Y|do(x_1, x_2)]$  is identifiable from  $P_{rand(X_1)}(C_1, W)$  and  $P_{rand(X_2)}(C_1, W, X_1, Y)$  and the expression is:

$$\mathbb{E}[Y|do(x_1, x_2)] = \mathbb{E}_{P_{x_1}}\left[\mathbb{E}_{P_{x_2}}[Y|C_1, W, x_1]\right].$$
(1)

**Proof of Theorem 1**.

$$\begin{split} \mathbb{E}\left[Y|do(x_{1},x_{2})\right] &= \mathbb{E}\left[\mathbb{E}\left[Y|do(x_{1},x_{2}),C_{1},W\right]|do(x_{1},x_{2})\right] \\ &= \frac{1}{2} \mathbb{E}\left[\mathbb{E}\left[Y|do(x_{2}),x_{1},C_{1},W\right]|do(x_{1},x_{2})\right] \\ &= \mathbb{E}\left[\mathbb{E}_{P_{x_{2}}}\left[Y|x_{1},C_{1},W\right]|do(x_{1},x_{2})\right] \\ &= \mathbb{E}\left[\mathbb{E}_{P_{x_{2}}}\left[Y|x_{1},C_{1},W\right]|do(x_{1})\right] \\ &= \mathbb{E}_{P_{x_{1}}}\left[\mathbb{E}_{P_{x_{2}}}\left[Y|x_{1},C_{1},W\right]\right], \end{split}$$

where  $\stackrel{1}{=}$  holds by the condition 2 which implies Rule 2 of *do*-calculus, and  $\stackrel{2}{=}$  holds by the condition 1 in AC-TTI which implies Rule 3 of *do*-calculus.

### C.3. Proof of Theorem 2 and Corollary 2

**Definition 2** (Nuisances for TTI). Nuisance functions for the AC-TTI functional in Eq. (1) are defined as follows: For a fixed  $x_1, x_2 \in \mathfrak{D}_{X_1, X_2}$  where  $x_1, x_2$  are specified in Eq. (1),  $\pi_0 \coloneqq \pi_0(C_1, X_1, W) \coloneqq \frac{P_{x_1}(W|C_1)}{P_{x_2}(W, X_1|C_1)}$ . Also,  $\mu_0 \coloneqq \mu_0(C_1, X_1, W) \coloneqq \mathbb{E}_{P_{x_2}}[Y|X_1, W, C_1]$ . We will use  $\pi \coloneqq \pi(C_1, X_1, W) > 0$  and  $\mu \coloneqq \mu(C_1, X_1, W)$  to denote arbitrary<sup>3</sup> finite functions.

**Definition 3** (Estimators for TTI). Let  $D_1$  and  $D_2$  denote two separate samples following the distributions  $P_{\operatorname{rand}(X_1)}(C_1, W)$  and  $P_{\operatorname{rand}(X_2)}(C_1, W, X_1, Y)$ , respectively. For fixed  $x_1, x_2 \in \mathfrak{D}_{X_1, X_2}$ , we define  $D_{x_1}$  and  $D_{x_2}$  as subsamples of  $D_1$  and  $D_2$  such that  $X_1 = x_1$  and  $X_2 = x_2$ . Let  $\mu$  and  $\pi$  denote the nuisances as defined in Definition 2. We now introduce the {REG, PW, DML} estimators for the AC-TTI-functional specified in Equation (1) as follows:

$$T^{reg} \coloneqq \mathbb{E}_{D_{x_1}} \left[ \mu(W, C_1, x_1) \right],$$
  

$$T^{pw} \coloneqq \mathbb{E}_{D_{x_2}} \left[ \pi(W, C_1, X_1) \mathbb{1}_{x_1}(X_1) Y \right],$$
  

$$T^{dml} \coloneqq \mathbb{E}_{D_{x_1}} \left[ \pi \mathbb{1}_{x_1}(X_1) \{Y - \mu\} \right] + \mathbb{E}_{D_1} \left[ \mu(W, C_1, x_1) \right) \right]$$

**Assumption 2** (Sample-splitting). Samples for training nuisances and evaluating the estimators equipped with the trained nuisance are separate and independent.

Assumption 3 ( $L_2$  consistency of nuisances). Estimated nuisances are  $L_2$  consistent; i.e.,  $\forall i \in \{1, 2\}, \forall x_i \in \mathfrak{D}_{X_i}$ ,

$$\begin{aligned} \|\mu(W,C_1,x_1) - \mu_0(W,C_1,x_1)\|_{P_{x_i}} &= o_{P_{x_i}}(1), \\ \|\pi(W,C_1,X_1) - \pi_0(W,C_1,X_1)\|_{P_{x_2}} &= o_{P_{x_2}}(1). \end{aligned}$$

**Theorem 2** (Error analysis of the estimators). Under Assumptions (1,2,3,4) and AC-TTI in Def. 1, the error of the estimators in Def. 3, denoted  $\epsilon^{est} := T^{est} - \mathbb{E}[Y|do(x_1, x_2)]$  for  $est \in \{reg, pw, dm\}$  are:

$$\begin{split} \epsilon^{reg} &= R_1 + O_{P_{x_1}} \left( \| \mu - \mu_0 \| \right), \\ \epsilon^{pw} &= R_2 + O_{P_{x_2}} \left( \| \pi - \pi_0 \| \right), \\ \epsilon^{dml} &= R_1 + R_2 + O_{P_{x_2}} \left( \| \pi - \pi_0 \| \| \mu - \mu_0 \| \right), \end{split}$$

where  $R_i$  is a random variable such that  $\sqrt{n_i}R_i$  converges in distribution to a zero-mean normal random variable, where  $n_i := |D_{x_i}|$  for  $i \in \{1, 2\}$ .

*Proof of Theorem 2.* We provide error analyses for each estimators:

Analysis for  $T^{reg}$ .

<sup>&</sup>lt;sup>3</sup>Throughout the paper,  $\mu$ ,  $\pi$  are understood as estimated nuisances for  $\mu_0$ ,  $\pi_0$ .

We first note that

$$\mathbb{E}_{P_{x_1}}\left[\mu_0(W, C_1, \mathbf{x})\right] = \mathbb{E}_{P_{x_1}}\left[\mathbb{E}_{P_{x_2}}\left[Y|W, C_1, x_1\right]\right] = \mathbb{E}\left[Y|do(x_1, x_2)\right],$$

where the last equation holds by Theorem 1. By Lemma C.3,

$$T^{reg} - \mathbb{E} \left[ Y | do(x_1, x_2) \right] = T^{reg} - \mathbb{E}_{P_{x_1}} \left[ \mu_0(W, C_1, \mathbf{x}) \right]$$
  
=  $\underbrace{\mathbb{E}_{P_{x_1} - D_{x_1}} \left[ \mu_0(W, C_1, \mathbf{x}) \right] + \mathbb{E}_{P_{x_1} - D_{x_1}} \left[ \mu(W, C_1, \mathbf{x}) - \mu_0(W, C_1, \mathbf{x}) \right]}_{:= R_1} + \mathbb{E}_{P_{x_1}} \left[ \mu(W, C_1, \mathbf{x}) - \mu_0(W, C_1, \mathbf{x}) \right]$   
=  $\frac{1}{R_1} + \mathbb{E}_{P_{x_1}} \left[ \mu(W, C_1, \mathbf{x}) - \mu_0(W, C_1, \mathbf{x}) \right]$   
=  $\frac{1}{R_1} + \mathbb{E}_{P_{x_1}} \left[ \mu(W, C_1, \mathbf{x}) - \mu_0(W, C_1, \mathbf{x}) \right]$   
=  $\frac{1}{R_1} + O_{P_{x_1}} \left( \| \mu_0 - \mu \| \right),$ 

where

- $\stackrel{1}{=}$  holds by Lemma C.3.
- $\stackrel{2}{=}$  holds by Cauchy-Schwartz inequality.

## Analysis for $T^{pw}$ .

We first note that

$$\begin{split} \mathbb{E}_{P_{x_2}} \left[ \pi_0(W, C_1, X_1) \mathbb{1}_{x_1}(X_1) Y \right] &= \mathbb{E}_{P_{x_2}} \left[ \pi_0(W, C_1, X_1) \mathbb{1}_{x_1}(X_1) \mu_0(W, C_1, X_1) \right] \\ &= \mathbb{E}_{P_{x_2}} \left[ \frac{P_{x_1}(W|C_1)}{P_{x_2}(W, X_1|C_1)} \mathbb{1}_{x_1}(X_1) \mu_0(W, C_1, X_1) \right] \\ &\stackrel{3}{=} \mathbb{E}_{P_{x_2}} \left[ \frac{P_{x_1}(W|C_1) P_{x_1}(C_1)}{P_{x_2}(W, X_1|C_1) P_{x_2}(C_1)} \mathbb{1}_{x_1}(X_1) \mu_0(W, C_1, X_1) \right] \\ &= \mathbb{E}_{P_{x_2}} \left[ \frac{P_{x_1}(W, C_1)}{P_{x_2}(X_1|W, C_1) P_{x_2}(W, C_1)} \mathbb{1}_{x_1}(X_1) \mu_0(W, C_1, X_1) \right] \\ &= \mathbb{E}_{P_{x_2}} \left[ \frac{P_{x_1}(W, C_1)}{P_{x_2}(W, C_1)} \mu_0(W, C_1, x_1) \right] \\ &= \mathbb{E}_{P_{x_1}} \left[ \mu_0(W, C_1, x_1) \right] \\ &= \mathbb{E}_{P_{x_1}} \left[ \mu_0(W, C_1, x_2) \right], \end{split}$$

where

•  $\stackrel{3}{=}$  holds by Assumption 4.

•  $\stackrel{4}{=}$  holds by Theorem 1.

By applying Lemma C.3,

$$\begin{split} T^{pw} &- \mathbb{E}\left[Y|do(x_{1},x_{2})\right] \\ &= T^{pw} - \mathbb{E}_{P_{x_{2}}}\left[\pi_{0}(W,C_{1},X_{1})\mathbbm{1}_{x_{1}}(X_{1})Y\right] \\ &= \underbrace{\mathbb{E}_{P_{x_{2}}-D_{x_{2}}}\left[\pi_{0}(W,C_{1},X_{1})\mathbbm{1}_{x_{1}}(X_{1})Y\right] + \mathbb{E}_{P_{x_{2}}-D_{x_{2}}}\left[\left\{\pi_{0}(W,C_{1},X_{1}) - \pi(W,C_{1},X_{1})\right\}\mathbbm{1}_{x_{1}}(X_{1})Y\right] \\ &= \underbrace{\mathbb{E}_{P_{x_{2}}}\left[\left\{\pi_{0}(W,C_{1},X_{1}) - \pi(W,C_{1},X_{1})\right\}\mathbbm{1}_{x_{1}}(X_{1})Y\right] \\ &= R_{2}} \\ &+ \mathbb{E}_{P_{x_{2}}}\left[\left\{\pi_{0}(W,C_{1},X_{1}) - \pi(W,C_{1},X_{1})\right\}\mathbbm{1}_{x_{1}}(X_{1})Y\right] \\ &= R_{2} + \mathbb{E}_{P_{x_{2}}}\left[\left\{\pi_{0}(W,C_{1},X_{1}) - \pi(W,C_{1},X_{1})\right\}\mathbbm{1}_{x_{1}}(X_{1})Y\right] \\ &= R_{2} + O_{P_{x_{2}}}\left(\|\pi_{0} - \pi\|\right), \end{split}$$

- $\stackrel{5}{=}$  holds by Assumption 4.
- $\stackrel{6}{=}$  holds by Cauchy-Schwartz inequality and the setting where Y has a finite variance.

Analysis for  $T^{dml}$ .

Let

$$T^{dml} \coloneqq \underbrace{\mathbb{E}_{D_{x_2}}\left[\pi(W, C_1, X_1) \mathbb{1}_{x_1}(X_1) \{Y - \mu(W, C_1, X_1)\}\right]}_{:=T^{dml, 1}} + \underbrace{\mathbb{E}_{D_{x_1}}\left[\mu(W, C_1, x_1)\right)}_{:=T^{dml, 2}}.$$

Let

$$T_0^{dml,1} \coloneqq \mathbb{E}_{P_{x_2}} \left[ \pi_0(W, C_1, X_1) \mathbb{1}_{x_1}(X_1) \{ Y - \mu_0(W, C_1, X_1) \} \right]$$
$$T_0^{dml,2} \coloneqq \mathbb{E}_{P_{x_1}} \left[ \mu_0(W, C_1, x_1)) \right].$$

We note that  $T_0^{dml} \coloneqq T_0^{dml,1} + T_0^{dml,2} = \mathbb{E}\left[Y | do(x_1, x_2)\right]$ . We first apply the Lemma C.3 to  $T^{dml,1}$  and  $T^{dml,2}$  separately.

$$\begin{split} T^{dml,1} &- T_0^{dml,1} \\ &= \mathbb{E}_{P_{x_2} - D_{x_2}} \left[ \pi_0(W, C_1, x_1) \mathbb{1}_{x_1}(X_1) \{ Y - \mu_0(W, C_1, X_1) \} \right] \\ &+ \mathbb{E}_{P_{x_2} - D_{x_2}} \left[ \pi_0(W, C_1, x_1) \mathbb{1}_{x_1}(X_1) \{ Y - \mu_0(W, C_1, X_1) \} - \pi(W, C_1, x_1) \mathbb{1}_{x_1}(X_1) \{ Y - \mu(W, C_1, X_1) \} \right] \\ &+ \mathbb{E}_{P_{x_2}} \left[ \pi(W, C_1, x_1) \mathbb{1}_{x_1}(X_1) \{ Y - \mu(W, C_1, X_1) \} \right] \\ &= R_2 + \mathbb{E}_{P_{x_2}} \left[ \pi(W, C_1, x_1) \mathbb{1}_{x_1}(X_1) \{ Y - \mu(W, C_1, X_1) \} \right], \end{split}$$

where

$$\begin{split} R_2 &\coloneqq \mathbb{E}_{P_{x_2} - D_{x_2}} \left[ \pi_0(W, C_1, x_1) \mathbb{1}_{x_1}(X_1) \{ Y - \mu_0(W, C_1, X_1) \} \right] \\ &+ \mathbb{E}_{P_{x_2} - D_{x_2}} \left[ \pi_0(W, C_1, x_1) \mathbb{1}_{x_1}(X_1) \{ Y - \mu_0(W, C_1, X_1) \} - \pi(W, C_1, x_1) \mathbb{1}_{x_1}(X_1) \{ Y - \mu(W, C_1, X_1) \} \right] \end{split}$$

Also, by the proof for analyzing the error of  $T^{reg}$ ,

$$T^{dml,2} - T_0^{dml,2} = R_1 + \mathbb{E}_{P_{x_1}} \left[ \mu(W, C_1, x_1) - \mu_0(W, C_1, x_1) \right].$$

Then,

$$\begin{split} T^{dml} &- \mathbb{E}\left[Y|do(x_1, x_2)\right] \\ &= T^{dml,1} + T^{dml,2} - T_0^{dml,1} - T_0^{dml,2} \\ &= R_1 + R_2 + \mathbb{E}_{P_{x_2}}\left[\pi(W, C_1, x_1)\mathbbm{1}_{x_1}(X_1)\{Y - \mu(W, C_1, X_1)\}\right] + \mathbb{E}_{P_{x_1}}\left[\mu(W, C_1, x_1) - \mu_0(W, C_1, x_1)\right]. \end{split}$$

Note that

$$\begin{split} & \mathbb{E}_{P_{x_1}} \left[ \mu(W, C_1, x_1) - \mu_0(W, C_1, x_1) \right] \\ &= \mathbb{E}_{P_{x_1}} \left[ \frac{\mathbbm{1}_{x_1}(X_1)}{P_{x_1}(X_1|W, C_1)} \left\{ \mu(W, C_1, X_1) - \mu_0(W, C_1, X_1) \right\} \right] \\ &= \mathbb{E}_{P_{x_2}} \left[ \frac{P_{x_1}(W, C_1, X_1)}{P_{x_2}(W, C_1, X_1)} \frac{\mathbbm{1}_{x_1}(X_1)}{P_{x_1}(X_1|W, C_1)} \left\{ \mu(W, C_1, X_1) - \mu_0(W, C_1, X_1) \right\} \right] \\ &= \mathbb{E}_{P_{x_2}} \left[ \frac{P_{x_1}(W, C_1)\mathbbm{1}_{x_1}(X_1)}{P_{x_2}(W, C_1, X_1)} \left\{ \mu(W, C_1, X_1) - \mu_0(W, C_1, X_1) \right\} \right] \\ &= \mathbb{E}_{P_{x_2}} \left[ \frac{P_{x_1}(W|C_1)\mathbbm{1}_{x_1}(X_1)}{P_{x_2}(W, X_1|C_1)} \left\{ \mu(W, C_1, X_1) - \mu_0(W, C_1, X_1) \right\} \right] \\ &= \mathbb{E}_{P_{x_2}} \left[ \pi_0(W, C_1, X_1)\mathbbm{1}_{x_1}(X_1) \left\{ \mu(W, C_1, X_1) - \mu_0(W, C_1, X_1) \right\} \right], \end{split}$$

where  $\stackrel{7}{=}$  holds by Assumption 4. Then,

$$\begin{split} & \mathbb{E}_{P_{x_2}}\left[\pi(W,C_1,x_1)\mathbbm{1}_{x_1}(X_1)\{Y-\mu(W,C_1,X_1)\}\right] + \mathbb{E}_{P_{x_1}}\left[\mu(W,C_1,x_1)-\mu_0(W,C_1,x_1)\right] \\ & = \mathbb{E}_{P_{x_2}}\left[\pi(W,C_1,x_1)\mathbbm{1}_{x_1}(X_1)\{Y-\mu(W,C_1,X_1)\}\right] + \mathbb{E}_{P_{x_2}}\left[\pi_0(W,C_1,X_1)\mathbbm{1}_{x_1}(X_1)\left\{\mu(W,C_1,X_1)-\mu_0(W,C_1,X_1)\right\}\right] \\ & = \mathbb{E}_{P_{x_2}}\left[\mathbbm{1}_{x_1}(X_1)\left(\pi(W,C_1,x_1)\{\mu_0(W,C_1,X_1)-\mu(W,C_1,X_1)\} + \pi_0(W,C_1,X_1)\left\{\mu(W,C_1,X_1)-\mu_0(W,C_1,X_1)\right\}\right)\right] \\ & = \mathbb{E}_{P_{x_2}}\left[\mathbbm{1}_{x_1}(X_1)\left\{\mu_0(W,C_1,X_1)-\mu(W,C_1,X_1)\right\} + \pi_0(W,C_1,X_1)\left\{\mu(W,C_1,X_1)-\mu_0(W,C_1,X_1)\right\}\right] \\ & = \mathbb{E}_{P_{x_2}}\left[\mathbbm{1}_{x_1}(X_1)\left\{\mu_0(W,C_1,X_1)-\mu(W,C_1,X_1)\right\} + \pi_0(W,C_1,X_1)\right\} \\ & = O_{P_{x_2}}\left(\mathbbm{1}_{x_1}-\mu_0\mathbbm{1}_{x_1}\mathbbm{1}_{x_1}-\mu_0\mathbbm{1}_{x_1}$$

Therefore,

$$T^{dml} - \mathbb{E}\left[Y|do(x_1, x_2)\right] = R_1 + R_2 + O_{P_{x_2}}\left(\left\|\mu - \mu_0\right\| \|\pi - \pi_0\|\right).$$

**Corollary 2** (Doubly robustness of the DML estimators (Corollary of Thm. 2)). Suppose Assumptions (1,2,3,4) and AC-TTI in Def. 1 hold. Suppose either  $\pi = \pi_0$  or  $\mu = \mu_0$ . Then,  $T^{dml}$  is an unbiased estimator of  $\mathbb{E}[Y|do(x_1, x_2)]$ .

**Proof of Corollary 2.** Let  $\pi$  and  $\mu$  denote the limiting estimator for  $\pi_0$  and  $\mu_0$ .

$$T^{dml} \coloneqq \underbrace{\mathbb{E}_{D_{x_2}}\left[\pi(W, C_1, X_1) \mathbbm{1}_{x_1}(X_1) \{Y - \mu(W, C_1, X_1)\}\right]}_{:=T^{dml, 1}} + \underbrace{\mathbb{E}_{D_{x_1}}\left[\mu(W, C_1, x_1))\right]}_{:=T^{dml, 2}}$$

Let

$$T_0^{dml,1} \coloneqq \mathbb{E}_{P_{x_2}} \left[ \pi_0(W, C_1, X_1) \mathbb{1}_{x_1}(X_1) \{ Y - \mu_0(W, C_1, X_1) \} \right]$$
$$T_0^{dml,2} \coloneqq \mathbb{E}_{P_{x_1}} \left[ \mu_0(W, C_1, x_1) \right].$$

Under the assumption that

$$\begin{split} \mathbb{E}_{P_{x_2}} \left[ T^{dml,1} \right] &= \mathbb{E}_{P_{x_2}} \left[ \pi(W, C_1, X_1) \mathbb{1}_{x_1}(X_1) \left\{ Y - \mu(W, C_1, X_1) \right\} \right] \\ \mathbb{E}_{P_{x_1}} \left[ T^{dml,2} \right] &= \mathbb{E}_{P_{x_1}} \left[ \mu(W, C_1, x_1) \right] \\ &= \mathbb{E}_{P_{x_2}} \left[ \pi_0(W, C_1, X_1) \mathbb{1}_{x_1}(X_1) \mu(W, C_1, X_1) \right]. \end{split}$$

Then,

$$\begin{split} & \mathbb{E}_{P_{x_2}} \left[ T^{dml,1} \right] + \mathbb{E}_{P_{x_1}} \left[ T^{dml,2} \right] - \mathbb{E}_{P_{x_2}} \left[ T_0^{dml,1} \right] + \mathbb{E}_{P_{x_1}} \left[ T_0^{dml,2} \right] \\ &= \mathbb{E}_{P_{x_2}} \left[ \pi(W, C_1, X_1) \mathbb{1}_{x_1}(X_1) \{ Y - \mu(W, C_1, X_1) \} + \pi_0(W, C_1, X_1) \mathbb{1}_{x_1}(X_1) \{ \mu(W, C_1, X_1) - \mu_0(W, C_1, X_1) \} \right] \\ &= \mathbb{E}_{P_{x_2}} \left[ \mathbb{1}_{x_1}(X_1) \{ \mu_0(W, C_1, X_1) - \mu(W, C_1, X_1) \} \{ \pi_0(W, C_1, X_1) - \pi(W, C_1, X_1) \} \right] \\ &= O_{P_{x_2}} \left( \| \mu - \mu_0 \| \| \pi - \pi_0 \| \right) \\ &= 0, \end{split}$$

where the last equation holds under the given condition.

#### C.4. Proof of Theorem 3

**Definition 4** (Adjustment criterion for combining two experiments (AC-gTTI)). A set of variables  $\mathbf{A}$  is said to satisfy *adjustment criterion for generalized TTI (AC-gTTI)* w.r.t ( $\mathbf{X}$ , Y) in G if

- 1.  $Z_1 \subseteq \mathbf{X}$  and  $(\mathbf{A} \perp \mathbf{X} \setminus Z_1 | Z_1)_{G_{\overline{\mathbf{X}}}}$ , i.e., there are no direct paths from  $\mathbf{X} \setminus Z_1$  to  $\mathbf{A}$  in  $G_{\overline{\mathbf{X}}}$ ; and
- 2.  $Z_2 \subseteq \mathbf{X}$  and  $(Y \perp \mathbf{X} \setminus Z_2 | \mathbf{A}, Z_2)_{G_{\mathbf{X} \setminus Z_2 \overline{Z_2}}}$ , i.e., the back-door paths from  $\mathbf{X} \setminus Z_2$  to Y are blocked by  $\mathbf{A}$  in  $G_{\overline{Z_2}}$ .

Assumption 5 (Positivity Assumption for AC-gTTI).  $P_{z_1}(\mathbf{A}), P_{z_2}(\mathbf{A}), P_{z_2}(\mathbf{X} \setminus Z_2 | \mathbf{A})$  are strictly positive distributions  $\forall z_1, z_2 \in \mathfrak{D}_{Z_1, Z_2}$ .

**Theorem 3** (Identification through AC-gTTI). Suppose AC-gTTI in Def. 4 and Assumption 5 hold. Then, the query  $\mathbb{E}[Y|do(\mathbf{x})]$  is identifiable from  $P_{rand(Z_1)}(\mathbf{A})$  and  $P_{rand(Z_2)}(\mathbf{A}, \mathbf{X}, Y)$  and given as follows:

$$\mathbb{E}\left[Y|do(\mathbf{x})\right] = \mathbb{E}_{P_{z_1}}\left[\mathbb{E}_{P_{z_2}}\left[Y|\mathbf{A}, \mathbf{x} \setminus \mathbf{z}_2\right]\right].$$
(2)

**Proof of Theorem 3**.

$$\begin{split} \mathbb{E}\left[Y|do(\mathbf{x})\right] &= \mathbb{E}\left[\mathbb{E}\left[Y|do(\mathbf{x}), \mathbf{A}\right]|do(\mathbf{x})\right] \\ &= \mathbb{E}\left[\mathbb{E}\left[Y|do(z_2), \mathbf{x} \setminus z_2, \mathbf{A}\right]|do(\mathbf{x})\right] \\ &= \mathbb{E}\left[\mathbb{E}_{P_{z_2}}\left[Y|\mathbf{x} \setminus z_2, \mathbf{A}\right]|do(\mathbf{x})\right] \\ &\stackrel{2}{=} \mathbb{E}\left[\mathbb{E}_{P_{z_2}}\left[Y|\mathbf{x} \setminus z_2, \mathbf{A}\right]|do(z_1)\right] \\ &= \mathbb{E}_{P_{z_1}}\left[\mathbb{E}_{P_{z_2}}\left[Y|\mathbf{x} \setminus z_2, \mathbf{A}\right]\right], \end{split}$$

where  $\stackrel{1}{=}$  holds by the condition 2 which implies Rule 2 of *do*-calculus, and  $\stackrel{2}{=}$  holds by the condition 1 in AC-gTTI which implies Rule 3 of *do*-calculus.

### C.5. Proof of Theorem 4 and Corollary 4

**Definition 5** (Nuisances for gTTI). Nuisance functions for estimating AC-gTTI functional in Eq. (2) are defined as follows: For a fixed  $z_1, z_2 \in \mathfrak{D}_{Z_1, Z_2}$  where  $z_1, z_2$  are specified in Eq. (2),  $\pi_0 \coloneqq \pi_0(\mathbf{A}, \mathbf{X}) \coloneqq \frac{P_{z_1}(\mathbf{A})}{P_{z_2}(\mathbf{A}, \mathbf{X} \setminus Z_2)}$ , and  $\mu_0 \coloneqq \mu_0(\mathbf{A}, \mathbf{X}) \coloneqq \mathbb{E}_{P_{z_2}}[Y|\mathbf{X} \setminus Z_2, \mathbf{A}]$ . We will use  $\pi \coloneqq \pi(\mathbf{A}, \mathbf{X}) > 0$  and  $\mu \coloneqq \mu(\mathbf{A}, \mathbf{X})$  to denote estimated nuisances.

**Definition 6** (Estimators for gTTI). Let  $D_1, D_2$  denote two sample sets following distributions  $P_{\text{rand}(Z_1)}(\mathbf{A})$  and  $P_{\text{rand}(Z_2)}(\mathbf{A}, \mathbf{X}, Y)$ , respectively. For a fixed  $z_1, z_2 \in \mathfrak{D}_{Z_1, Z_2}$ , we define  $D_{z_1}$  and  $D_{z_2}$  as subsamples of  $D_1$  and  $D_2$  such that  $Z_1 = z_1$  and  $Z_2 = z_2$ . Let  $\mu, \pi$  denote nuisances defined in Def. 5. Then, {REG, PW, DML} estimators for the

AC-gTTI functional are defined as follows:

$$\begin{split} T^{reg} &\coloneqq \mathbb{E}_{D_{z_1}} \left[ \mu(\mathbf{A}, \mathbf{x}) \right], \\ T^{pw} &\coloneqq \mathbb{E}_{D_{z_2}} \left[ \pi(\mathbf{A}, \mathbf{X}) \mathbb{1}_{\mathbf{x}}(\mathbf{X}) Y \right], \\ T^{dml} &\coloneqq \mathbb{E}_{D_{z_2}} \left[ \pi \mathbb{1}_{\mathbf{x}}(\mathbf{X}) \{Y - \mu\} \right] + \mathbb{E}_{D_{z_1}} \left[ \mu(\mathbf{A}, \mathbf{x}) \right] . \end{split}$$

**Assumption 2** (Sample-splitting). Samples for training nuisances and evaluating the estimators equipped with the trained nuisance are separate and independent

Assumption 6 ( $L_2$  consistency of nuisances). Estimated nuisances are  $L_2$  consistent; i.e.,  $\forall i \in \{1,2\}, \forall z_i \in \mathfrak{D}_{Z_i}$ ,

$$\|\mu(\mathbf{A}, \mathbf{x}) - \mu_0(\mathbf{A}, \mathbf{x})\|_{P_{z_i}} = o_{P_{z_i}}(1),$$
  
$$\|\pi(\mathbf{A}, \mathbf{X}) - \pi_0(\mathbf{A}, \mathbf{X})\|_{P_{z_0}} = o_{P_{z_2}}(1).$$

**Theorem 4** (Error analysis of the estimators). Under Assumptions (2,5,6) and AC-gTTI in Def. 4, the errors of the estimators in Def. 6, denoted  $\epsilon^{est} := T^{est} - \mathbb{E}[Y|do(\mathbf{x})]$  for  $est \in \{reg, pw, dml\}$ , are:

$$\begin{split} \epsilon^{reg} &= R_1 + O_{P_{z_1}} \left( \| \mu - \mu_0 \| \right), \\ \epsilon^{pw} &= R_2 + O_{P_{z_2}} \left( \| \pi - \pi_0 \| \right), \\ \epsilon^{dml} &= R_1 + R_2 + O_{P_{z_2}} \left( \| \pi - \pi_0 \| \| \mu - \mu_0 \| \right), \end{split}$$

where  $R_i$  is a random variable such that  $\sqrt{n_i}R_i$  converges in distribution to a zero-mean normal random variable, where  $n_i := |D_{z_i}|$ .

*Proof of Theorem 4.* We provide error analyses for each estimators:

Analysis for  $T^{reg}$ .

We first note that

$$\mathbb{E}_{P_{z_1}}\left[\mu_0(\mathbf{A}, \mathbf{x})\right] = \mathbb{E}_{P_{z_1}}\left[\mathbb{E}_{P_{z_2}}\left[Y | \mathbf{A}, \mathbf{x} \setminus \mathbf{z}_2\right]\right] = \mathbb{E}\left[Y | do(\mathbf{x})\right],$$

where the last equation holds by Theorem 3. By Lemma C.3,

$$\begin{split} T^{reg} &- \mathbb{E}\left[Y|do(\mathbf{x})\right] \\ &= T^{reg} - \mathbb{E}_{P_{z_1}}\left[\mu_0(\mathbf{A}, \mathbf{x})\right] \\ &= \underbrace{\mathbb{E}_{P_{z_1} - D_1}\left[\mu_0(\mathbf{A}, \mathbf{x})\right] + \mathbb{E}_{P_{z_1} - D_1}\left[\mu(\mathbf{A}, \mathbf{x}) - \mu_0(\mathbf{A}, \mathbf{x})\right]}_{:= R_1} + \mathbb{E}_{P_{z_1}}\left[\mu(\mathbf{A}, \mathbf{x}) - \mu_0(\mathbf{A}, \mathbf{x})\right] \\ &\stackrel{1}{=} R_1 + \mathbb{E}_{P_{z_1}}\left[\mu(\mathbf{A}, \mathbf{x}) - \mu_0(\mathbf{A}, \mathbf{x})\right] \\ &\stackrel{2}{=} R_1 + O_{P_{z_1}}\left(\|\mu_0 - \mu\|\right), \end{split}$$

where

- $\stackrel{1}{=}$  holds by Lemma C.3.
- $\stackrel{2}{=}$  holds by Cauchy-Schwartz inequality.

Analysis for  $T^{pw}$ .

We first note that

$$\begin{split} \mathbb{E}_{P_{z_2}} \left[ \pi_0(\mathbf{A}, \mathbf{X}) \mathbb{1}_{\mathbf{x}}(\mathbf{X}) Y \right] &= \mathbb{E}_{P_{z_2}} \left[ \pi_0(\mathbf{A}, \mathbf{X}) \mathbb{1}_{\mathbf{x}}(\mathbf{X}) \mu_0(\mathbf{A}, \mathbf{X}) \right] \\ &= \mathbb{E}_{P_{z_2}} \left[ \frac{P_{z_1}(\mathbf{A})}{P_{z_2}(\mathbf{A}, \mathbf{X} \backslash Z_2)} \mathbb{1}_{\mathbf{x}}(\mathbf{X}) \mu_0(\mathbf{A}, \mathbf{X}) \right] \\ &= \mathbb{E}_{P_{z_1}} \left[ \mathbb{1}_{\mathbf{x}}(\mathbf{X}) \mu_0(\mathbf{A}, \mathbf{X}) \right] \\ &= \mathbb{E}_{P_{z_1}} \left[ \mu_0(\mathbf{A}, \mathbf{x}) \right] \\ &\stackrel{3}{=} \mathbb{E} \left[ Y | do(\mathbf{x}) \right], \end{split}$$

where  $\stackrel{3}{=}$  holds by Theorem 3. By applying Lemma C.3,

$$\begin{split} T^{pw} &- \mathbb{E} \left[ Y | do(\mathbf{x}) \right] \\ &= T^{pw} - \mathbb{E}_{P_{z_2}} \left[ \pi_0(\mathbf{A}, \mathbf{X}) \mathbb{1}_{\mathbf{X}}(\mathbf{X}) Y \right] \\ &= \underbrace{\mathbb{E}_{P_{z_2} - D_2} \left[ \pi_0(\mathbf{A}, \mathbf{X}) \mathbb{1}_{\mathbf{x}}(\mathbf{X}) Y \right] + \mathbb{E}_{P_{z_2} - D_2} \left[ \left\{ \pi_0(\mathbf{A}, \mathbf{X}) - \pi(\mathbf{A}, \mathbf{X}) \right\} \mathbb{1}_{\mathbf{x}}(\mathbf{X}) Y \right] \\ &= \underbrace{\mathbb{E}_{P_{z_2}} \left[ \left\{ \pi_0(\mathbf{A}, \mathbf{X}) - \pi(\mathbf{A}, \mathbf{X}) \right\} \mathbb{1}_{\mathbf{x}}(\mathbf{X}) Y \right] \\ &= R_2 \\ &+ \mathbb{E}_{P_{z_2}} \left[ \left\{ \pi_0(\mathbf{A}, \mathbf{X}) - \pi(\mathbf{A}, \mathbf{X}) \right\} \mathbb{1}_{\mathbf{x}}(\mathbf{X}) Y \right] \\ &= R_2 \\ &= R_2 \\ &= R_2 + \mathbb{E}_{P_{z_2}} \left[ \left\{ \pi_0(\mathbf{A}, \mathbf{X}) - \pi(\mathbf{A}, \mathbf{X}) \right\} \mathbb{1}_{\mathbf{x}}(\mathbf{X}) Y \right] \\ &= R_2 + \mathbb{E}_{P_{z_2}} \left[ \left\{ \pi_0(\mathbf{A}, \mathbf{X}) - \pi(\mathbf{A}, \mathbf{X}) \right\} \mathbb{1}_{\mathbf{x}}(\mathbf{X}) Y \right] \\ &= R_2 + O_{P_{z_2}} \left( \| \pi_0 - \pi \| \right), \end{split}$$

- $\stackrel{4}{=}$  holds by Lemmas (C.1, C.3).
- $\stackrel{5}{=}$  holds by Cauchy-Schwartz inequality and the setting where Y has a finite variance.

Analysis for  $T^{dml}$ .

Let

$$T^{dml} \coloneqq \underbrace{\mathbb{E}_{D_2}\left[\pi(\mathbf{A}, \mathbf{X}) \mathbb{1}_{\mathbf{x}}(\mathbf{X}) \{Y - \mu(\mathbf{A}, \mathbf{X})\}\right]}_{:=T^{dml, 1}} + \underbrace{\mathbb{E}_{D_1}\left[\mu(\mathbf{A}, \mathbf{x})\right)\right]}_{:=T^{dml, 2}}.$$

Let

$$T_0^{dml,1} \coloneqq \mathbb{E}_{P_{z_2}} \left[ \pi_0(\mathbf{A}, \mathbf{X}) \mathbb{1}_{\mathbf{x}}(\mathbf{X}) \{ Y - \mu_0(\mathbf{A}, \mathbf{X}) \} \right]$$
$$T_0^{dml,2} \coloneqq \mathbb{E}_{P_{z_1}} \left[ \mu_0(\mathbf{A}, \mathbf{x}) \right].$$

We note that  $T_0^{dml} \coloneqq T_0^{dml,1} + T_0^{dml,2} = \mathbb{E}[Y|do(\mathbf{x})]$ . We first apply the Lemma C.3 to  $T^{dml,1}$  and  $T^{dml,2}$  separately.

$$T^{dml,1} - T_0^{dml,1}$$

$$= \mathbb{E}_{P_{z_2} - D_2} \left[ \pi_0(\mathbf{A}, \mathbf{x}) \mathbb{1}_{\mathbf{x}}(\mathbf{X}) \{Y - \mu_0(\mathbf{A}, \mathbf{X})\} \right]$$

$$+ \mathbb{E}_{P_{z_2} - D_2} \left[ \pi_0(\mathbf{A}, \mathbf{x}) \mathbb{1}_{\mathbf{x}}(\mathbf{X}) \{Y - \mu_0(\mathbf{A}, \mathbf{X})\} - \pi(\mathbf{A}, \mathbf{x}) \mathbb{1}_{\mathbf{x}}(\mathbf{X}) \{Y - \mu(\mathbf{A}, \mathbf{X})\} \right]$$

$$+ \mathbb{E}_{P_{z_2}} \left[ \pi(\mathbf{A}, \mathbf{x}) \mathbb{1}_{\mathbf{x}}(\mathbf{X}) \{Y - \mu(\mathbf{A}, \mathbf{X})\} \right]$$

$$= R_2 + \mathbb{E}_{P_{z_2}} \left[ \pi(\mathbf{A}, \mathbf{x}) \mathbb{1}_{\mathbf{x}}(\mathbf{X}) \{Y - \mu(\mathbf{A}, \mathbf{X})\} \right],$$

where

$$\begin{aligned} R_2 &\coloneqq \mathbb{E}_{P_{z_2} - D_2} \left[ \pi_0(\mathbf{A}, \mathbf{x}) \mathbb{1}_{\mathbf{x}}(\mathbf{X}) \{ Y - \mu_0(\mathbf{A}, \mathbf{X}) \} \right] \\ &+ \mathbb{E}_{P_{z_2} - D_2} \left[ \pi_0(\mathbf{A}, \mathbf{x}) \mathbb{1}_{\mathbf{x}}(\mathbf{X}) \{ Y - \mu_0(\mathbf{A}, \mathbf{X}) \} - \pi(\mathbf{A}, \mathbf{x}) \mathbb{1}_{\mathbf{x}}(\mathbf{X}) \{ Y - \mu(\mathbf{A}, \mathbf{X}) \} \right]. \end{aligned}$$

Also, by the proof for analyzing the error of  $T^{reg}$ ,

$$T^{dml,2} - T_0^{dml,2} = R_1 + \mathbb{E}_{P_{z_1}} \left[ \mu(\mathbf{A}, \mathbf{x}) - \mu_0(\mathbf{A}, \mathbf{x}) \right].$$

Finally,

$$\begin{split} T^{dml} &- \mathbb{E}\left[Y|do(\mathbf{x})\right] \\ &= T^{dml,1} + T^{dml,2} - T_0^{dml,1} - T_0^{dml,2} \\ &= R_1 + R_2 + \mathbb{E}_{P_{z_2}}\left[\pi(\mathbf{A}, \mathbf{x})\mathbb{1}_{\mathbf{x}}(\mathbf{X})\{Y - \mu(\mathbf{A}, \mathbf{X})\}\right] + \mathbb{E}_{P_{z_1}}\left[\mu(\mathbf{A}, \mathbf{x}) - \mu_0(\mathbf{A}, \mathbf{x})\right]. \end{split}$$

We note that  $R_i$  for  $i \in \{1, 2\}$  is a variable such that  $\sqrt{n_i}R_i$  converges in distribution to the normal random variable, where  $n_i := |D_{z_i}|$ , by Lemmas (C.1, C.3). Note that

$$\begin{split} & \mathbb{E}_{P_{z_1}} \left[ \mu(\mathbf{A}, \mathbf{x}) - \mu_0(\mathbf{A}, \mathbf{x}) \right] \\ &= \mathbb{E}_{P_{z_1}} \left[ \frac{\mathbb{1}_{\mathbf{x}}(\mathbf{X})}{P_{z_1}(\mathbf{X}|\mathbf{A})} \left\{ \mu(\mathbf{A}, \mathbf{X}) - \mu_0(\mathbf{A}, \mathbf{X}) \right\} \right] \\ &= \mathbb{E}_{P_{z_2}} \left[ \frac{P_{z_1}(\mathbf{A}, \mathbf{X})}{P_{z_2}(\mathbf{A}, \mathbf{X})} \frac{\mathbb{1}_{\mathbf{x}}(\mathbf{X})}{P_{z_1}(\mathbf{X}|\mathbf{A})} \left\{ \mu(\mathbf{A}, \mathbf{X}) - \mu_0(\mathbf{A}, \mathbf{X}) \right\} \right] \\ &= \mathbb{E}_{P_{z_2}} \left[ \frac{P_{z_1}(\mathbf{A})\mathbb{1}_{\mathbf{x}}(\mathbf{X})}{P_{z_2}(\mathbf{A}, \mathbf{X} \setminus \mathbf{Z}_2)} \left\{ \mu(\mathbf{A}, \mathbf{X}) - \mu_0(\mathbf{A}, \mathbf{X}) \right\} \right] \\ &= \mathbb{E}_{P_{z_2}} \left[ \pi_0(\mathbf{A}, \mathbf{X})\mathbb{1}_{\mathbf{x}}(\mathbf{X}) \left\{ \mu(\mathbf{A}, \mathbf{X}) - \mu_0(\mathbf{A}, \mathbf{X}) \right\} \right]. \end{split}$$

Then,

$$\begin{split} & \mathbb{E}_{P_{z_2}}\left[\pi(\mathbf{A}, \mathbf{x})\mathbb{1}_{\mathbf{x}}(\mathbf{X})\{Y - \mu(\mathbf{A}, \mathbf{X})\}\right] + \mathbb{E}_{P_{z_1}}\left[\mu(\mathbf{A}, \mathbf{x}) - \mu_0(\mathbf{A}, \mathbf{x})\right] \\ &= \mathbb{E}_{P_{z_2}}\left[\pi(\mathbf{A}, \mathbf{x})\mathbb{1}_{\mathbf{x}}(\mathbf{X})\{Y - \mu(\mathbf{A}, \mathbf{X})\}\right] + \mathbb{E}_{P_{z_2}}\left[\pi_0(\mathbf{A}, \mathbf{X})\mathbb{1}_{\mathbf{x}}(\mathbf{X})\{\mu(\mathbf{A}, \mathbf{X}) - \mu_0(\mathbf{A}, \mathbf{X})\}\right] \\ &= \mathbb{E}_{P_{z_2}}\left[\mathbb{1}_{\mathbf{x}}(\mathbf{X})\left(\pi(\mathbf{A}, \mathbf{x})\{\mu_0(\mathbf{A}, \mathbf{X}) - \mu(\mathbf{A}, \mathbf{X})\} + \pi_0(\mathbf{A}, \mathbf{X})\{\mu(\mathbf{A}, \mathbf{X}) - \mu_0(\mathbf{A}, \mathbf{X})\}\right)\right] \\ &= \mathbb{E}_{P_{z_2}}\left[\mathbb{1}_{\mathbf{x}}(\mathbf{X})\{\mu_0(\mathbf{A}, \mathbf{X}) - \mu(\mathbf{A}, \mathbf{X})\} + \pi_0(\mathbf{A}, \mathbf{X})\{\mu(\mathbf{A}, \mathbf{X}) - \mu_0(\mathbf{A}, \mathbf{X})\}\right)\right] \\ &= \mathbb{E}_{P_{z_2}}\left[\mathbb{1}_{\mathbf{x}}(\mathbf{X})\{\mu_0(\mathbf{A}, \mathbf{X}) - \mu(\mathbf{A}, \mathbf{X})\}\{\pi_0(\mathbf{A}, \mathbf{X}) - \pi(\mathbf{A}, \mathbf{X})\}\right] \\ &= O_{P_{z_0}}\left(\|\mu - \mu_0\|\|\pi - \pi_0\|\right). \end{split}$$

Therefore,

$$T^{dml} - \mathbb{E}\left[Y|do(\mathbf{x})\right] = R_1 + R_2 + O_{P_{z_2}}\left(\|\mu - \mu_0\| \|\pi - \pi_0\|\right).$$

**Corollary 4** (Doubly robustness of the DML estimators (Corollary of Thm. 4)). Suppose Assumptions (2,5,6) and *AC-gTTI in Def. 4 hold. Suppose either*  $\pi = \pi_0$  or  $\mu = \mu_0$ . Then,  $T^{dml}$  is an unbiased estimator of  $\mathbb{E}[Y|do(\mathbf{x})]$ .

**Proof of Corollary 4.** Let  $\pi$  and  $\mu$  denote the limiting estimator for  $\pi_0$  and  $\mu_0$ .

$$T^{dml} \coloneqq \underbrace{\mathbb{E}_{D_2}\left[\pi(\mathbf{A}, \mathbf{X}) \mathbb{1}_{\mathbf{x}}(\mathbf{X}) \{Y - \mu(\mathbf{A}, \mathbf{X})\}\right]}_{:=T^{dml, 1}} + \underbrace{\mathbb{E}_{D_1}\left[\mu(\mathbf{A}, \mathbf{x})\right)}_{:=T^{dml, 2}}$$

Let

$$T_0^{dml,1} \coloneqq \mathbb{E}_{P_{z_2}} \left[ \pi_0(\mathbf{A}, \mathbf{X}) \mathbb{1}_{\mathbf{x}}(\mathbf{X}) \{ Y - \mu_0(\mathbf{A}, \mathbf{X}) \} \right]$$
$$T_0^{dml,2} \coloneqq \mathbb{E}_{P_{z_1}} \left[ \mu_0(\mathbf{A}, \mathbf{x}) \right].$$
Under the assumption that samples are i.i.d.,

$$\begin{split} \mathbb{E}_{P_{z_2}} \left[ T^{dml,1} \right] &= \mathbb{E}_{P_{z_2}} \left[ \pi(\mathbf{A}, \mathbf{X}) \mathbb{1}_{\mathbf{x}}(\mathbf{X}) \left\{ Y - \mu(\mathbf{A}, \mathbf{X}) \right\} \right] \\ \mathbb{E}_{P_{z_1}} \left[ T^{dml,2} \right] &= \mathbb{E}_{P_{z_1}} \left[ \mu(\mathbf{A}, \mathbf{x}) \right] \\ &= \mathbb{E}_{P_{z_0}} \left[ \pi_0(\mathbf{A}, \mathbf{X}) \mathbb{1}_{\mathbf{x}}(\mathbf{X}) \mu(\mathbf{A}, \mathbf{X}) \right]. \end{split}$$

Then,

$$\begin{split} & \mathbb{E}_{P_{z_2}}\left[T^{dml,1}\right] + \mathbb{E}_{P_{z_1}}\left[T^{dml,2}\right] - \mathbb{E}_{P_{z_2}}\left[T_0^{dml,1}\right] + \mathbb{E}_{P_{z_1}}\left[T_0^{dml,2}\right] \\ & = \mathbb{E}_{P_{z_2}}\left[\pi(\mathbf{A}, \mathbf{X})\mathbb{1}_{\mathbf{x}}(\mathbf{X})\{Y - \mu(\mathbf{A}, \mathbf{X})\} + \pi_0(\mathbf{A}, \mathbf{X})\mathbb{1}_{\mathbf{x}}(\mathbf{X})\{\mu(\mathbf{A}, \mathbf{X}) - \mu_0(\mathbf{A}, \mathbf{X})\}\} \\ & = \mathbb{E}_{P_{z_2}}\left[\mathbb{1}_{\mathbf{x}}(\mathbf{X})\{\mu_0(\mathbf{A}, \mathbf{X}) - \mu(\mathbf{A}, \mathbf{X})\}\{\pi_0(\mathbf{A}, \mathbf{X}) - \pi(\mathbf{A}, \mathbf{X})\}\right] \\ & = 0, \end{split}$$

where the last equation holds under the given condition.

### C.6. Proof of Theorem 5

**Definition 7** (Adjustment criterion for Multiple Treatment Interaction (AC-MTI)). An ordered set  $\{C_1, W_1, C_2, W_2, \dots, C_{m-1}, W_{m-1}\}$  satisfies adjustment criterion for multiple treatment interaction (AC-MTI) w.r.t.  $\{\mathbf{X}, Y\}$  for  $\mathbf{X} = \{X_i\}_{i=1}^m$  in G if, for  $i = 1, 2, \dots, m$ ,

1.  $\{X_j\}_{j>i}$  is non-ancestor of  $\{\mathbf{X}^{(i)}, \mathbf{W}^{(i)}, \mathbf{C}^{(i)}\}$ ; and

2.  $(Y \perp X_i | \mathbf{C}^{(i-1)}, \mathbf{X}^{(i-1)}, \mathbf{W}^{(i)}, \mathbf{X}^{>i})_{G_{\underline{X_i}, \overline{\mathbf{X}^{>i}}}}$ , i.e., the back-door paths from  $X_i$  to Y are blocked by  $\mathbf{C}^{(i-1)}, \mathbf{X}^{(i-1)}, \mathbf{W}^{(i)}, \mathbf{X}^{>i}$  in the graph  $G_{\overline{\mathbf{X}^{>i}}}$ .

Assumption 7 (Positivity Assumption for AC-MTI).  $\{P_{x_i}(W_i, C_i | \mathbf{W}^{(i-1)}, \mathbf{C}^{(i-1)}, \mathbf{X}^{(i-1)})\}_{i=1}^m$ ,  $P_{x_{i+1}}(X_i | \mathbf{W}^{(i)}, \mathbf{C}^{(i)}, \mathbf{X}^{(i-1)})$  for  $i = 1, \dots, m-1$  are strictly positive  $\forall \mathbf{x} \in \mathfrak{D}_{\mathbf{X}}$ .

**Theorem 5** (Identification through AC-MTI). Suppose AC-MTI in Def. 7 and Assumption 7 hold. Then,  $\mathbb{E}[Y(\mathbf{x})]$  is identifiable from  $\{P_{rand(X_i)}(\mathbf{C}^{(i)}, \mathbf{W}^{(i)}, \mathbf{X}^{(i-1)})\}_{i=1}^m$  as follows: Let  $\mu_0^m \coloneqq \mathbb{E}_{P_{x_m}}[Y|\mathbf{W}^{(m-1)}, \mathbf{C}^{(m-1)}, \mathbf{X}^{(m-1)}]$ , and for  $i = m - 1, \dots, 2$ ,

$$\mu_0^i \coloneqq \mathbb{E}_{P_{x_i}}\left[\overline{\mu}_0^{i+1} | \mathbf{W}^{(i-1)}, \mathbf{C}^{(i-1)}, \mathbf{X}^{(i-1)}\right],$$

where  $\overline{\mu}_{0}^{i+1} \coloneqq \mu_{0}^{i+1}(\mathbf{W}^{(i)}, \mathbf{C}^{(i)}, x_{i}, \mathbf{X}^{(i-1)})$ . Then,

$$\mathbb{E}[Y(\mathbf{x})] = \mathbb{E}_{P_{x_1}}\left[\mu_0^2(W_1, C_1, x_1)\right].$$
(3)

**Proof of Theorem 5.** Let  $A_i := \{W_i, C_i\}$  in this proof. Then, it suffices to show the following equation: For all  $i = 1, 2, \dots, m-1$ ,

$$\mathbb{E}\left[Y\middle|do(\mathbf{x}^{\geq i}), \mathbf{A}^{(i-1)}, \mathbf{x}^{(i-1)}\right] = \mathbb{E}\left[\mathbb{E}\left[Y\middle|do(\mathbf{x}^{\geq i+1}), \mathbf{A}^{(i)}, \mathbf{x}^{(i)}\right]\middle|do(x_i), \mathbf{A}^{(i-1)}, \mathbf{x}^{(i-1)}\right].$$

It holds as follow:

$$\begin{split} & \mathbb{E}\left[Y \middle| do(\mathbf{x}^{\geq i}), \mathbf{A}^{(i-1)}, \mathbf{x}^{(i-1)}\right] \\ & \stackrel{1}{=} \mathbb{E}\left[\mathbb{E}\left[Y \middle| do(\mathbf{x}^{\geq i}), \mathbf{A}^{(i)}, \mathbf{x}^{(i-1)}\right] \middle| do(\mathbf{x}^{\geq i}), \mathbf{A}^{(i-1)}, \mathbf{x}^{(i-1)}\right] \\ & \stackrel{2}{=} \mathbb{E}\left[\mathbb{E}\left[Y \middle| do(\mathbf{x}^{\geq i+1}), \mathbf{A}^{(i)}, \mathbf{x}^{(i)}\right] \middle| do(\mathbf{x}^{\geq i}), \mathbf{A}^{(i-1)}, \mathbf{x}^{(i-1)}\right] \\ & \stackrel{3}{=} \mathbb{E}\left[\mathbb{E}\left[Y \middle| do(\mathbf{x}^{\geq i+1}), \mathbf{A}^{(i)}, \mathbf{x}^{(i)}\right] \middle| do(x_i), \mathbf{A}^{(i-1)}, \mathbf{x}^{(i-1)}\right], \end{split}$$

where

- $\stackrel{1}{=}$  holds by marginalizing over  $\mathbf{A}_i$ .
- $\stackrel{2}{=}$  holds as follow:

$$\mathbb{E}\left[Y\middle|do(\mathbf{x}^{\geq i}), \mathbf{A}^{(i)}, \mathbf{x}^{(i-1)}\right] = \mathbb{E}\left[Y\middle|do(\mathbf{x}^{\geq i+1}), \mathbf{A}^{(i)}, \mathbf{x}^{(i)}\right],$$

since  $(Y \perp X_i | \mathbf{A}^{(i)}, \mathbf{X}^{(i-1)}, \mathbf{X}^{\geq i+1})_{G_{\underline{X_i}, \overline{\mathbf{X}^{\geq i+1}}}}$  by the given condition and the positivity condition.

•  $\stackrel{3}{=}$  holds because  $\mathbf{X}^{\geq i+1}$  is not an ancestor of  $\mathbf{A}^{(i)}, \mathbf{X}^{(i)}$ .

### C.7. Proof of Theorem 6 and Corollary 6

**Definition 8** (Nuisances for MTI). Nuisance functions for AC-MTI are defined as follows: For a fixed  $\mathbf{x} := \{x_1, \dots, x_m\} \in \mathfrak{D}_{\mathbf{X}}, \text{ let } \{\mu^i\}_{i=2}^m \text{ and } \{\overline{\mu}^i\}_{i=2}^m \text{ be the nuisances defined in Thm. 5. For } i = 1, \dots, m-1, \pi_0^i := \frac{P_{x_i}(W_i|C_i, \mathbf{C}^{(i-1)}, \mathbf{X}^{(i-1)}, \mathbf{W}^{(i-1)})}{P_{x_m}(W_i, X_i|C_i, \mathbf{C}^{(i-1)}, \mathbf{X}^{(i-1)}, \mathbf{W}^{(i-1)})}, \text{ and } \pi_0^{(i)} := \prod_{j=1}^i \pi_0^j(\mathbf{W}^{(j)}, \mathbf{C}^{(j)}, \mathbf{X}^{(j)}). \text{ We will use } \pi^i(\mathbf{W}^{(i)}, \mathbf{C}^{(i)}, \mathbf{X}^{(i)}) > 0$  and  $\mu^i(\mathbf{W}^{(i-1)}, \mathbf{C}^{(i-1)}, \mathbf{X}^{(i-1)})$  to denote estimated nuisances.

**Definition 9** (AC-MTI estimators). Let  $D_i$  denote samples following  $P_{\operatorname{rand}(X_i)}(\mathbf{C}^{(i)}, \mathbf{W}^{(i)}, \mathbf{X}^{(i)})$  for  $i = 1, 2, \dots, m$ . For a fixed  $x_i \in \mathfrak{D}_{X_i}$ , let  $D_{x_i}$  denote the subsamples of  $D_i$  such that  $X_i = x_i$ . Let  $A_i \coloneqq \{W_i, C_i\}$  and  $V_i \coloneqq \{A_i, X_i\}$ . Let  $\mu^{m+1} \coloneqq Y$ . Let  $\mathbb{1}_{\mathbf{x}}^{i-1} \coloneqq \mathbb{1}_{\mathbf{x}^{(i-1)}}(\mathbf{X}^{(i-1)})$  for  $i = 2, \dots, m$ . Then {REG, PW, DML} estimators are defined as follows:

$$\begin{split} T^{reg} &\coloneqq \mathbb{E}_{D_{x_1}} \left[ \mu^2(W_1, C_1, x_1)) \right], \\ T^{pw} &\coloneqq \mathbb{E}_{D_{x_m}} \left[ \pi^{(m-1)} \mathbb{1}_{\mathbf{x}}(\mathbf{X}) Y \right], \\ T^{dml} &\coloneqq \sum_{i=2}^m \mathbb{E}_{D_{x_i}} \left[ \pi^{(i-1)} \mathbb{1}_{\mathbf{x}}^{i-1} \{ \overline{\mu}^{i+1} - \mu^i \} \right] + \mathbb{E}_{D_{x_1}} \left[ \overline{\mu}^2 \right] \end{split}$$

**Assumption 2** (Sample-splitting). Samples for training nuisances and evaluating the estimators equipped with the trained nuisance are separate and independent

**Assumption 8** ( $L_2$  consistency of nuisances). Estimated nuisances are  $L_2$ -consistent; specifically,

$$\begin{aligned} \|\mu^{i+1} - \mu_0^{i+1}\|_{P_{x_i}} &= o_{P_{x_i}}(1), \ \forall i \in \{1, 2, \cdots, m-1\} \\ \|\mu^i - \mu_0^i\|_{P_{x_i}} &= o_{P_{x_i}}(1), \ \forall i \in \{2, \cdots, m\} \\ \|\pi^i - \pi^i\|_{P_{x_{i+1}}} &= o_{P_{x_{i+1}}}(1), \ \forall i \in \{1, \cdots, m-1\}. \end{aligned}$$

Assumption 9 (Multiple experiments represent the same population). For any fixed  $i, j \in \{1, 2, \dots, m-1\}$ s.t. j > i and any fixed  $x_i, x_j \in \mathfrak{D}_{X_i, X_j}$ , the baseline covariates  $C_i$ 's distribution satisfies the following:  $P_{x_i}(C_i|\mathbf{C}^{(i-1)}, \mathbf{X}^{(j-1)}, \mathbf{W}^{(j-1)}) = P_{x_j}(C_i|\mathbf{C}^{(i-1)}, \mathbf{X}^{(i-1)}, \mathbf{W}^{(i-1)}).$ 

**Lemma C.4** (Error analysis of the REG estimator for MTI). Suppose Assumptions (2,8) hold. Let  $T^{reg}$  denote the estimator defined in Def. 9. Then,

$$T^{reg} - \mathbb{E}\left[Y|do(\mathbf{x})\right] = R_1 + O_{P_{x_1}}(\|\mu^2 - \mu_0^2\|),$$

where  $R_1$  is the random variable such that  $\sqrt{n_1}R_1$  converges in distribution to the mean-zero normal random variable, where  $n_1 \coloneqq |D_{x_1}|$ .

*Proof of Lemma C.4.* We first note that, by Theorem 5,

$$\mathbb{E}_{P_{x_1}}\left[\mu_0^2(W_1, C_1, x_1)\right] = \mathbb{E}\left[Y|do(\mathbf{x})\right].$$

By Lemma C.3,

$$\begin{split} T^{reg} &- \mathbb{E}\left[Y|do(\mathbf{x})\right] \\ &= T^{reg} - \mathbb{E}_{P_{x_1}}\left[\mu_0^2(W_1, C_1, x_1)\right] \\ &\stackrel{1}{=} \underbrace{\mathbb{E}_{P_{x_1} - D_{x_1}}\left[\mu_0^2(W_1, C_1, x_1)\right] + \mathbb{E}_{P_{x_1} - D_{x_1}}\left[\mu_0^2(W_1, C_1, x_1) - \mu^2(W_1, C_1, x_1)\right]}_{:= R_1} \\ &+ \mathbb{E}_{P_{x_1}}\left[\mu_0^2(W_1, C_1, x_1) - \mu^2(W_1, C_1, x_1)\right] \\ &= R_1 + \mathbb{E}_{P_{x_1}}\left[\mu_0^2(W_1, C_1, x_1) - \mu^2(W_1, C_1, x_1)\right] \\ &= R_1 + O_{P_{x_1}}(\|\mu^2 - \mu_0^2\|). \end{split}$$

 $\stackrel{1}{=}$  holds by Lemmas (C.1, C.3), and the last equation holds by Cauchy-Schwartz inequality.

**Lemma C.5** (Error analysis of the PW estimator for MTI). Suppose Assumptions (2,8,9) hold. Let  $T^{pw}$  denote the estimator defined in Def. 9. Then,

$$T^{pw} - \mathbb{E}\left[Y|do(\mathbf{x})\right] = R_m + O_{P_{x_m}}(\|\pi^{(m-1)} - \pi_0^{(m-1)}\|),$$

where  $R_m$  is the random variable such that  $\sqrt{n_m}R_m$  converges in distribution to the mean-zero normal random variable, where  $n_m := |D_{x_m}|$ .

**Proof of Lemma C.5.** Throughout the proof, we set  $A_i := \{C_i, W_i\}$  for all  $i = 1, 2, \dots, m$ . We will use  $V_i := \{A_i, X_i\}$ . In the proof, we tentatively assume

$$\mathbb{E}_{P_{x_m}}\left[\pi_0^{(m-1)}(\mathbf{W}^{(m-1)}, \mathbf{C}^{(m-1)}, \mathbf{X}^{(m-1)})\mathbb{1}_{\mathbf{x}}(\mathbf{X})Y\right] = \mathbb{E}\left[Y|do(\mathbf{x})\right].$$
(C.3)

Then, by Lemma C.3,

$$T^{pw} - \mathbb{E}\left[Y|do(\mathbf{x})\right] = T^{pw} - \mathbb{E}_{P_{x_m}}\left[\pi_0^{(m-1)}(\mathbf{W}^{(m-1)}, \mathbf{C}^{(m-1)}, \mathbf{X}^{(m-1)})\mathbb{1}_{\mathbf{x}}(\mathbf{X})Y\right] \\ \stackrel{1}{=} \underbrace{\mathbb{E}_{P_{x_m} - D_{x_m}}\left[\pi_0^{(m-1)}(\mathbf{V}^{(m-1)})\mathbb{1}_{\mathbf{x}}(\mathbf{X})Y\right] + \mathbb{E}_{P_{x_m} - D_{x_m}}\left[\left\{\pi_0^{(m-1)}(\mathbf{V}^{(m-1)}) - \pi^{(m-1)}(\mathbf{V}^{(m-1)})\right\}\mathbb{1}_{\mathbf{x}}(\mathbf{X})Y\right] \\ \stackrel{:= R_m}{:= R_m} \\ + \mathbb{E}_{P_{x_m}}\left[\left\{\pi_0^{(m-1)}(\mathbf{V}^{(m-1)}) - \pi^{(m-1)}(\mathbf{V}^{(m-1)})\right\}\mathbb{1}_{\mathbf{x}}(\mathbf{X})Y\right] \\ = R_m + \mathbb{E}_{P_{x_m}}\left[\left\{\pi_0^{(m-1)}(\mathbf{V}^{(m-1)}) - \pi^{(m-1)}(\mathbf{V}^{(m-1)})\right\}\mathbb{1}_{\mathbf{x}}(\mathbf{X})Y\right] \\ = R_m + O_{P_{x_m}}(\|\pi^{(m-1)} - \pi_0^{(m-1)}\|),$$

where  $\stackrel{1}{=}$  holds by Lemmas (C.1, C.3). The last equation holds by Cauchy-Schwartz inequality.

We now prove Eq. (C.3). We first show the following: For  $i = 2, \dots, m$ ,

$$\mathbb{E}\left[Y|do(\mathbf{x})\right] = \mathbb{E}_{P_{x_i}}\left[\prod_{j=1}^{i-1} \frac{P_{x_j}(\mathbf{A}^{(j)}, \mathbf{X}^{(j-1)})}{P_{x_{j+1}}(\mathbf{A}^{(j)}, \mathbf{X}^{(j)})} \mu_0^i(\mathbf{A}^{(i-1)}, \mathbf{X}^{(i-1)}) \mathbb{1}_{\mathbf{x}^{(i-1)}}(\mathbf{X}^{(i-1)})\right].$$
(C.4)

It holds for i = 2 as follow:

$$\mathbb{E}_{P_{x_2}}\left[\frac{P_{x_1}(A_1)}{P_{x_2}(A_1, X_1)}\mu_0^2(A_1, X_1)\mathbb{1}_{x_1}(X_1)\right] = \mathbb{E}_{P_{x_1}}\left[\mu_0^2(A_1, x_1)\right] = \mathbb{E}\left[Y|do(\mathbf{x})\right]$$

where the last equation holds by Lemma C.4. Now, we make the following induction hypothesis: For some  $i - 1 \in \{2, 3, \dots, m - 1\}$ , suppose

$$\mathbb{E}[Y|do(\mathbf{x})] \stackrel{\text{induction hypothesis}}{=} \mathbb{E}_{P_{x_{i-1}}} \left[ \prod_{j=1}^{i-2} \frac{P_{x_j}(\mathbf{A}^{(j)}, \mathbf{X}^{(j-1)})}{P_{x_{j+1}}(\mathbf{A}^{(j)}, \mathbf{X}^{(j)})} \mu_0^{i-1}(\mathbf{A}^{(i-2)}, \mathbf{X}^{(i-2)}) \mathbb{1}_{\mathbf{x}^{(i-2)}}(\mathbf{X}^{(i-2)}) \right].$$

Then,

$$\begin{split} & \mathbb{E}_{P_{x_{i-1}}} \left[ \prod_{j=1}^{i-2} \frac{P_{x_j}(\mathbf{A}^{(j)}, \mathbf{X}^{(j-1)})}{P_{x_{j+1}}(\mathbf{A}^{(j)}, \mathbf{X}^{(j)})} \mu_0^{i-1}(\mathbf{A}^{(i-2)}, \mathbf{X}^{(i-2)}) \mathbb{1}_{\mathbf{x}^{(i-2)}}(\mathbf{X}^{(i-2)}) \right] \\ &= \mathbb{E}_{P_{x_{i-1}}} \left[ \prod_{j=1}^{i-2} \frac{P_{x_j}(\mathbf{A}^{(j)}, \mathbf{X}^{(j-1)})}{P_{x_{j+1}}(\mathbf{A}^{(j)}, \mathbf{X}^{(j)})} \mathbb{E}_{P_{x_{i-1}}} \left[ \mu_0^i(\mathbf{A}^{(i-1)}, x_{i-1}, \mathbf{X}^{(i-2)}) | \mathbf{A}^{(i-2)}, \mathbf{X}^{(i-2)} \right] \mathbb{1}_{\mathbf{x}^{(i-2)}}(\mathbf{X}^{(i-2)}) \right] \\ &= \mathbb{E}_{P_{x_{i-1}}} \left[ \prod_{j=1}^{i-2} \frac{P_{x_j}(\mathbf{A}^{(j)}, \mathbf{X}^{(j-1)})}{P_{x_{j+1}}(\mathbf{A}^{(j)}, \mathbf{X}^{(j)})} \mu_0^i(\mathbf{A}^{(i-1)}, x_{i-1}, \mathbf{X}^{(i-2)}) \mathbb{1}_{\mathbf{x}^{(i-2)}}(\mathbf{X}^{(i-2)}) \right] \\ &= \mathbb{E}_{P_{x_{i-1}}} \left[ \prod_{j=1}^{i-2} \frac{P_{x_j}(\mathbf{A}^{(j)}, \mathbf{X}^{(j-1)})}{P_{x_{j+1}}(\mathbf{A}^{(j)}, \mathbf{X}^{(j)})} \mu_0^i(\mathbf{A}^{(i-1)}, \mathbf{X}^{(i-1)}) \mathbb{1}_{\mathbf{x}^{(i-1)}}(\mathbf{X}^{(i-1)}) \right] \\ &= \mathbb{E}_{P_{x_i}} \left[ \prod_{j=1}^{i-2} \frac{P_{x_j}(\mathbf{A}^{(j)}, \mathbf{X}^{(j-1)})}{P_{x_{j+1}}(\mathbf{A}^{(j)}, \mathbf{X}^{(j)})} \mu_0^i(\mathbf{A}^{(i-1)}, \mathbf{X}^{(i-1)}) \right] \\ &= \mathbb{E}_{P_{x_i}} \left[ \prod_{j=1}^{i-2} \frac{P_{x_j}(\mathbf{A}^{(j)}, \mathbf{X}^{(j-1)})}{P_{x_{j+1}}(\mathbf{A}^{(j)}, \mathbf{X}^{(j)})} \mu_0^i(\mathbf{A}^{(i-1)}, \mathbf{X}^{(i-1)}) \right] \\ &= \mathbb{E}_{P_{x_i}} \left[ \prod_{j=1}^{i-1} \frac{P_{x_j}(\mathbf{A}^{(j)}, \mathbf{X}^{(j-1)})}{P_{x_{j+1}}(\mathbf{A}^{(j)}, \mathbf{X}^{(j)})} \mu_0^i(\mathbf{A}^{(i-1)}, \mathbf{X}^{(i-1)}) \right] \\ &= \mathbb{E}_{P_{x_i}} \left[ \prod_{j=1}^{i-1} \frac{P_{x_j}(\mathbf{A}^{(j)}, \mathbf{X}^{(j-1)})}{P_{x_{j+1}}(\mathbf{A}^{(j)}, \mathbf{X}^{(j)})} \mu_0^i(\mathbf{A}^{(i-1)}, \mathbf{X}^{(i-1)}) \right] \\ &= \mathbb{E}_{P_{x_i}} \left[ \prod_{j=1}^{i-1} \frac{P_{x_j}(\mathbf{A}^{(j)}, \mathbf{X}^{(j-1)})}{P_{x_{j+1}}(\mathbf{A}^{(j)}, \mathbf{X}^{(j)})} \mu_0^i(\mathbf{A}^{(i-1)}, \mathbf{X}^{(i-1)}) \right] \\ &= \mathbb{E}_{i-1} \left[ \prod_{j=1}^{i-1} \frac{P_{x_j}(\mathbf{A}^{(j)}, \mathbf{X}^{(j-1)})}{P_{x_{j+1}}(\mathbf{A}^{(j)}, \mathbf{X}^{(j)})} \mu_0^i(\mathbf{A}^{(i-1)}, \mathbf{X}^{(i-1)}) \right] \\ &= \mathbb{E}_{i-1} \left[ \prod_{j=1}^{i-1} \frac{P_{x_j}(\mathbf{A}^{(j)}, \mathbf{X}^{(j-1)})}{P_{x_{j+1}}(\mathbf{A}^{(j)}, \mathbf{X}^{(j)})} \right] \\ \end{bmatrix}$$

where

•  $\stackrel{1}{=}$  holds by the law of total expectation.

•  $\stackrel{2}{=}$  holds since the expectation is over  $P_{x_{i-1}}$ .

Therefore, Eq. (C.4) holds. By plugging i = m, we have

$$\begin{split} \mathbb{E}\left[Y|do(\mathbf{x})\right] &= \mathbb{E}_{P_{x_m}} \left[\prod_{j=1}^{m-1} \frac{P_{x_j}(\mathbf{A}^{(j)}, \mathbf{X}^{(j-1)})}{P_{x_{j+1}}(\mathbf{A}^{(j)}, \mathbf{X}^{(j)})} \mu_0^m(\mathbf{A}^{(m-1)}, \mathbf{X}^{(m-1)}) \mathbb{1}_{\mathbf{x}^{(m-1)}}(\mathbf{X}^{(m-1)})\right] \\ &= \mathbb{E}_{P_{x_m}} \left[\prod_{j=1}^{m-1} \frac{P_{x_j}(\mathbf{A}^{(j)}, \mathbf{X}^{(j-1)})}{P_{x_{j+1}}(\mathbf{A}^{(j)}, \mathbf{X}^{(j)})} \mathbb{E}_{P_{x_m}}\left[Y|\mathbf{A}^{(m-1)}, \mathbf{X}^{(m-1)}\right] \mathbb{1}_{\mathbf{x}^{(m-1)}}(\mathbf{X}^{(m-1)})\right] \\ &= \mathbb{E}_{P_{x_m}} \left[\prod_{j=1}^{m-1} \frac{P_{x_j}(\mathbf{A}^{(j)}, \mathbf{X}^{(j-1)})}{P_{x_{j+1}}(\mathbf{A}^{(j)}, \mathbf{X}^{(j)})} \mathbb{1}_{\mathbf{x}^{(m-1)}}(\mathbf{X}^{(m-1)})Y\right]. \end{split}$$

Finally,

$$\begin{split} \prod_{j=1}^{m-1} \frac{P_{x_j}(\mathbf{A}^{(j)}, \mathbf{X}^{(j-1)})}{P_{x_{j+1}}(\mathbf{A}^{(j)}, \mathbf{X}^{(j)})} &= \frac{P_{x_1}(\mathbf{A}^{(1)}) P_{x_2}(\mathbf{A}^{(2)}, \mathbf{X}^{(1)}) \cdots P_{x_{m-1}}(\mathbf{A}^{(m-1)}, \mathbf{X}^{(m-2)})}{P_{x_2}(\mathbf{A}^{(1)}, \mathbf{X}^{(1)}) P_{x_3}(\mathbf{A}^{(2)}, \mathbf{X}^{(2)}) \cdots P_{x_{m-1}}(\mathbf{A}^{(m-2)}, \mathbf{X}^{(m-2)}) P_{x_m}(\mathbf{A}^{(m-1)}, \mathbf{X}^{(m-1)})}{\mathbf{P}_{x_m}(\mathbf{A}^{(m-1)}, \mathbf{X}^{(m-1)})} \\ &= \frac{1}{P_{x_m}(\mathbf{A}^{(m-1)}, \mathbf{X}^{(m-1)})} \prod_{j=1}^{m-1} P_{x_j}(A_j | \mathbf{A}^{(j-1)}, \mathbf{X}^{(j-1)}) \\ &= \prod_{j=1}^{m-1} \frac{P_{x_j}(A_j | \mathbf{A}^{(j-1)}, \mathbf{X}^{(j-1)})}{P_{x_m}(C_j, W_j | \mathbf{A}^{(j-1)}, \mathbf{X}^{(j-1)})} \\ &= \prod_{j=1}^{m-1} \frac{P_{x_j}(W_j | C_j, \mathbf{A}^{(j-1)}, \mathbf{X}^{(j-1)})}{P_{x_m}(W_j, X_j | C_j, \mathbf{A}^{(j-1)}, \mathbf{X}^{(j-1)}) P_{x_m}(C_j | \mathbf{A}^{(j-1)}, \mathbf{X}^{(j-1)})} \\ &= \prod_{j=1}^{m-1} \frac{P_{x_j}(W_j | C_j, \mathbf{A}^{(j-1)}, \mathbf{X}^{(j-1)}) P_{x_m}(C_j | \mathbf{A}^{(j-1)}, \mathbf{X}^{(j-1)})}{P_{x_m}(W_j, X_j | C_j, \mathbf{A}^{(j-1)}, \mathbf{X}^{(j-1)})} \\ &= \prod_{j=1}^{m-1} \frac{P_{x_j}(W_j | C_j, \mathbf{A}^{(j-1)}, \mathbf{X}^{(j-1)}) P_{x_m}(C_j | \mathbf{A}^{(j-1)}, \mathbf{X}^{(j-1)})}{P_{x_m}(W_j, X_j | C_j, \mathbf{A}^{(j-1)}, \mathbf{X}^{(j-1)})} \\ &= \prod_{j=1}^{m-1} \frac{P_{x_j}(W_j | C_j, \mathbf{A}^{(j-1)}, \mathbf{X}^{(j-1)})}{P_{x_m}(W_j, X_j | C_j, \mathbf{A}^{(j-1)}, \mathbf{X}^{(j-1)})} \\ &= \pi_0^{(m-1)}(\mathbf{A}^{(m-1)}, \mathbf{X}^{(m-1)}). \end{split}$$

**Lemma C.6 (Bias Analysis of the DML estimator for MTI).** Suppose Assumptions (2,8,9) hold. For  $i = 1, 2, \dots, m$ , let  $A_i := \{C_i, W_i\}$  and  $V_i := \{A_i, X_i\}$ . For  $i = 1, \dots, m$ , let  $B_i := \{A_i, X_{i-1}\}$  where  $X_0 := \emptyset$ . Let  $T^{dml}(\{\pi^k\}_{k=1}^{m-1}, \{\mu^k\}_{k=2}^m)$  be defined as follow:

$$T^{dml}(\{\pi^k\}_{k=1}^{m-1}, \{\mu^k\}_{k=2}^m) = \sum_{i=2}^m \mathbb{E}_{P_{x_i}}\left[\pi^{(i-1)}(\mathbf{V}^{(i-1)})\mathbb{1}_{\mathbf{x}^{(i-1)}}(\mathbf{X}^{(i-1)})\left\{\mu^{i+1}(\mathbf{B}^{(i)}, x_i) - \mu^i(\mathbf{B}^{(i-1)}, X_{i-1})\right\}\right] + \mathbb{E}_{P_{x_1}}\left[\mu^2(\mathbf{B}^{(1)}, x_1)\right]. \quad (C.5)$$

Then,

$$T^{dml}(\{\pi^k\}_{k=1}^{m-1}, \{\mu^k\}_{k=2}^m) - \mathbb{E}\left[Y|do(\mathbf{x})\right] = \sum_{i=2}^m O_{P_{x_i}}\left(\left\|\mu^i - \mu_0^i\right\| \left\|\pi^{i-1} - \pi_0^{i-1}\right\|\right)$$
(C.6)

Proof. We follow the proof technique used in (Rotnitzky et al., 2017). We first note that

$$T^{dml}(\{\pi_0^k\}_{k=1}^{m-1}, \{\mu_0^k\}_{k=2}^m) = \mathbb{E}\left[Y|do(\mathbf{x})\right].$$
(C.7)

It's easy to witness Eq. (C.7) because, for  $i = 2, 3, \dots, m$ ,

$$\begin{split} & \mathbb{E}_{P_{x_i}} \left[ \pi_0^{(i-1)}(\mathbf{V}^{(i-1)}) \mathbb{1}_{\mathbf{x}^{(i-1)}}(\mathbf{X}^{(i-1)}) \left\{ \mu_0^{i+1}(\mathbf{B}^{(i)}, x_i) - \mu_0^i(\mathbf{B}^{(i-1)}, X_{i-1}) \right\} \right] \\ & = \mathbb{E}_{P_{x_i}} \left[ \mathbb{E}_{P_{x_i}} \left[ \pi_0^{(i-1)}(\mathbf{V}^{(i-1)}) \mathbb{1}_{\mathbf{x}^{(i-1)}}(\mathbf{X}^{(i-1)}) \left\{ \mu_0^{i+1}(\mathbf{B}^{(i)}, x_i) - \mu_0^i(\mathbf{B}^{(i-1)}, X_{i-1}) \right\} \left| \mathbf{B}^{(i-1)}, X_{i-1} \right] \right] \right] \\ & = \mathbb{E}_{P_{x_i}} \left[ \pi_0^{(i-1)}(\mathbf{V}^{(i-1)}) \mathbb{1}_{\mathbf{x}^{(i-1)}}(\mathbf{X}^{(i-1)}) \left\{ \mathbb{E}_{P_{x_i}} \left[ \mu_0^{i+1}(\mathbf{B}^{(i)}, x_i) | \mathbf{B}^{(i-1)}, X_{i-1} \right] - \mu_0^i(\mathbf{B}^{(i-1)}, X_{i-1}) \right\} \right] \\ & = \mathbb{E}_{P_{x_i}} \left[ \pi_0^{(i-1)}(\mathbf{V}^{(i-1)}) \mathbb{1}_{\mathbf{x}^{(i-1)}}(\mathbf{X}^{(i-1)}) \left\{ \mu_0^i(\mathbf{B}^{(i-1)}, X_{i-1}) - \mu_0^i(\mathbf{B}^{(i-1)}, X_{i-1}) \right\} \right] \\ & = 0, \end{split}$$

where the equation  $\stackrel{1}{=}$  holds by the law of total expectation. Therefore,

$$T^{dml}(\{\pi_0^k\}_{k=1}^{m-1}, \{\mu_0^k\}_{k=2}^m) = \mathbb{E}_{P_{x_1}}\left[\mu_0^2(\mathbf{B}^{(1)}, x_1)\right] = \mathbb{E}\left[Y|do(\mathbf{x})\right],$$

where the second equation holds by Lemma C.4. Therefore, it suffices to prove the following to show Eq. (C.6):

$$T^{dml}(\{\pi^k\}_{k=1}^{m-1}, \{\mu^k\}_{k=2}^m) - T^{dml}(\{\pi_0^k\}_{k=1}^{m-1}, \{\mu_0^k\}_{k=2}^m) = \sum_{i=2}^m O_{P_{x_i}}\left(\left\|\mu^i - \mu_0^i\right\| \left\|\pi^{i-1} - \pi_0^{i-1}\right\|\right).$$
(C.8)

For  $i = 1, 2, \cdots, m - 1$ , we define a quantity

$$\omega_0^i(\mathbf{B}^{(i)}) \coloneqq \frac{P_{x_i}(\mathbf{B}^{(i)})}{P_{x_m}(\mathbf{B}^{(i)})}.$$

We note that  $\omega_0^i(\mathbf{B}^{(i)})$  is related with  $\pi$  as follow:

$$\omega_0^i(\mathbf{B}^{(i)}) = \pi_0^i(\mathbf{V}^{(i)}) P_{x_m}(X_i | \mathbf{B}^{(i)}).$$
(C.9)

To witness, consider the following:

$$\begin{split} \omega_0^i(\mathbf{B}^{(i)}) &= \frac{P_{x_i}(A_i|\mathbf{V}^{(i-1)})P_{x_i}(\mathbf{V}^{(i-1)})}{P_{x_m}(A_i|\mathbf{V}^{(i-1)})P_{x_m}(\mathbf{V}^{(i-1)})} \\ & \stackrel{2}{=} \frac{P_{x_i}(A_i|\mathbf{V}^{(i-1)})}{P_{x_m}(A_i|\mathbf{V}^{(i-1)})} \\ & = \frac{P_{x_i}(W_i,C_i|\mathbf{W}^{(i-1)},\mathbf{C}^{(i-1)},\mathbf{X}^{(i-1)})}{P_{x_m}(W_i,C_i|\mathbf{W}^{(i-1)},\mathbf{C}^{(i-1)},\mathbf{X}^{(i-1)})} \\ & \stackrel{3}{=} \frac{P_{x_i}(W_i|C_i,\mathbf{W}^{(i-1)},\mathbf{C}^{(i-1)},\mathbf{X}^{(i-1)})}{P_{x_m}(W_i|C_i,\mathbf{W}^{(i-1)},\mathbf{C}^{(i-1)},\mathbf{X}^{(i-1)})} \\ & = \pi_0^i(\mathbf{W}^{(i)},\mathbf{C}^{(i)},\mathbf{X}^{(i)})\frac{P_{x_m}(W_i,X_i|C_i,\mathbf{W}^{(i-1)},\mathbf{C}^{(i-1)},\mathbf{X}^{(i-1)})}{P_{x_m}(W_i|C_i,\mathbf{W}^{(i-1)},\mathbf{C}^{(i-1)},\mathbf{X}^{(i-1)})} \\ & = \pi_0^i(\mathbf{V}^{(i)})P_{x_m}(X_i|\mathbf{B}^{(i)}), \end{split}$$

where

- $\stackrel{2}{=}$  holds since  $X_i$  is non-descendent to  $\mathbf{V}^{(i-1)}$ , so that  $P_{x_i}(\mathbf{V}^{(i-1)}) = P_{x_m}(\mathbf{V}^{(i-1)})$ .
- $\stackrel{3}{=}$  holds by Assumption 9.

To simplify the notation, we sometimes simply denote  $\omega_0^i(\mathbf{B}^{(i)})$  as  $\omega_0^i$ ;  $\mu^i(\mathbf{B}^{(i-1)}, X_{i-1})$  as  $\mu^i$ ;  $\mu^i(\mathbf{B}^{(i-1)}, x_{i-1})$  as  $\overline{\mu^i}$ ;

and  $\pi^i(\mathbf{V}^{(i)})$  as  $\pi^i$ .

Then,  $T^{dml}(\{\pi^k\}_{k=1}^{m-1},\{\mu^k\}_{k=2}^m)$  in Eq. (C.5) can be rewritten as

$$T^{dml}(\{\pi^k\}_{k=1}^{m-1}, \{\mu^k\}_{k=2}^m) = \sum_{i=2}^m \mathbb{E}_{P_{x_m}} \left[ \omega_0^i \pi^{(i-1)} \mathbb{1}_{\mathbf{x}^{(i-1)}} (\mathbf{X}^{(i-1)}) \left\{ \overline{\mu}^{i+1} - \mu^i \right\} + \omega_0^1 \overline{\mu}^2 \right],$$
(C.10)

where  $\overline{\mu}^{m+1} \coloneqq Y$ .

For each  $k = 1, 2, \dots, m$ , we define a quantity  $Q_k$  as follow:

$$Q_k \coloneqq Q_k (\{\pi^j\}_{j=k}^{m-1}, \{\mu^j\}_{j=k+1}^m) \coloneqq \omega_0^k \overline{\mu}^{k+1} + \sum_{i=k+1}^m \omega_0^i \pi^{(k:i-1)} \mathbb{1}_{\mathbf{x}^{(k:i-1)}} (\mathbf{X}^{(k:i-1)}) \{\overline{\mu}^{i+1} - \mu^i\}.$$
(C.11)

Note  $Q_m = Y$  and  $\mathbb{E}_{P_{x_m}}[Q_1] = T^{dml}(\{\pi^k\}_{k=1}^{m-1}, \{\mu^k\}_{k=2}^m)$  defined in Eq. (C.10). We note that

$$\mathbb{E}_{P_{x_m}} \left[ Q_1 - \omega_0^1 \overline{\mu}_0^2 \right] \stackrel{4}{=} T^{dml}(\{\pi^k\}_{k=1}^{m-1}, \{\mu^k\}_{k=2}^m) - \mathbb{E}_{P_{x_1}} \left[ \mu_0^2(\mathbf{B}^{(1)}, x_1) \right]$$
  
$$\stackrel{5}{=} T^{dml}(\{\pi^k\}_{k=1}^{m-1}, \{\mu^k\}_{k=2}^m) - \mathbb{E}\left[ Y | do(\mathbf{x}) \right]$$
  
$$= \mathbf{l.h.s. of Eq. (C.6),}$$

where

- $\stackrel{4}{=}$  holds since  $\mathbb{E}_{P_{x_m}} \left[ \omega_0^1(\mathbf{B}^{(1)}) \mu^2(\mathbf{B}^{(1)}, x_1) \right] = \mathbb{E}_{P_{x_1}} \left[ \mu^1(\mathbf{B}^{(1)}, x_1) \right].$
- $\stackrel{5}{=}$  holds by Lemma C.4.

Motivating from the fact that  $\mathbb{E}_{P_{x_m}} \left[ Q_1 - \omega_0^1 \overline{\mu}_0^2 \right] = 1.h.s.$  of Eq. (C.6), we establish a following induction hypothesis. For  $\overline{P}_{x_m}^{i-1} \coloneqq P_{x_m}(\cdot | \mathbf{V}^{(i-1)})$ , the induction hypothesis is given as follow:

Hypothesis: 
$$\mathbb{E}_{\overline{P}_{x_m}^{k-1}} \left[ Q_k - \omega_0^k \overline{\mu}_0^{k+1} \right] = \sum_{i=k+1}^m O_{\overline{P}_{x_i}^{k-1}} \left( \left\| \mu^i - \mu_0^i \right\| \left\| \pi^{i-1} - \pi_0^{i-1} \right\| \right), \text{ for } k \in \{2, \cdots, m-1\}$$
 (C.12)

We first verify the hypothesis Eq. (C.12) for k = m - 1.

$$\begin{split} & \mathbb{E}_{\overline{P}_{x_m}^{m-2}} \left[ Q_{m-1} - \omega_0^{m-1} \overline{\mu}_0^m \right] \\ &= \mathbb{E}_{\overline{P}_{x_m}^{m-2}} \left[ \omega_0^{m-1} \overline{\mu}^m + \pi^{m-1} \mathbb{1}_{x_{m-1}} (X_{m-1}) \left\{ Y - \mu^m \right\} - \omega_0^{m-1} \overline{\mu}_0^m \right] \\ & \stackrel{6}{=} \mathbb{E}_{\overline{P}_{x_m}^{m-2}} \left[ \omega_0^{m-1} \overline{\mu}^m + \pi^{m-1} \mathbb{1}_{x_{m-1}} (X_{m-1}) \left\{ \mu_0^m - \mu^m \right\} - \omega_0^{m-1} \overline{\mu}_0^m \right] \\ &= \mathbb{E}_{\overline{P}_{x_m}^{m-2}} \left[ \omega_0^{m-1} \left\{ \overline{\mu}^m - \overline{\mu}_0^m \right\} + \pi^{m-1} \mathbb{1}_{x_{m-1}} (X_{m-1}) \left\{ \mu_0^m - \mu^m \right\} \right] \\ & \stackrel{7}{=} \mathbb{E}_{\overline{P}_{x_m}^{m-2}} \left[ \omega_0^{m-1} \frac{\mathbb{1}_{x_{m-1}} (X_{m-1})}{P_{x_m} (X_{m-1} | \mathbf{B}^{(m-1)})} \left\{ \mu^m - \mu_0^m \right\} + \pi^{m-1} \mathbb{1}_{x_{m-1}} (X_{m-1}) \left\{ \mu_0^m - \mu^m \right\} \right] \\ & \stackrel{8}{=} \mathbb{E}_{\overline{P}_{x_m}^{m-2}} \left[ \pi_0^{m-1} \mathbb{1}_{x_{m-1}} (X_{m-1}) \left\{ \mu^m - \mu_0^m \right\} + \pi^{m-1} \mathbb{1}_{x_{m-1}} (X_{m-1}) \left\{ \mu_0^m - \mu^m \right\} \right] \\ &= \mathbb{E}_{\overline{P}_{x_m}^{m-2}} \left[ \mathbb{1}_{x_{m-1}} (X_{m-1}) \left\{ \mu^m - \mu_0^m \right\} \left\{ \pi_0^{m-1} - \pi^{m-1} \right\} \right] \\ &= \mathbb{E}_{P_{x_m}} \left[ \mathbb{1}_{x_{m-1}} (X_{m-1}) \left\{ \mu^m - \mu_0^m \right\} \left\{ \pi_0^{m-1} - \pi^{m-1} \right\} \right] \\ &= \mathbb{E}_{\overline{P}_{x_m}^{m-2}} \left( \left\| \mu^m - \mu_0^m \right\| \left\| \pi^{m-1} - \pi_0^{m-1} \right\| \right), \end{split}$$

where

- $\stackrel{6}{=}$  holds by the total law of expectation.
- $\stackrel{7}{=}$  holds since

$$\mathbb{E}_{P_{x_m}}\left[\mu^{m-1}(\mathbf{B}^{(m-1)}, x_{m-1}) \middle| \mathbf{V}^{(m-2)}\right] = \mathbb{E}_{P_{x_m}}\left[\mu^{m-1}(\mathbf{B}^{(m-1)}, X_{m-1}) \frac{\mathbb{1}_{x_{m-1}}(X_{m-1})}{P_{x_m}(X_{m-1}|\mathbf{B}^{(m-1)})} \middle| \mathbf{V}^{(m-2)}\right]$$

- $\stackrel{8}{=}$  holds by the definition of  $\omega_0^{m-1}$ .
- $\stackrel{9}{=}$  holds by applying Cauchy-Schwarz inequality.

Now, we suppose Eq. (C.12) holds for some  $k + 1 \in \{2, \dots, m-1\}$ . Then, we will show that Eq. (C.12) holds for k. Toward this end, we first rewrite  $Q_k$  in Eq. (C.11) in a recursive form. For any  $k + 1 \in \{2, \dots, m-1\}$ , the following relation can be derived from Eq. (C.11):

$$\pi^{k} \mathbb{1}_{x_{k}}(X_{k}) \left\{ Q_{k+1} - \omega_{0}^{k+1} \overline{\mu}^{k+2} \right\} = \sum_{i=k+2}^{m} \omega_{0}^{i} \pi^{(k:i-1)} \mathbb{1}_{\mathbf{x}^{(k:i-1)}} \left( \mathbf{X}^{(k:i-1)} \right) \left\{ \overline{\mu}^{i+1} - \mu^{i} \right\}.$$

Therefore, for each  $k = 1, 2, \cdots, m - 1$ ,

$$Q_k(\{\pi^j\}_{j=k}^{m-1},\{\mu^j\}_{j=k+1}^m) = \omega_0^k \overline{\mu}^{k+1} + \omega_0^{k+1} \pi^k \mathbb{1}_{x_k}(X_k) \left\{ \overline{\mu}^{k+2} - \mu^{k+1} \right\} + \pi^k \mathbb{1}_{x_k}(X_k) \left\{ Q_{k+1} - \omega_0^{k+1} \overline{\mu}^{k+2} \right\}.$$

Then,

$$\begin{split} &\mathbb{E}_{\mathcal{F}_{xm}^{k-1}}\left[Q_{k}-\omega_{0}^{k}\overline{\mu}^{k+1}\right] \\ &=\mathbb{E}_{\mathcal{F}_{xm}^{k-1}}\left[\omega_{0}^{k}\overline{\mu}^{k+1}+\omega_{0}^{k+1}\pi^{k}\mathbbm{1}_{x_{k}}(X_{k})\left\{\overline{\mu}^{k+2}-\mu^{k+1}\right\}+\pi^{k}\mathbbm{1}_{x_{k}}(X_{k})\left\{Q_{k+1}-\omega_{0}^{k+1}\overline{\mu}^{k+2}\right\}-\omega_{0}^{k}\overline{\mu}_{0}^{k+1}\right] \\ &=\mathbb{E}_{\mathcal{F}_{xm}^{k-1}}\left[\omega_{0}^{k}\overline{\mu}^{k+1}+\omega_{0}^{k+1}\pi^{k}\mathbbm{1}_{x_{k}}(X_{k})\left\{\overline{\mu}^{k+2}-\mu^{k+1}\right\}+\omega_{0}^{k+1}\pi^{k}\mathbbm{1}_{x_{k}}(X_{k})\left\{\overline{\mu}^{k+2}-\overline{\mu}^{k+2}\right\}-\omega_{0}^{k}\overline{\mu}_{0}^{k+1}\right] \\ &=\mathbb{E}_{\mathcal{F}_{xm}^{k-1}}\left[\omega_{0}^{k}\overline{\mu}^{k+1}+\omega_{0}^{k+1}\pi^{k}\mathbbm{1}_{x_{k}}(X_{k})\left\{\overline{\mu}^{k+2}-\mu^{k+1}\right\}-\omega_{0}^{k}\overline{\mu}_{0}^{k+1}\right] \\ &=\mathbb{E}_{\mathcal{F}_{xm}^{k-1}}\left[\omega_{0}^{k}\left\{\overline{\mu}^{k+1}+\omega_{0}^{k+1}\pi^{k}\mathbbm{1}_{x_{k}}(X_{k})\left\{\overline{\mu}^{k+2}-\mu^{k+1}\right\}\right] \\ &=\mathbb{E}_{\mathcal{F}_{xm}^{k-1}}\left[\omega_{0}^{k}\left\{\overline{\mu}^{k+1}-\overline{\mu}_{0}^{k+1}\right\}+\omega_{0}^{k+1}\pi^{k}\mathbbm{1}_{x_{k}}(X_{k})\left\{\overline{\mu}^{k+2}-\mu^{k+1}\right\}\right] \\ &=\mathbb{E}_{\mathcal{F}_{xm}^{k-1}}\left[\omega_{0}^{k}\left\{\overline{\mu}^{k+1}-\overline{\mu}_{0}^{k+1}\right\}+\omega_{0}^{k+1}\pi^{k}\mathbbm{1}_{x_{k}}(X_{k})\left\{\overline{\mu}^{k+2}-\mu^{k+1}\right\}\right] \\ &=\mathbb{E}_{\mathcal{F}_{xm}^{k-1}}\left[\omega_{0}^{k}\left\{\overline{\mu}^{k+1}-\overline{\mu}_{0}^{k+1}\right\}+\pi^{k}\mathbbm{1}_{x_{k}}(X_{k})\left\{\mu^{k+1}-\mu^{k+1}_{0}\right\}\right] \\ &=\mathbb{E}_{\mathcal{F}_{xm}^{k-1}}\left[\omega_{0}^{k}\left\{\overline{\mu}^{k+1}-\overline{\mu}_{0}^{k+1}\right\}+\pi^{k}\mathbbm{1}_{x_{k}}(X_{k})\left\{\mu^{k+1}-\mu^{k+1}_{0}\right\}\right] \\ &=\mathbb{E}_{\mathcal{F}_{xm}^{k-1}}\left[\pi_{0}^{k}\mathbbm_{x_{k}}(X_{k})\left\{\mu^{k+1}-\mu^{k+1}_{0}\right\}+\pi^{k}\mathbbm_{x_{k}}(X_{k})\left\{\mu^{k+1}-\mu^{k+1}_{0}\right\}\right] \\ &=\mathbb{E}_{\mathcal{F}_{xm}^{k-1}}\left[\pi_{0}^{k}\mathbbm_{x_{k}}(X_{k})\left\{\mu^{k+1}-\mu^{k+1}_{0}\right\}\right] \\ &=\mathbb{E}_{\mathcal{F}_{xm}^{k-1}}\left[\omega_{0}^{k+1}\mathbbm_{x_{k}}(X_{k})\left\{\mu^{k+1}-\mu^{k+1}_{0}\right\}\left\{\pi_{0}^{k}-\pi^{k}\right\}\right] \\ &=\mathbb{E}_{\mathcal{F}_{xm}^{k-1}}\left[\pi_{x_{k}}(X_{k})\left\{\mu^{k+1}-\mu^{k+1}_{0}\right\}\left\{\pi_{0}^{k}-\pi^{k}\right\}\right\right] \\ &=\mathbb{E}_{\mathcal{F}_{xm}^{k-1}}\left[\pi_{x_{k}}(X_{k})\left\{\mu^{k+1}-\mu^{k+1}_{0}\right\}\left\{\pi_{0}^{k}-\pi^{k}\right\}\right] \\ &=\mathbb{E}_{\mathcal{F}_{xm}^{k-1}}\left[\pi_{x_{k}}(X_{k})\left\{\mu^{k+1}-\mu^{k+1}_{0}\right\}\left\{\pi_{0}^{k}-\pi^{k}\right\}\right\right] \\ &=\mathbb{E}_{\mathcal{F}_{xm}^{k-1}}\left[\pi_{x_{k}}(X_{k})\left\{\mu^{k+1}-\mu^{k+1}_{0}\right\}\left\{\pi_{0}^{k}-\pi^{k}\right\}\right] \\ &=\mathbb{E}_{\mathcal{F}_{xm}^{k-1}}\left[\pi_{x_{k}}(X_{k})\left\{\mu^{k+1}-\mu^{k+1}_{0}\right\}\left\{\pi_{0}^{k}-\pi^{k}\right\}\right\right] \\ &=\mathbb{E}_{\mathcal{F}_{xm}^{k-1}}\left[\mathbb{E}_{\mathcal{F}_{xm}^{k-1}}\left[\pi_{x_{k}}(X_{k})\left\{\mu^{k+1}-\mu^{k+1}_{0}\right\}\right] \\ &=\mathbb{E}_{\mathcal{F}_{x$$

where

•  $\stackrel{10}{=}$  holds since

$$\begin{split} & \mathbb{E}_{P_{x_m}} \left[ \omega_0^{k+1}(\mathbf{B}^{(k+1)}) \pi^k(\mathbf{V}^{(k)}) \mathbb{1}_{x_k}(X_k) \mu_0^{k+2}(\mathbf{B}^{(k+1)}, x_{k+1}) \Big| \mathbf{V}^{(k-1)} \right] \\ &= \mathbb{E}_{P_{x_m}} \left[ \mathbb{E}_{P_{x_m}} \left[ \omega_0^{k+1}(\mathbf{B}^{(k+1)}) \pi^k(\mathbf{V}^{(k)}) \mathbb{1}_{x_k}(X_k) \mu_0^{k+2}(\mathbf{B}^{(k+1)}, x_{k+1}) \Big| \mathbf{V}^{(k)} \right] \Big| \mathbf{V}^{(k-1)} \right] \\ &= \mathbb{E}_{P_{x_m}} \left[ \pi^k(\mathbf{V}^{(k)}) \mathbb{1}_{x_k}(X_k) \mathbb{E}_{P_{x_m}} \left[ \omega_0^{k+1}(\mathbf{B}^{(k+1)}) \mu_0^{k+2}(\mathbf{B}^{(k+1)}, x_{k+1}) \Big| \mathbf{V}^{(k)} \right] \Big| \mathbf{V}^{(k-1)} \right] \\ &= \mathbb{E}_{P_{x_m}} \left[ \pi^k(\mathbf{V}^{(k)}) \mathbb{1}_{x_k}(X_k) \mathbb{E}_{P_{x_m}} \left[ \frac{P_{x_{k+1}}(\mathbf{B}^{(k+1)})}{P_{x_m}(\mathbf{B}^{(k+1)})} \mu_0^{k+2}(\mathbf{B}^{(k+1)}, x_{k+1}) \Big| \mathbf{V}^{(k)} \right] \Big| \mathbf{V}^{(k-1)} \right] \\ &= \mathbb{E}_{P_{x_m}} \left[ \pi^k(\mathbf{V}^{(k)}) \mathbb{1}_{x_k}(X_k) \mathbb{E}_{P_{x_{k+1}}} \left[ \mu_0^{k+2}(\mathbf{B}^{(k+1)}, x_{k+1}) \Big| \mathbf{V}^{(k)} \right] \Big| \mathbf{V}^{(k-1)} \right] \\ &= \mathbb{E}_{P_{x_m}} \left[ \pi^k(\mathbf{V}^{(k)}) \mathbb{1}_{x_k}(X_k) \mu_0^{k+1}(\mathbf{V}^{(k)}) \Big| \mathbf{V}^{(k-1)} \right]. \end{split}$$

•  $\stackrel{11}{=}$  holds since

$$\begin{split} & \mathbb{E}_{P_{x_m}} \left[ \omega_0^k(\mathbf{B}^{(k)}) \left\{ \mu^{k+1}(\mathbf{B}^{(k)}, x_k) - \mu_0^{k+1}(\mathbf{B}^{(k)}, x_k) \right\} \left| \mathbf{V}^{(k-1)} \right] \\ &= \mathbb{E}_{P_{x_m}} \left[ \omega_0^k(\mathbf{B}^{(k)}) \frac{\mathbb{1}_{x_k}(X_k)}{P_{x_m}(X_k | \mathbf{B}^{(k)})} \left\{ \mu^{k+1}(\mathbf{B}^{(k)}, X_k) - \mu_0^{k+1}(\mathbf{B}^{(k)}, X_k) \right\} \left| \mathbf{V}^{(k-1)} \right] \right] \\ &= \mathbb{E}_{P_{x_m}} \left[ \pi_0^k(\mathbf{V}^{(k)}) P_{x_m}(X_k | \mathbf{B}^{(k)}) \frac{\mathbb{1}_{x_k}(X_k)}{P_{x_m}(X_k | \mathbf{B}^{(k)})} \left\{ \mu^{k+1}(\mathbf{B}^{(k)}, X_k) - \mu_0^{k+1}(\mathbf{B}^{(k)}, X_k) \right\} \left| \mathbf{V}^{(k-1)} \right] \right] \\ &= \mathbb{E}_{P_{x_m}} \left[ \pi_0^k(\mathbf{V}^{(k)}) \mathbb{1}_{x_k}(X_k) \left\{ \mu^{k+1}(\mathbf{B}^{(k)}, X_k) - \mu_0^{k+1}(\mathbf{B}^{(k)}, X_k) \right\} \left| \mathbf{V}^{(k-1)} \right]. \end{split}$$

•  $\stackrel{12}{=}$  and  $\stackrel{13}{=}$  hold since

$$\begin{split} & \mathbb{E}_{P_{x_m}} \left[ \mathbbm{1}_{x_k}(X_k) \left\{ \mu^{k+1}(\mathbf{V}^{(k)}) - \mu_0^{k+1}(\mathbf{V}^{(k)}) \right\} \left\{ \pi_0^k(\mathbf{V}^{(k)}) - \pi^k(\mathbf{V}^{(k)}) \right\} \left| \mathbf{V}^{(k-1)} \right] \\ &= \mathbb{E}_{P_{x_m}} \left[ \omega_0^{k+1}(\mathbf{V}^{(k)}) \mathbbm{1}_{x_k}(X_k) \left\{ \mu^{k+1}(\mathbf{V}^{(k)}) - \mu_0^{k+1}(\mathbf{V}^{(k)}) \right\} \left\{ \pi_0^k(\mathbf{V}^{(k)}) - \pi^k(\mathbf{V}^{(k)}) \right\} \left| \mathbf{V}^{(k-1)} \right] \\ &= \mathbb{E}_{P_{x_m}} \left[ \frac{P_{x_{k+1}}(\mathbf{V}^{(k)})}{P_{x_m}(\mathbf{V}^{(k)})} \mathbbm{1}_{x_k}(X_k) \left\{ \mu^{k+1}(\mathbf{V}^{(k)}) - \mu_0^{k+1}(\mathbf{V}^{(k)}) \right\} \left\{ \pi_0^k(\mathbf{V}^{(k)}) - \pi^k(\mathbf{V}^{(k)}) \right\} \left| \mathbf{V}^{(k-1)} \right] \\ &= \mathbb{E}_{P_{x_{k+1}}} \left[ \mathbbm{1}_{x_k}(X_k) \left\{ \mu^{k+1}(\mathbf{V}^{(k)}) - \mu_0^{k+1}(\mathbf{V}^{(k)}) \right\} \left\{ \pi_0^k(\mathbf{V}^{(k)}) - \pi^k(\mathbf{V}^{(k)}) \right\} \left| \mathbf{V}^{(k-1)} \right], \end{split}$$

where the second equation hold since

$$\omega_0^{k+1}(\mathbf{V}^{(k)}) = \frac{P_{x_{k+1}}(\mathbf{V}^{(k)})}{P_{x_m}(\mathbf{V}^{(k)})} = 1$$

since  $X_{k+1}, X_m$  are non-descendants of  $\mathbf{V}^{(k)}$  so that  $P_{x_{k+1}}(\mathbf{V}^{(k)}) = P_{x_m}(\mathbf{V}^{(k)})$ .

- $\stackrel{14}{=}$  holds by the induction hypothesis.
- $\stackrel{15}{=}$  holds by Cauchy-Schwarz inequality.

Therefore, the induction hypothesis in Eq. (C.12) holds for all  $k = 1, 2, \dots, m-1$ . Therefore,

l.h.s. of Eq. (C.6) = 
$$\mathbb{E}_{P_{x_m}} \left[ Q_1 - \omega_0^1 \overline{\mu}_0^2 \right] = \sum_{i=2}^m O_{P_{x_i}} \left( \left\| \mu^i - \mu_0^i \right\| \left\| \pi^{i-1} - \pi_0^{i-1} \right\| \right),$$

where the second equation holds by plugging k = 1 into the verified hypothesis in Eq. (C.12). This completes the proof.  $\Box$ 

**Lemma C.7** (Error analysis of the DML estimator for MTI). Suppose Assumptions (2,8,9) hold. Let  $T^{dml}$  denote the estimator defined in Def. 9. Then,

$$T^{dml} - \mathbb{E}\left[Y|do(\mathbf{x})\right] = \sum_{i=1}^{m} R_i + \sum_{i=2}^{m} O_{P_{x_i}}\left(\left\|\mu^i - \mu_0^i\right\| \left\|\pi^{i-1} - \pi_0^{i-1}\right\|\right)$$

where  $R_i$  for  $i = 1, 2, \dots, m$  are variables converging in mean-zero normal distribution at  $n_i^{-1/2}$  rates.

**Proof of Lemma C.7.** Throughout the proof, we set  $\mathbf{A}_i := \{C_i, W_i\}$  for all  $i = 1, 2, \dots, m$ . We will use  $V_i := \{A_i, X_i\}$  or all  $i = 1, 2, \dots, m$ . We will use  $B_i := \{A_i, X_{i-1}\}$  or all  $i = 1, 2, \dots, m$ , where  $X_0 := \emptyset$ . To simplify the notation, we sometimes simply denote  $\mu^i(\mathbf{B}^{(i-1)}, X_{i-1})$  as  $\mu^i$ ;  $\mu^i(\mathbf{B}^{(i-1)}, x_{i-1})$  as  $\overline{\mu^i}$ ; and  $\pi^i(\mathbf{V}^{(i)})$  as  $\pi^i$ .

Let  $T^{dml}({\pi^k}_{k=1}^{m-1}, {\mu^k}_{k=2}^m)$  be a quantity defined in Eq. (C.5). We first note that

$$T^{dml}(\{\pi_0^k\}_{k=1}^{m-1}, \{\mu_0^k\}_{k=2}^m) = \mathbb{E}\left[Y|do(\mathbf{x})\right]$$

by Eq. (C.7). Then, by Lemma C.3,

$$T^{dml} - \mathbb{E}\left[Y|do(\mathbf{x})\right] = T^{dml} - T^{dml}(\{\pi_0^k\}_{k=1}^{m-1}, \{\mu_0^k\}_{k=2}^m)$$

$$= \sum_{i=2}^m \mathbb{E}_{P_{x_i} - D_{x_i}} \left[\pi_0^{(i-1)} \mathbb{1}_{\mathbf{x}^{(i-1)}}(\mathbf{X}^{(i-1)}) \left\{\overline{\mu}_0^{i+1} - \mu_0^i\right\}\right] + \mathbb{E}_{P_{x_1} - D_{x_1}} \left[\overline{\mu}_0^2\right]$$

$$+ \sum_{i=2}^m \mathbb{E}_{P_{x_i} - D_{x_i}} \left[\pi_0^{(i-1)} \mathbb{1}_{\mathbf{x}^{(i-1)}}(\mathbf{X}^{(i-1)}) \left\{\overline{\mu}_0^{i+1} - \mu_0^i\right\} - \pi^{(i-1)} \mathbb{1}_{\mathbf{x}^{(i-1)}}(\mathbf{X}^{(i-1)}) \left\{\overline{\mu}_0^{i+1} - \mu^i\right\}\right] + \mathbb{E}_{P_{x_1} - D_{x_1}} \left[\overline{\mu}_0^2 - \overline{\mu}_0^2\right]$$

$$(C.14)$$

$$(C.14)$$

$$+\sum_{i=2}^{m} \mathbb{E}_{P_{x_i}} \left[ \pi^{(i-1)} \mathbb{1}_{\mathbf{x}^{(i-1)}} (\mathbf{X}^{(i-1)}) \left\{ \overline{\mu}^{i+1} - \mu^i \right\} \right] + \mathbb{E}_{P_{x_1}} \left[ \overline{\mu}^2 - \overline{\mu}_0^2 \right].$$
(C.15)

We first note that

Eq. (C.14) = 
$$\sum_{i=1}^{m} O_{P_{x_i}}(n_i^{-1/2})$$

under Assumptions (2,8) by Lemma C.3.

Then,

Eq. (C.13) + Eq. (C.14) = 
$$\sum_{i=1}^{m} R_i$$
,

where  $R_i$  for  $i = 1, 2, \dots, m$  are variables converging in mean-zero normal distribution, by the central limit theorem and Slutsky's theorem.

Finally

Eq. (C.15) = 
$$T^{dml}(\{\pi^k\}_{k=1}^{m-1}, \{\mu^k\}_{k=2}^m) - \mathbb{E}[Y|do(\mathbf{x})]$$
  
=  $\sum_{i=k+1}^m O_{\overline{P}_{x_i}^{k-1}}(\|\mu^i - \mu_0^i\| \|\pi^{i-1} - \pi_0^{i-1}\|),$ 

where the second equation holds by Lemma C.6.

**Theorem 6** (Error analysis of the estimators for MTI). Under Assumptions (2,7,8,9) and AC-MTI in Def. 7, the errors of the estimators in Def. 9, denoted  $\epsilon^{est} \coloneqq T^{est} - \mathbb{E}[Y|do(\mathbf{x})]$  for  $est \in \{reg, pw, dml\}$ , are:

$$\begin{split} \epsilon^{reg} &= R_1 + O_{P_{x_1}} \left( \| \mu^1 - \mu_0^1 \| \right), \\ \epsilon^{pw} &= R_m + O_{P_{x_m}} (\| \pi^{(m-1)} - \pi_0^{(m-1)} \|), \\ \epsilon^{dml} &= \sum_{i=1}^m R_i + \sum_{i=2}^m O_{P_{x_i}} \left( \| \mu^i - \mu_0^i \| \| \pi^{i-1} - \pi_0^{i-1} \| \right), \end{split}$$

where  $R_i$  is a random variable such that  $\sqrt{n_i}R_i$  converges in distribution to a zero-mean normal random variable, where  $n_i \coloneqq |D_{x_i}| \text{ for } i \in \{1, \cdots, m\}.$ 

*Proof of Theorem 6.* The proof is complete by Lemmas (C.4, C.5, C.7).

Corollary 6 (Multiply robustness of the DML estimators (Corollary of Thm. 6)). Suppose Asumptions (2,7,8,9) and AC-MTI in Def. 7 hold. For  $i = 2, \dots, m-1$ , suppose either  $\pi^{i-1} = \pi_0^{i-1}$  or  $\mu^i = \mu_0^i$ . Then,  $T^{dml}$  in Def. 9 is an unbiased estimator of  $\mathbb{E}[Y|do(\mathbf{x})]$ .

**Proof of Corollary 6.** Let  $T^{dml}(\{\pi^k\}_{k=1}^{m-1}, \{\mu^k\}_{k=2}^m)$  be a quantity defined in Eq. (C.5). Let

$$T^{dml,i} \coloneqq \mathbb{E}_{D_i} \left[ \pi^{(i-1)}(\mathbf{V}^{(i-1)}) \mathbb{1}_{\mathbf{x}^{(i-1)}}(\mathbf{X}^{(i-1)}) \left\{ \mu^{i+1}(\mathbf{B}^{(i)}, x_i) - \mu^i(\mathbf{V}^{(i)}) \right\} \right], \ i = 2, \cdots, m$$

and

$$T^{dml,1} \coloneqq \mathbb{E}_{D_1}\left[\mu^2(B_1, x_1)\right].$$

Under the assumption that samples are i.i.d.,

$$\sum_{i=1}^{m} \mathbb{E}_{P_{x_i}} \left[ T^{dml,i} \right] = T^{dml}(\{\pi^k\}_{k=1}^{m-1}, \{\mu^k\}_{k=2}^m).$$

Then,

$$\begin{split} &\sum_{i=1}^{m} \mathbb{E}_{P_{x_i}} \left[ T^{dml,i} \right] - \mathbb{E} \left[ Y | do(\mathbf{x}) \right] \\ &= T^{dml}(\{\pi^k\}_{k=1}^{m-1}, \{\mu^k\}_{k=2}^m) - \mathbb{E} \left[ Y | do(\mathbf{x}) \right] \\ &= \sum_{i=2}^{m} O_{P_{x_i}} \left( \left\| \mu^i - \mu_0^i \right\| \left\| \pi^{i-1} - \pi_0^{i-1} \right\| \right) \\ &= 0, \end{split}$$

where the third equation holds by Lemma C.6, and the last equation holds under the given condition.

#### C.8. Proof of Theorem 7

Definition 10 (Adjustment criterion for gMTI (AC-gMTI)). Let  $\mathbf{Z} \coloneqq \{Z_1, \cdots, Z_m\} \subseteq \mathbf{X}$  denote the subset of treatments. Let  $\{\ell_i\}_{i=1}^m \subseteq \{1, 2, \cdots, |\mathbf{X}|\}$  denote the index of  $\mathbf{Z}$ ; i.e.,  $\mathbf{Z} = \{X_{\ell_1}, \cdots, X_{\ell_m}\}$ . Let  $\overline{X}_1 \coloneqq \{X_j\}_{j \leq \ell_1}$ ,  $\overline{X}_{m+1} \coloneqq \{X_j\}_{j > \ell_m}$ , and  $\overline{X}_i \coloneqq \{X_j\}_{\ell_{i-1} < j \le \ell_i}$  for  $i = 2, 3, \cdots, m$ . An ordered set  $\mathbf{A} \coloneqq \{\mathbf{A}_1, \mathbf{A}_2, \cdots, \mathbf{A}_m\}$  satisfies adjustment criterion for combining multiple experiments (AC-gMTI) w.r.t.  $(\mathbf{X}, Y)$  in G if, for  $i = 1, 2, \dots, m-1$ ,

(. 1)

3. 
$$(Y \perp \mathbf{\overline{X}}^{\geq m} \setminus Z_m | \mathbf{A}^{(m-1)}, \mathbf{\overline{X}}^{(m-1)}, Z_m)_{G_{\overline{Zm}, \mathbf{\overline{X}}^{\geq m} \setminus Z_m}}$$

Assumption 10 (Positivity Assumption for AC-gMTI).  $P_{z_m}(\overline{X}_m \setminus Z_m, \overline{X}_{m+1} | \mathbf{A}^{(m-1)}, \overline{\mathbf{X}}^{(m-1)})$  and  $\{P_{z_i}(\mathbf{A}_i | \mathbf{A}^{(i-1)}, \overline{\mathbf{X}}^{(i-1)}), P_{z_{i+1}}(\mathbf{A}_i | \mathbf{A}^{(i-1)}, \overline{\mathbf{X}}^{(i-1)})\}_{i=1}^{m-1}$ ,  $\{P^{i+1}(\overline{X}_i | \mathbf{A}^{(i)}, \overline{\mathbf{X}}^{(i-1)})\}_{i=1}^{m-1}$  are strictly positive distributions  $\forall i \in \{1, \cdots, m\}, \forall z_i \in \mathfrak{D}_{Z_i}$ .

**Theorem 7** (Identification through AC-gMTI). Suppose AC-gMTI in Def. 10 and Assumption 10 hold. Then,  $\mathbb{E}[Y|do(\mathbf{x})]$  is identifiable from  $\{P_{rand(Z_i)}(\mathbf{A}^{(i)}, \overline{\mathbf{X}}^{(i)})\}_{i=1}^{m}$  and given as follows. Denote

$$\mu_0^m \coloneqq \mathbb{E}_{P_{z_m}} \left[ Y | \mathbf{A}^{(m-1)}, \mathbf{X} \backslash Z_m \right]$$
  
$$\overline{\mu}_0^m \coloneqq \mathbb{E}_{P_{z_m}} \left[ Y | \mathbf{A}^{(m-1)}, \overline{\mathbf{x}}_{m-1:m+1}, \overline{\mathbf{X}}^{(m-2)} \right]$$
  
$$\iota_0^{m-1} \coloneqq \mathbb{E}_{P_{z_{m-1}}} \left[ \overline{\mu}_0^m \middle| \mathbf{A}^{(m-2)}, \overline{\mathbf{X}}^{(m-2)} \right],$$

where  $\overline{X}_{m-1:m+1} \coloneqq \{\overline{X}_{m-1}, \overline{X}_m, \overline{X}_{m+1}\}$ . For  $i = m - 2, \cdots, 2$ ,

ŀ

$$\mu_0^i \coloneqq \mathbb{E}_{P_{z_i}} \left[ \mu^{i+1}(\mathbf{A}^{(i)}, \overline{x}_i, \overline{\mathbf{X}}^{(i-1)}) \middle| \mathbf{A}^{(i-1)}, \overline{\mathbf{X}}^{(i-1)} \right],$$

and  $\overline{\mu}_0^{i+1} \coloneqq \mu_0^{i+1}(\mathbf{A}^{(i)}, \overline{x}_i, \overline{\mathbf{X}}^{(i-1)})$ . Then,

$$\mathbb{E}\left[Y(\mathbf{x})\right] = \mathbb{E}_{P_{z_1}}\left[\overline{\mu}_0^2\right].\tag{4}$$

Proof of Theorem 7. We first note that

$$\mathbb{E}\left[Y|do(\mathbf{x}\setminus\overline{\mathbf{x}}^{(m-1)}), \overline{\mathbf{x}}^{(m-1)}, \mathbf{A}^{(m-1)}\right] = \mathbb{E}\left[Y|do(\overline{x}_m\setminus z_m, z_m, \overline{x}_{m+1}), \overline{\mathbf{x}}^{(m-1)}, \mathbf{A}^{(m-1)}\right]$$
$$\stackrel{1}{=} \mathbb{E}\left[Y|do(z_m), \overline{x}_m\setminus z_m, \overline{x}_{m+1}, \overline{\mathbf{x}}^{(m-1)}, \mathbf{A}^{(m-1)}\right]$$
$$= \mathbb{E}\left[Y|do(z_m), \mathbf{x}\setminus z_m, \mathbf{A}^{(m-1)}\right]$$
$$= \mu_0^m(\mathbf{A}^{(m-1)}, \mathbf{x}\setminus z_m),$$

where

•  $\stackrel{1}{=}$  holds by the condition  $\left(Y \perp \{\overline{X}_m \setminus Z_m, \overline{X}_{m+1}\} | \mathbf{A}^{(m-1)}(\mathbf{x}), \overline{\mathbf{X}}^{(m-1)}, Z_m\right)_{G_{\overline{Z_m}, \overline{X_m \setminus Z_m, \overline{X_{m+1}}}}}$  in Def. 10. Specifically, the condition is an application of Rule 2 of do-calculus (Pearl, 2000).

We also note that

$$\begin{split} & \mathbb{E}\left[Y|do(\mathbf{x}\backslash\overline{\mathbf{x}}^{(m-2)}),\overline{\mathbf{x}}^{(m-2)},\mathbf{A}^{(m-2)}\right] \\ &= \mathbb{E}\left[\mathbb{E}\left[Y|do(\mathbf{x}\backslash\overline{\mathbf{x}}^{(m-2)}),\overline{\mathbf{x}}^{(m-2)},\mathbf{A}^{(m-1)}\right] \left| do(\mathbf{x}\backslash\overline{\mathbf{x}}^{(m-2)}),\overline{\mathbf{x}}^{(m-2)},\mathbf{A}^{(m-2)}\right. \right] \\ &\stackrel{2}{=} \mathbb{E}\left[\mathbb{E}\left[Y|do(\mathbf{x}\backslash\overline{\mathbf{x}}^{(m-1)}),\overline{\mathbf{x}}^{(m-1)},\mathbf{A}^{(m-1)}\right] \left| do(\mathbf{x}\backslash\overline{\mathbf{x}}^{(m-2)}),\overline{\mathbf{x}}^{(m-2)},\mathbf{A}^{(m-2)}\right. \right] \\ &= \mathbb{E}\left[\mu_{0}^{m}(\mathbf{A}^{(m-1)},\mathbf{x}\backslash z_{m}) \left| do(\overline{\mathbf{x}}\backslash\overline{\mathbf{x}}^{(m-2)}),\overline{\mathbf{x}}^{(m-2)},\mathbf{A}^{(m-2)}\right] \right] \\ &= \mathbb{E}\left[\mu_{0}^{m}(\mathbf{A}^{(m-1)},\mathbf{x}\backslash z_{m}) \left| do(\overline{\mathbf{x}}^{>m-1}\backslash z_{m-1},z_{m-1}),\overline{\mathbf{x}}^{(m-2)},\mathbf{A}^{(m-2)}\right] \right] \\ &\stackrel{2}{=} \mathbb{E}\left[\mu_{0}^{m}(\mathbf{A}^{(m-1)},\mathbf{x}\backslash z_{m}) \left| do(z_{m-1}),\overline{\mathbf{x}}^{(m-2)},\mathbf{A}^{(m-2)}\right] \right] \\ &= \mu_{0}^{m-1}(\mathbf{A}^{(m-2)},\overline{\mathbf{x}}^{(m-2)}), \end{split}$$

where

- $\stackrel{2}{=}$  holds since  $(Y \perp \overline{X}_i | \mathbf{A}^{(i)}, \overline{\mathbf{X}}^{(i-1)}, \overline{\mathbf{X}}^{>i})_{G_{\underline{\overline{X}_i}, \overline{\mathbf{X}}^{>i}}}$  in Def. 10 and the given positivity condition.
- $\stackrel{3}{=}$  hold since  $\left(\mathbf{A}_{i} \perp \mathbf{\overline{X}}^{>i-1} \setminus Z_{i} | \mathbf{\overline{X}}^{(i-1)}, \mathbf{A}^{(i-1)}, Z_{i} \right)_{G_{\overline{\mathbf{\overline{X}}}^{\geq i-1}}}$  in Def. 10,  $\mathbf{\overline{X}}^{>i-1} \setminus Z_{i}$  is non-ancestral to  $\mathbf{\overline{X}}^{(i-1)}, \mathbf{A}^{(i-1)}$  and the given positivity condition.

Finally, for  $i + 1 \in \{m - 2, \dots, 3\}$ , suppose

$$\mathbb{E}\left[Y|do(\mathbf{x}\setminus\overline{\mathbf{x}}^{(i)}),\overline{\mathbf{x}}^{(i)},\mathbf{A}^{(i)}\right] = \mu_0^{i+1}(\mathbf{A}^{(i)},\overline{\mathbf{x}}^{(i)}).$$

Then,

$$\begin{split} & \mathbb{E}\left[Y|do(\mathbf{x}\backslash\overline{\mathbf{x}}^{(i-1)}),\overline{\mathbf{x}}^{(i-1)},\mathbf{A}^{(i-1)}\right] \\ &= \mathbb{E}\left[\mathbb{E}\left[Y|do(\mathbf{x}\backslash\overline{\mathbf{x}}^{(i-1)}),\overline{\mathbf{x}}^{(i-1)},\mathbf{A}^{(i)}\right] \left| do(\mathbf{x}\backslash\overline{\mathbf{x}}^{(i-1)}),\overline{\mathbf{x}}^{(i-1)},\mathbf{A}^{(i-1)}\right] \right] \\ &\stackrel{4}{=} \mathbb{E}\left[\mathbb{E}\left[Y|do(\mathbf{x}\backslash\overline{\mathbf{x}}^{(i)}),\overline{\mathbf{x}}^{(i)},\mathbf{A}^{(i)}\right] \left| do(\mathbf{x}\backslash\overline{\mathbf{x}}^{(i-1)}),\overline{\mathbf{x}}^{(i-1)},\mathbf{A}^{(i-1)}\right] \right] \\ &= \mathbb{E}\left[\mu_{0}^{i+1}(\mathbf{A}^{(i)},\overline{\mathbf{x}}^{(i)}) \left| do(\mathbf{x}\backslash\overline{\mathbf{x}}^{(i-1)}),\overline{\mathbf{x}}^{(i-1)},\mathbf{A}^{(i-1)}\right] \right] \\ &= \mathbb{E}\left[\mu_{0}^{i+1}(\mathbf{A}^{(i)},\overline{\mathbf{x}}^{(i)}) \left| do(\overline{\mathbf{x}}^{\geq i}\backslash z_{i},z_{i}),\overline{\mathbf{x}}^{(i-1)},\mathbf{A}^{(i-1)}\right] \right] \\ &\stackrel{5}{=} \mathbb{E}\left[\mu_{0}^{i+1}(\mathbf{A}^{(i)},\overline{\mathbf{x}}^{(i)}) \left| do(z_{i}),\overline{\mathbf{x}}^{(i-1)},\mathbf{A}^{(i-1)}\right] \right] \\ &= \mu_{0}^{i}(\mathbf{A}^{(i-1)},\overline{\mathbf{x}}^{(i-1)}), \end{split}$$

where

- $\stackrel{4}{=}$  holds since  $(Y \perp \overline{X}_i | \mathbf{A}^{(i)}, \overline{\mathbf{X}}^{(i-1)}, \overline{\mathbf{X}}^{>i})_{G_{\underline{\overline{X}_i}, \overline{\mathbf{X}}^{>i}}}$  in Def. 10 and the given positivity condition.
- $\stackrel{5}{=}$  holds since  $\left(\mathbf{A}_{i} \perp \mathbf{\overline{X}}^{>i-1} \setminus Z_{i} | \mathbf{\overline{X}}^{(i-1)}, \mathbf{A}^{(i-1)}, Z_{i} \right)_{G_{\overline{\mathbf{X}}^{\geq i-1}}}$  in Def. 10,  $\mathbf{\overline{X}}^{>i-1} \setminus Z_{i}$  is non-ancestral to  $\mathbf{\overline{X}}^{(i-1)}, \mathbf{A}^{(i-1)}$  and the given positivity condition.

Therefore, for all  $i = m - 2, \cdots, 2$ .

$$\mathbb{E}\left[Y|do(\mathbf{x}\setminus\overline{\mathbf{x}}^{(i)}),\overline{\mathbf{x}}^{(i)},\mathbf{A}^{(i)}\right] = \mu_0^{i+1}(\mathbf{A}^{(i)},\overline{\mathbf{x}}^{(i)}).$$

Finally,

$$\begin{split} \mathbb{E}\left[Y|do(\mathbf{x})\right] &= \mathbb{E}\left[\mathbb{E}\left[Y|do(\mathbf{x}), A_{1}\right]|do(\mathbf{x})\right] \\ &\stackrel{6}{=} \mathbb{E}\left[\mathbb{E}\left[Y|do(\mathbf{x}\backslash x_{1}), x_{1}, A_{1}\right]|do(\mathbf{x})\right] \\ &= \mathbb{E}\left[\mu_{0}^{2}(A_{1}, \overline{x}_{1})|do(\mathbf{x})\right] \\ &= \mathbb{E}\left[\mu_{0}^{2}(A_{1}, \overline{x}_{1})|do(\mathbf{x}\backslash z_{1}, z_{1})\right] \\ &\stackrel{7}{=} \mathbb{E}\left[\mu_{0}^{2}(A_{1}, \overline{x}_{1})|do(z_{1})\right] \\ &= \mathbb{E}_{P_{z_{1}}}\left[\mu_{0}^{2}(A_{1}, \overline{x}_{1})\right], \end{split}$$

where

•  $\stackrel{6}{=}$  holds since  $(Y \perp \overline{X}_i | \mathbf{A}^{(i)}, \overline{\mathbf{X}}^{(i-1)}, \overline{\mathbf{X}}^{>i})_{G_{\underline{X}_i, \overline{\mathbf{X}}^{>i}}}$  in Def. 10 and the given positivity condition.

•  $\stackrel{7}{=}$  holds since  $\left(\mathbf{A}_{i} \perp \mathbf{\overline{X}}^{>i-1} \setminus Z_{i} | \mathbf{\overline{X}}^{(i-1)}, \mathbf{A}^{(i-1)}, Z_{i} \right)_{G_{\overline{\mathbf{X}}^{\geq i-1}}}$  in Def. 10,  $\mathbf{\overline{X}}^{>i-1} \setminus Z_{i}$  is non-ancestral to  $\mathbf{\overline{X}}^{(i-1)}, \mathbf{A}^{(i-1)}$  and the given positivity condition.

#### C.9. Proof of Theorem 8 and Corollary 8

**Definition 11 (Nuisances for AC-gMTI).** Nuisance functions for AC-gMTI are defined as follows: For a fixed  $\mathbf{z} := \{z_1, \cdots, z_m\} \in \mathfrak{D}_{\mathbf{Z}}, \text{ let } \{\mu_0^i\}_{i=2}^m$  be the nuisances defined in Thm. 7. For  $i = 1, \cdots, m-2, \pi_0^i := \frac{P_{z_i}(A_i | \mathbf{A}^{(i-1)}, \overline{\mathbf{X}}^{(i-1)})}{P_{z_m}(A_i, \overline{X}_i | \mathbf{A}^{(i-1)}, \overline{\mathbf{X}}^{(i-1)})},$ and  $\pi_0^{(i)} := \prod_{j=1}^i \pi_0^j(\mathbf{A}^{(j)}, \overline{\mathbf{X}}^{(j)})$ . Also,  $\pi_0^{m-1} := \frac{P_{z_{m-1}}(A_{m-1} | \mathbf{A}^{(m-2)}, \overline{\mathbf{X}}^{(m-2)})}{P_{z_m}(A_{m-1}, \overline{X}_{m-1:m+1} | \mathbf{A}^{(m-2)}, \overline{\mathbf{X}}^{(m-2)})},$  and  $p_0^{(m-1)} := \pi_0^{(m-2)} \times \pi_0^{m-1},$ where  $\overline{X}_{m-1:m+1} := \{\overline{X}_{m-1}, \overline{X}_m, \overline{X}_{m+1}\}$ . For all  $i = 1, 2, \cdots, m-1$ , we will use  $\pi^i(\mathbf{W}^{(i)}, \mathbf{C}^{(i)}, \mathbf{X}^{(i)}) > 0$  and  $\mu^i$ and  $\overline{\mu}^i$  to denote estimated nuisances.

**Definition 12 (AC-gMTI estimators).** Let  $D_i$  denote samples following  $P_{\operatorname{rand}(Z_i)}(\mathbf{V})$  for  $i = 1, 2, \dots, m$ . For a fixed  $z_i \in \mathfrak{D}_{Z_i}$ , let  $D_{z_i}$  denote the subsamples of  $D_i$  such that  $Z_i = z_i$ . Let  $\mu^{m+1} \coloneqq Y$ . Let  $\mathbb{1}_{\mathbf{x}}^{i-1} \coloneqq \mathbb{1}_{\overline{\mathbf{x}}^{(i-1)}}(\overline{\mathbf{X}}^{(i-1)})$ . Then {REG, PW, DML} estimators are defined as:

$$T^{reg} \coloneqq \mathbb{E}_{D_{z_1}} \left[ \mu^2(A_1, \overline{x}_1)) \right],$$
  

$$T^{pw} \coloneqq \mathbb{E}_{D_{z_m}} \left[ \pi^{(m-1)}(\mathbf{A}^{(m-1)}, \mathbf{X}) \mathbb{1}_{\mathbf{X}}(\mathbf{X})Y \right],$$
  

$$T^{dml} \coloneqq \sum_{i=2}^m \mathbb{E}_{D_{z_i}} \left[ \pi^{(i-1)} \mathbb{1}_{\mathbf{x}}^{i-1} \{ \overline{\mu}^{i+1} - \mu^i \} \right] + \mathbb{E}_{D_{z_1}} \left[ \overline{\mu}^2 \right]$$

Assumption 11 ( $L_2$  consistency of nuisances). Estimated nuisances  $\{\mu^i\}_{i=2}^m$  and  $\{\pi^i\}_{i=1}^{m-1}$  are  $L_2$  consistent; specifically,

$$\begin{aligned} \|\mu^{i+1} - \mu_0^{i+1}\|_{P_{z_i}} &= o_{P_{z_i}}(1), \ \forall i \in \{1, 2, \cdots, m-1\} \\ \|\mu^i - \mu_0^i\|_{P_{z_i}} &= o_{P_{z_i}}(1), \ \forall i \in \{2, \cdots, m\} \\ \|\pi^i - \pi^i\|_{P_{z_{i+1}}} &= o_{P_{z_{i+1}}}(1), \ \forall i \in \{1, \cdots, m-1\}. \end{aligned}$$

**Lemma C.8** (Error analysis of the REG estimator for AC-gMTI). Suppose Assumptions (2,11) hold. Let  $T^{reg}$  denote the estimator defined in Def. 12. Then,

$$T^{reg} - \mathbb{E}\left[Y|do(\mathbf{x})\right] = R_1 + O_{P_{z_1}}(\left\|\mu^2 - \mu_0^2\right\|)$$

Proof of Lemma C.8. We first note that, by Theorem 7,

$$\mathbb{E}_{P_{z_1}}\left[\mu_0^2(A_1, \overline{x}_1)\right] = \mathbb{E}\left[Y|do(\mathbf{x})\right].$$

By Lemma C.3,

$$T^{reg} - \mathbb{E} \left[ Y | do(\mathbf{x}) \right] = T^{reg} - \mathbb{E}_{P_{z_1}} \left[ \mu_0^2(A_1, \overline{x}_1) \right]$$

$$= \underbrace{\mathbb{E}_{P_{z_1} - D_1} \left[ \mu_0^2(A_1, \overline{x}_1) \right] + \mathbb{E}_{P_{z_1} - D_1} \left[ \mu_0^2(A_1, \overline{x}_1) - \mu^2(A_1, \overline{x}_1) \right]}_{:= R_1}$$

$$+ \mathbb{E}_{P_{z_1}} \left[ \mu_0^2(A_1, \overline{x}_1) - \mu^2(A_1, \overline{x}_1) \right]$$

$$= R_1 + \mathbb{E}_{P_{z_1}} \left[ \mu_0^2(A_1, \overline{x}_1) - \mu^2(A_1, \overline{x}_1) \right]$$

$$= R_1 + O_{P_{z_1}} ( \left\| \mu^2 - \mu_0^2 \right\| ),$$

where  $R_1$  is a variable such that  $\sqrt{n_1}R_1$  converges in distribution to the normal random variable. The last equation holds by Cauchy-Schwartz inequality.

**Lemma C.9** (Error analysis of the PW estimator for AC-gMTI). Suppose Assumptions (2,11) hold. Let  $T^{pw}$  denote the estimator defined in Def. 12. Then,

$$T^{pw} - \mathbb{E}\left[Y|do(\mathbf{x})\right] = R_m + O_{P_{z_m}}(\|\pi^{(m-1)} - \pi_0^{(m-1)}\|)$$

**Proof of Lemma C.9.** In the proof, we will use  $\tilde{X}_i := \overline{X}_i$  for  $i = 1, 2, \dots, m-2$ , and  $\tilde{X}_{m-1} := \{\overline{X}_{m-1}, \overline{X}_m, \overline{X}_{m+1}\}$ . Therefore,  $\tilde{X}_i$  partitions **X**. In the proof, we tentatively assume

$$\mathbb{E}_{P_{z_m}}\left[\pi_0^{(m-1)}(\mathbf{V}^{(m-1)})\mathbb{1}_{\mathbf{x}}(\mathbf{X})Y\right] = \mathbb{E}\left[Y|do(\mathbf{x})\right].$$
(C.16)

\_

Then, by Lemma C.3,

$$T^{pw} - \mathbb{E} \left[ Y | do(\mathbf{x}) \right] = T^{pw} - \mathbb{E}_{P_{z_m}} \left[ \pi_0^{(m-1)} (\mathbf{V}^{(m-1)}) \mathbb{1}_{\mathbf{x}}(\mathbf{X}) Y \right]$$
  
=  $\underbrace{\mathbb{E}_{P_{z_m} - D_m} \left[ \pi_0^{(m-1)} (\mathbf{V}^{(m-1)}) \mathbb{1}_{\mathbf{x}}(\mathbf{X}) Y \right] + \mathbb{E}_{P_{z_m} - D_m} \left[ \left\{ \pi_0^{(m-1)} (\mathbf{V}^{(m-1)}) - \pi^{(m-1)} (\mathbf{V}^{(m-1)}) \right\} \mathbb{1}_{\mathbf{x}}(\mathbf{X}) Y \right]$   
:=  $R_m$   
+  $\mathbb{E}_{P_{z_m}} \left[ \left\{ \pi_0^{(m-1)} (\mathbf{V}^{(m-1)}) - \pi^{(m-1)} (\mathbf{V}^{(m-1)}) \right\} \mathbb{1}_{\mathbf{x}}(\mathbf{X}) Y \right]$   
=  $R_m + \mathbb{E}_{P_{x_m}} \left[ \left\{ \pi_0^{(m-1)} (\mathbf{V}^{(m-1)}) - \pi^{(m-1)} (\mathbf{V}^{(m-1)}) \right\} \mathbb{1}_{\mathbf{x}}(\mathbf{X}) Y \right]$   
=  $R_1 + O_{P_{z_m}} ( \left\| \pi^{(m-1)} - \pi_0^{(m-1)} \right\| ),$ 

where  $R_m$  is a variable converging in distribution to the normal distribution at  $\sqrt{n_m}$ -rate. The last equation holds by Cauchy-Schwartz inequality.

We now prove Eq. (C.16). We first show the following: For  $i = 2, \dots, m$ ,

$$\mathbb{E}\left[Y|do(\mathbf{x})\right] = \mathbb{E}_{P_{z_i}}\left[\prod_{j=1}^{i-1} \frac{P_{z_j}(A_j|\mathbf{A}^{(j-1)}, \tilde{\mathbf{X}}^{(j-1)})}{P_{z_{j+1}}(A_j, \tilde{X}_j|\mathbf{A}^{(j-1)}, \tilde{\mathbf{X}}^{(j-1)})} \mu_0^i(\mathbf{A}^{(i-1)}, \tilde{\mathbf{X}}^{(i-1)}) \mathbb{1}_{\tilde{\mathbf{x}}^{(i-1)}}(\tilde{\mathbf{X}}^{(i-1)})\right].$$
(C.17)

It holds for i = 2 as follow:

$$\mathbb{E}_{P_{z_2}}\left[\frac{P_{z_1}(A_1)}{P_{z_2}(A_1,\tilde{X}_1)}\mu_0^2(A_1,\tilde{X}_1)\mathbb{1}_{\tilde{x}_1}(\tilde{X}_1)\right] = \mathbb{E}_{P_{z_1}}\left[\mu_0^2(A_1,\tilde{x}_1)\right] = \mathbb{E}\left[Y|do(\mathbf{x})\right],$$

where the last equation holds by Lemma C.8. Now, we make a following induction hypothesis: For some  $i \in \{2, \dots, m-1\}$ , suppose

$$\mathbb{E}\left[Y|do(\mathbf{x})\right] \stackrel{\text{induction hypothesis}}{=} \mathbb{E}_{P_{z_{i-1}}}\left[\prod_{j=1}^{i-2} \frac{P_{z_j}(A_j|\mathbf{A}^{(j-1)}, \tilde{\mathbf{X}}^{(j-1)})}{P_{z_{j+1}}(A_j, \tilde{X}_j|\mathbf{A}^{(j-1)}, \tilde{\mathbf{X}}^{(j-1)})} \mu_0^{i-1}(\mathbf{A}^{(i-2)}, \tilde{\mathbf{X}}^{(i-2)})\mathbb{1}_{\tilde{\mathbf{X}}^{(i-2)}}(\tilde{\mathbf{X}}^{(i-2)})\right].$$

Then,

$$\begin{split} & \mathbb{E}_{P_{z_{i-1}}}\left[\prod_{j=1}^{i-2} \frac{P_{z_{j}}(A_{j}|\mathbf{A}^{(j-1)},\tilde{\mathbf{X}}^{(j-1)})}{P_{z_{j+1}}(A_{j},\tilde{X}_{j}|\mathbf{A}^{(j-1)},\tilde{\mathbf{X}}^{(j-1)})} \mu_{0}^{i-1}(\mathbf{A}^{(i-2)},\tilde{\mathbf{X}}^{(i-2)})\mathbb{1}_{\tilde{\mathbf{x}}^{(i-2)}}(\tilde{\mathbf{X}}^{(i-2)})\right] \\ &= \mathbb{E}_{P_{z_{i-1}}}\left[\prod_{j=1}^{i-2} \frac{P_{z_{j}}(A_{j}|\mathbf{A}^{(j-1)},\tilde{\mathbf{X}}^{(j-1)})}{P_{z_{j+1}}(A_{j},\tilde{X}_{j}|\mathbf{A}^{(j-1)},\tilde{\mathbf{X}}^{(j-1)})} \mathbb{E}_{P_{z_{i-1}}}\left[\mu_{0}^{i}(\mathbf{A}^{(i-1)},\tilde{\mathbf{x}}^{(i-2)})|\mathbf{A}^{(i-2)},\tilde{\mathbf{X}}^{(i-2)}\right]\mathbb{1}_{\tilde{\mathbf{x}}^{(i-2)}}(\tilde{\mathbf{X}}^{(i-2)})\right] \\ &= \mathbb{E}_{P_{z_{i-1}}}\left[\prod_{j=1}^{i-2} \frac{P_{z_{j}}(A_{j}|\mathbf{A}^{(j-1)},\tilde{\mathbf{X}}^{(j-1)})}{P_{z_{j+1}}(A_{j},\tilde{X}_{j}|\mathbf{A}^{(j-1)},\tilde{\mathbf{X}}^{(j-1)})} \frac{\mu_{0}^{i}(\mathbf{A}^{(i-1)},\tilde{\mathbf{x}}^{(i-2)})\mathbb{1}_{\tilde{\mathbf{x}}^{(i-2)}}(\tilde{\mathbf{X}}^{(i-2)})\right] \\ &= \mathbb{E}_{P_{z_{i-1}}}\left[\prod_{j=1}^{i-2} \frac{P_{z_{j}}(A_{j}|\mathbf{A}^{(j-1)},\tilde{\mathbf{X}}^{(j-1)})}{P_{z_{j+1}}(A_{j},\tilde{X}_{j}|\mathbf{A}^{(j-1)},\tilde{\mathbf{X}}^{(j-1)})} \frac{\mu_{z_{i-1}}(\tilde{X}_{i-1}|\mathbf{A}^{(i-1)},\tilde{\mathbf{X}}^{(i-2)})}{P_{z_{i-1}}(\tilde{X}_{i-1}|\mathbf{A}^{(i-2)},\tilde{\mathbf{X}}^{(i-2)})} \mu_{0}^{i}(\mathbf{A}^{(i-1)},\tilde{\mathbf{X}}^{(i-1)})\mathbb{1}_{\tilde{\mathbf{x}}^{(i-2)}}(\tilde{\mathbf{X}}^{(i-2)})\right] \\ &= \mathbb{E}_{P_{z_{i}}}\left[\prod_{j=1}^{i-2} \frac{P_{z_{j}}(A_{j}|\mathbf{A}^{(j-1)},\tilde{\mathbf{X}}^{(j-1)})}{P_{z_{j+1}}(A_{j},\tilde{X}_{j}|\mathbf{A}^{(j-1)},\tilde{\mathbf{X}}^{(j-1)})} \frac{P_{z_{i-1}}(A_{i-1}|\mathbf{A}^{(i-2)},\tilde{\mathbf{X}}^{(i-2)})}{P_{z_{i-1}}(\tilde{X}_{i-1}|\mathbf{A}^{(i-2)},\tilde{\mathbf{X}}^{(i-2)})} \mu_{0}^{i}(\mathbf{A}^{(i-1)},\tilde{\mathbf{X}}^{(i-1)})\mathbb{1}_{\tilde{\mathbf{x}}^{(i-1)}}(\tilde{\mathbf{X}}^{(i-1)})\right] \\ &= \mathbb{E}_{P_{z_{i}}}\left[\prod_{j=1}^{i-2} \frac{P_{z_{j}}(A_{j}|\mathbf{A}^{(j-1)},\tilde{\mathbf{X}}^{(j-1)})}{P_{z_{i-1}}(\tilde{X}_{i-1}|\mathbf{A}^{(i-2)},\tilde{\mathbf{X}}^{(i-2)})} \mu_{0}^{i}(\mathbf{A}^{(i-1)},\tilde{\mathbf{X}}^{(i-1)})\mathbb{1}_{\tilde{\mathbf{x}}^{(i-1)}}(\tilde{\mathbf{X}}^{(i-1)})\right] \\ &= \mathbb{E}_{P_{z_{i}}}\left[\prod_{j=1}^{i-1} \frac{P_{z_{j}}(A_{j}|\mathbf{A}^{(j-1)},\tilde{\mathbf{X}}^{(j-1)})}{P_{z_{i-1}}(\tilde{X}_{i-1})} \mu_{0}^{i}(\mathbf{A}^{(i-1)},\tilde{\mathbf{X}}^{(i-1)})} \mu_{0}^{i}(\mathbf{A}^{(i-1)},\tilde{\mathbf{X}}^{(i-1)})\right] \\ &= \mathbb{E}_{P_{z_{i}}}\left[\prod_{j=1}^{i-1} \frac{P_{z_{j}}(A_{j}|\mathbf{A}^{(j-1)},\tilde{\mathbf{X}}^{(j-1)})}{P_{z_{j+1}}(A_{j},\tilde{X}^{(j-1)})} \mu_{0}^{i}(\mathbf{A}^{(i-1)},\tilde{\mathbf{X}}^{(i-1)})} \mu_{0}^{i}(\mathbf{A}^{(i-1)})\right) \\ &= \mathbb{E}_{P_{z_{i}}}\left[\prod_{j=1}^{i-1} \frac{$$

where  $\stackrel{1}{=}$  holds by the law of total expectation. Therefore, Eq. (C.17) holds. By plugging i = m, we have

$$\begin{split} \mathbb{E}\left[Y|do(\mathbf{x})\right] &= \mathbb{E}_{P_{z_m}}\left[\prod_{j=1}^{m-1} \frac{P_{z_j}(A_j|\mathbf{A}^{(j-1)}, \tilde{\mathbf{X}}^{(j-1)})}{P_{z_m}(A_j, \tilde{X}_j|\mathbf{A}^{(j-1)}, \tilde{\mathbf{X}}^{(j-1)})} \mu_0^m(\mathbf{A}^{(m-1)}, \tilde{\mathbf{X}}^{(m-1)}) \mathbb{1}_{\tilde{\mathbf{x}}^{(m-1)}}(\tilde{\mathbf{X}}^{(m-1)})\right] \\ &= \mathbb{E}_{P_{z_m}}\left[\pi_0^{m-1}(\mathbf{A}^{(m-1)}, \tilde{\mathbf{X}}^{(m-1)}) \mu_0^m(\mathbf{A}^{(m-1)}, \tilde{\mathbf{X}}^{(m-1)}) \mathbb{1}_{\mathbf{x}}(\mathbf{X})\right] \\ &= \mathbb{E}_{P_{z_m}}\left[\pi_0^{m-1}\mathbf{A}^{(m-1)}, \tilde{\mathbf{X}}^{(m-1)}) \mathbb{E}_{z_m}\left[Y|\mathbf{A}^{(m-1)}, \tilde{\mathbf{X}}^{(m-1)}\right] \mathbb{1}_{\mathbf{x}}(\mathbf{X})\right] \\ &= \mathbb{E}_{P_{z_m}}\left[\pi_0^{m-1}(\mathbf{A}^{(m-1)}, \tilde{\mathbf{X}}^{(m-1)}) \mathbb{1}_{\mathbf{x}}(\mathbf{X})Y\right]. \end{split}$$

Lemma C.10 (Bias Analysis of the DML estimator for AC-gMTI). Suppose Assumptions (2,8) hold. Let  $\mu^{m+1} \coloneqq Y$ . Let  $\tilde{X}_i \coloneqq \overline{X}_i$  for  $i = 1, 2, \dots, m-2$ , and  $\tilde{X}_{m-1} \coloneqq \{\overline{X}_{m-1}, \overline{X}_m, \overline{X}_{m+1}\}$ . Let  $V_i \coloneqq \{A_i, \tilde{X}_i\}$  for  $i = 1, 2, \dots, m-1$ . Let  $B_i \coloneqq \{A_i, \tilde{X}_{i-1}\}$  for  $i = 1, 2, \dots, m$  where  $\tilde{X}_0 \coloneqq \emptyset$ . Let  $T^{dml}(\{\pi^k\}_{k=1}^{m-1}, \{\mu^k\}_{k=2}^m)$  be defined as follow:

$$T^{dml}(\{\pi^k\}_{k=1}^{m-1}, \{\mu^k\}_{k=2}^m) \approx \sum_{i=2}^m \mathbb{E}_{P_{z_i}}\left[\pi^{(i-1)}(\mathbf{V}^{(i-1)})\mathbb{1}_{\mathbf{\tilde{x}}^{(i-1)}}(\mathbf{\tilde{X}}^{(i-1)})\left\{\mu^{i+1}(\mathbf{B}^{(i)}, \tilde{x}_i) - \mu^i(\mathbf{B}^{(i-1)}, \tilde{X}_{i-1})\right\}\right] + \mathbb{E}_{P_{z_1}}\left[\mu^2(\mathbf{B}^{(1)}, \tilde{x}_1)\right].$$
(C.18)

Then,

$$T^{dml}(\{\pi^k\}_{k=1}^{m-1}, \{\mu^k\}_{k=2}^m) - \mathbb{E}\left[Y|do(\mathbf{x})\right] = \sum_{i=2}^m O_{P_{z_i}}\left(\left\|\mu^i - \mu_0^i\right\| \left\|\pi^{i-1} - \pi_0^{i-1}\right\|\right)$$
(C.19)

**Proof of Lemma C.10.** We follow the proof technique used in (Rotnitzky et al., 2017). To simplify the notation, we sometimes simply denote  $\mu^i(\mathbf{B}^{(i-1)}, \tilde{X}_{i-1})$  as  $\mu^i$ ;  $\mu^i(\mathbf{B}^{(i-1)}, \tilde{x}_{i-1})$  as  $\overline{\mu^i}$ ; and  $\pi^i(\mathbf{V}^{(i)})$  as  $\pi^i$ .

We first note that

$$T^{dml}(\{\pi_0^k\}_{k=1}^{m-1}, \{\mu_0^k\}_{k=2}^m) = \mathbb{E}\left[Y|do(\mathbf{x})\right].$$
(C.20)

It's easy to witness Eq. (C.20) because, for  $i = 2, 3, \dots, m$ ,

$$\begin{split} & \mathbb{E}_{P_{z_i}} \left[ \pi_0^{(i-1)}(\mathbf{V}^{(i-1)}) \mathbb{1}_{\tilde{\mathbf{X}}^{(i-1)}}(\tilde{\mathbf{X}}^{(i-1)}) \left\{ \mu_0^{i+1}(\mathbf{B}^{(i)}, \tilde{x}_i) - \mu_0^i(\mathbf{B}^{(i-1)}, \tilde{X}_{i-1}) \right\} \right] \\ & = \mathbb{E}_{P_{z_i}} \left[ \mathbb{E}_{P_{z_i}} \left[ \pi_0^{(i-1)}(\mathbf{V}^{(i-1)}) \mathbb{1}_{\tilde{\mathbf{X}}^{(i-1)}}(\tilde{\mathbf{X}}^{(i-1)}) \left\{ \mu_0^{i+1}(\mathbf{B}^{(i)}, \tilde{x}_i) - \mu_0^i(\mathbf{B}^{(i-1)}, \tilde{X}_{i-1}) \right\} \left| \mathbf{B}^{(i-1)}, \tilde{X}_{i-1} \right] \right] \\ & = \mathbb{E}_{P_{z_i}} \left[ \pi_0^{(i-1)}(\mathbf{V}^{(i-1)}) \mathbb{1}_{\tilde{\mathbf{X}}^{(i-1)}}(\tilde{\mathbf{X}}^{(i-1)}) \left\{ \mathbb{E}_{P_{z_i}} \left[ \mu_0^{i+1}(\mathbf{B}^{(i)}, \tilde{x}_i) | \mathbf{B}^{(i-1)}, \tilde{X}_{i-1} \right] - \mu_0^i(\mathbf{B}^{(i-1)}, X_{i-1}) \right\} \right] \\ & = \mathbb{E}_{P_{z_i}} \left[ \pi_0^{(i-1)}(\mathbf{V}^{(i-1)}) \mathbb{1}_{\tilde{\mathbf{X}}^{(i-1)}}(\tilde{\mathbf{X}}^{(i-1)}) \left\{ \mu_0^i(\mathbf{B}^{(i-1)}, \tilde{X}_{i-1}) - \mu_0^i(\mathbf{B}^{(i-1)}, \tilde{X}_{i-1}) \right\} \right] \\ & = 0, \end{split}$$

where the equation  $\stackrel{1}{=}$  holds by the law of total expectation. Therefore,

$$T^{dml}(\{\pi_0^k\}_{k=1}^{m-1}, \{\mu_0^k\}_{k=2}^m) = \mathbb{E}_{P_{x_1}}\left[\mu_0^2(\mathbf{B}^{(1)}, \overline{x}_1)\right] = \mathbb{E}\left[Y|do(\mathbf{x})\right],$$

where the second equation holds by Lemma C.8. Therefore, it suffices to prove the following to show Eq. (C.19):

$$T^{dml}(\{\pi^k\}_{k=1}^{m-1}, \{\mu^k\}_{k=2}^m) - T^{dml}(\{\pi_0^k\}_{k=1}^{m-1}, \{\mu_0^k\}_{k=2}^m) = \sum_{i=2}^m O_{P_{z_i}}\left(\left\|\mu^i - \mu_0^i\right\| \left\|\pi^{i-1} - \pi_0^{i-1}\right\|\right).$$
(C.21)

For  $i = 1, 2, \cdots, m - 1$ , we define a quantity

$$\omega_0^i(\mathbf{B}^{(i)}) \coloneqq \frac{P_{z_i}(\mathbf{B}^{(i)})}{P_{z_m}(\mathbf{B}^{(i)})}.$$

We note that  $\omega_0^i(\mathbf{B}^{(i)})$  is related with  $\pi$  as follow:

$$\omega_0^i(\mathbf{B}^{(i)}) = \pi_0^i(\mathbf{V}^{(i)}) P_{z_m}(\tilde{X}_i | \mathbf{B}^{(i)}).$$
(C.22)

To witness, consider the following:

$$\begin{split} \omega_0^i(\mathbf{B}^{(i)}) &= \frac{P_{z_i}(A_i | \mathbf{V}^{(i-1)}) P_{z_i}(\mathbf{V}^{(i-1)})}{P_{z_m}(A_i | \mathbf{V}^{(i-1)}) P_{z_m}(\mathbf{V}^{(i-1)})} \\ &\stackrel{2}{=} \frac{P_{z_i}(A_i | \mathbf{V}^{(i-1)})}{P_{z_m}(A_i | \mathbf{V}^{(i-1)})} \\ &= \pi_0^i(\mathbf{V}^{(i)}) \frac{P_{z_m}(A_i, \tilde{X}_i | \mathbf{V}^{(i-1)})}{P_{z_m}(A_i | \mathbf{V}^{(i-1)})} \\ &= \pi_0^i(\mathbf{V}^{(i)}) P_{z_m}(\tilde{X}_i | \mathbf{B}^{(i)}), \end{split}$$

where

•  $\stackrel{2}{=}$  holds since  $\tilde{X}_i$  is non-descendent to  $\mathbf{V}^{(i-1)}$ , so that  $P_{z_i}(\mathbf{V}^{(i-1)}) = P_{z_m}(\mathbf{V}^{(i-1)})$ .

To simplify the notation, we sometimes simply denote  $\omega_0^i(\mathbf{B}^{(i)})$  as  $\omega_0^i$ ;  $\mu^i(\mathbf{B}^{(i-1)}, X_{i-1})$  as  $\mu^i$ ;  $\mu^i(\mathbf{B}^{(i-1)}, \tilde{x}_{i-1})$  as  $\overline{\mu}^i$ ; and  $\pi^i(\mathbf{V}^{(i)})$  as  $\pi^i$ .

Then,  $T^{dml}(\{\pi^k\}_{k=1}^{m-1}, \{\mu^k\}_{k=2}^m)$  in Eq. (C.18) can be rewritten as

$$T^{dml}(\{\pi^k\}_{k=1}^{m-1}, \{\mu^k\}_{k=2}^m) = \sum_{i=2}^m \mathbb{E}_{P_{z_m}} \left[ \omega_0^i \pi^{(i-1)} \mathbb{1}_{\tilde{\mathbf{x}}^{(i-1)}}(\tilde{\mathbf{X}}^{(i-1)}) \left\{ \overline{\mu}^{i+1} - \mu^i \right\} + \omega_0^1 \overline{\mu}^2 \right],$$
(C.23)

where  $\overline{\mu}^{m+1} \coloneqq Y$ .

For each  $k = 1, 2, \dots, m$ , we define a quantity  $Q_k$  as follow:

$$Q_{k} \coloneqq Q_{k} (\{\pi^{j}\}_{j=k}^{m-1}, \{\mu^{j}\}_{j=k+1}^{m}) \coloneqq \omega_{0}^{k} \overline{\mu}^{k+1} + \sum_{i=k+1}^{m} \omega_{0}^{i} \pi^{(k:i-1)} \mathbb{1}_{\tilde{\mathbf{X}}^{(k:i-1)}} (\tilde{\mathbf{X}}^{(k:i-1)}) \{\overline{\mu}^{i+1} - \mu^{i}\}.$$
(C.24)

Note  $Q_m = Y$  and  $\mathbb{E}_{P_{x_m}}[Q_1] = T^{dml}(\{\pi^k\}_{k=1}^{m-1}, \{\mu^k\}_{k=2}^m)$  defined in Eq. (C.23). We note that

$$\mathbb{E}_{P_{x_m}} \left[ Q_1 - \omega_0^1 \overline{\mu}_0^2 \right] \stackrel{4}{=} T^{dml}(\{\pi^k\}_{k=1}^{m-1}, \{\mu^k\}_{k=2}^m) - \mathbb{E}_{P_{z_1}} \left[ \mu_0^2(\mathbf{B}^{(1)}, \tilde{x}_1) \right]$$
  
$$\stackrel{5}{=} T^{dml}(\{\pi^k\}_{k=1}^{m-1}, \{\mu^k\}_{k=2}^m) - \mathbb{E}\left[ Y | do(\mathbf{x}) \right]$$
  
$$= \text{l.h.s. of Eq. (C.19),}$$

where

- $\stackrel{4}{=}$  holds since  $\mathbb{E}_{P_{z_m}}\left[\omega_0^1(\mathbf{B}^{(1)})\mu^2(\mathbf{B}^{(1)},\tilde{x}_1)\right] = \mathbb{E}_{P_{z_1}}\left[\mu^1(\mathbf{B}^{(1)},\tilde{x}_1)\right].$
- $\stackrel{5}{=}$  holds by Lemma C.8.

Motivating from the fact that  $\mathbb{E}_{P_{z_m}}\left[Q_1 - \omega_0^1 \overline{\mu}_0^2\right] = 1.$ h.s. of Eq. (C.19), we establish a following induction hypothesis. For  $\overline{P}_{z_m}^{i-1} \coloneqq P_{z_m}(\cdot | \mathbf{V}^{(i-1)})$ , the induction hypothesis is given as follow:

Hypothesis: 
$$\mathbb{E}_{\overline{P}_{z_m}^{k-1}}\left[Q_k - \omega_0^k \overline{\mu}_0^{k+1}\right] = \sum_{i=k+1}^m O_{\overline{P}_{z_i}^{k-1}}\left(\left\|\mu^i - \mu_0^i\right\| \left\|\pi^{i-1} - \pi_0^{i-1}\right\|\right), \text{ for } k \in \{2, \cdots, m-1\}$$
 (C.25)

We first verify the hypothesis Eq. (C.25) for k = m - 1.

$$\begin{split} & \mathbb{E}_{\overline{P}_{z_m}^{m-2}} \left[ Q_{m-1} - \omega_0^{m-1} \overline{\mu}_0^m \right] \\ &= \mathbb{E}_{\overline{P}_{z_m}^{m-2}} \left[ \omega_0^{m-1} \overline{\mu}^m + \pi^{m-1} \mathbb{1}_{\tilde{x}_{m-1}} (\tilde{X}_{m-1}) \left\{ Y - \mu^m \right\} - \omega_0^{m-1} \overline{\mu}_0^m \right] \\ & \stackrel{6}{=} \mathbb{E}_{\overline{P}_{z_m}^{m-2}} \left[ \omega_0^{m-1} \overline{\mu}^m + \pi^{m-1} \mathbb{1}_{\tilde{x}_{m-1}} (\tilde{X}_{m-1}) \left\{ \mu_0^m - \mu^m \right\} - \omega_0^{m-1} \overline{\mu}_0^m \right] \\ &= \mathbb{E}_{\overline{P}_{z_m}^{m-2}} \left[ \omega_0^{m-1} \left\{ \overline{\mu}^m - \overline{\mu}_0^m \right\} + \pi^{m-1} \mathbb{1}_{\tilde{x}_{m-1}} (\tilde{X}_{m-1}) \left\{ \mu_0^m - \mu^m \right\} \right] \\ & \stackrel{7}{=} \mathbb{E}_{\overline{P}_{z_m}^{m-2}} \left[ \omega_0^{m-1} \frac{\mathbb{1}_{\tilde{x}_{m-1}} (\tilde{X}_{m-1})}{P_{z_m} (\tilde{X}_{m-1} | \mathbf{B}^{(m-1)})} \left\{ \mu^m - \mu_0^m \right\} + \pi^{m-1} \mathbb{1}_{\tilde{x}_{m-1}} (\tilde{X}_{m-1}) \left\{ \mu_0^m - \mu^m \right\} \right] \\ & \stackrel{8}{=} \mathbb{E}_{\overline{P}_{z_m}^{m-2}} \left[ \pi_0^{m-1} \mathbb{1}_{\tilde{x}_{m-1}} (\tilde{X}_{m-1}) \left\{ \mu^m - \mu_0^m \right\} + \pi^{m-1} \mathbb{1}_{\tilde{x}_{m-1}} (\tilde{X}_{m-1}) \left\{ \mu_0^m - \mu^m \right\} \right] \\ &= \mathbb{E}_{\overline{P}_{z_m}^{m-2}} \left[ \mathbb{1}_{\tilde{x}_{m-1}} (\tilde{X}_{m-1}) \left\{ \mu^m - \mu_0^m \right\} \left\{ \pi_0^{m-1} - \pi^{m-1} \right\} \right] \\ &= \mathbb{E}_{P_{x_m}} \left[ \mathbb{1}_{\tilde{x}_{m-1}} (\tilde{X}_{m-1}) \left\{ \mu^m - \mu_0^m \right\} \left\{ \pi_0^{m-1} - \pi^{m-1} \right\} \right] \\ &= \mathbb{E}_{\overline{P}_{z_m}^{m-2}} \left( \left\| \mu^m - \mu_0^m \right\| \left\| \pi^{m-1} - \pi_0^{m-1} \right\| \right), \end{split}$$

where

- $\stackrel{6}{=}$  holds by the total law of expectation.
- $\stackrel{7}{=}$  holds since

$$\mathbb{E}_{P_{z_m}}\left[\mu^{m-1}(\mathbf{B}^{(m-1)}, \tilde{x}_{m-1}) \middle| \mathbf{V}^{(m-2)}\right] = \mathbb{E}_{P_{z_m}}\left[\mu^{m-1}(\mathbf{B}^{(m-1)}, \tilde{X}_{m-1}) \frac{\mathbb{1}_{\tilde{x}_{m-1}}(\tilde{X}_{m-1})}{P_{z_m}(\tilde{X}_{m-1}|\mathbf{B}^{(m-1)})} \middle| \mathbf{V}^{(m-2)}\right].$$

- $\stackrel{8}{=}$  holds by the definition of  $\omega_0^{m-1}$ .
- $\stackrel{9}{=}$  holds by applying Cauchy-Schwartz inequality.

Now, we suppose Eq. (C.25) holds for some  $k + 1 \in \{2, \dots, m-1\}$ . Then, we will show that Eq. (C.25) holds for k. Toward this end, we first rewrite  $Q_k$  in Eq. (C.24) in a recursive form. For any  $k + 1 \in \{2, \dots, m-1\}$ , the following relation can be derived from Eq. (C.24):

$$\pi^{k} \mathbb{1}_{\tilde{x}_{k}}(\tilde{X}_{k}) \left\{ Q_{k+1} - \omega_{0}^{k+1} \overline{\mu}^{k+2} \right\} = \sum_{i=k+2}^{m} \omega_{0}^{i} \pi^{(k:i-1)} \mathbb{1}_{\tilde{\mathbf{x}}^{(k:i-1)}}(\tilde{\mathbf{X}}^{(k:i-1)}) \left\{ \overline{\mu}^{i+1} - \mu^{i} \right\}.$$

Therefore, for each  $k = 1, 2, \cdots, m - 1$ ,

$$Q_{k}(\{\pi^{j}\}_{j=k}^{m-1},\{\mu^{j}\}_{j=k+1}^{m}) = \omega_{0}^{k}\overline{\mu}^{k+1} + \omega_{0}^{k+1}\pi^{k}\mathbb{1}_{\tilde{x}_{k}}(\tilde{X}_{k})\left\{\overline{\mu}^{k+2} - \mu^{k+1}\right\} + \pi^{k}\mathbb{1}_{\tilde{x}_{k}}(\tilde{X}_{k})\left\{Q_{k+1} - \omega_{0}^{k+1}\overline{\mu}^{k+2}\right\}.$$

Then,

$$\begin{split} &\mathbb{E}_{\overline{P}^{k-1}_{z_m}}\left[Q_k - \omega_0^k \overline{\mu}^{k+1}\right] \\ &= \mathbb{E}_{\overline{P}^{k-1}_{z_m}}\left[\omega_0^k \overline{\mu}^{k+1} + \omega_0^{k+1} \pi^k \mathbb{1}_{\bar{x}_k}(\bar{X}_k) \left\{\overline{\mu}^{k+2} - \mu^{k+1}\right\} + \pi^k \mathbb{1}_{\bar{x}_k}(\bar{X}_k) \left\{Q_{k+1} - \omega_0^{k+1} \overline{\mu}^{k+2}\right\} - \omega_0^k \overline{\mu}^{k+1}_0\right] \\ &= \mathbb{E}_{\overline{P}^{k-1}_{z_m}}\left[\omega_0^k \overline{\mu}^{k+1} + \omega_0^{k+1} \pi^k \mathbb{1}_{\bar{x}_k}(\bar{X}_k) \left\{\overline{\mu}^{k+2} - \mu^{k+1}\right\} + \omega_0^{k+1} \pi^k \mathbb{1}_{\bar{x}_k}(\bar{X}_k) \left\{\overline{\mu}^{k+2} - \overline{\mu}^{k+2}\right\} - \omega_0^k \overline{\mu}^{k+1}_0\right] \\ &+ \mathbb{E}_{\overline{P}^{k-1}_{z_m}}\left[\pi^k \mathbb{1}_{\bar{x}_k}(\bar{X}_k) \left\{Q_{k+1} - \omega_0^{k+1} \overline{\mu}^{k+2}_0\right\}\right] \\ &= \mathbb{E}_{\overline{P}^{k-1}_{z_m}}\left[\omega_0^k \overline{\mu}^{k+1} + \omega_0^{k+1} \pi^k \mathbb{1}_{\bar{x}_k}(\bar{X}_k) \left\{\overline{\mu}^{k+2} - \mu^{k+1}\right\} - \omega_0^k \overline{\mu}^{k+1}_0\right] \\ &= \mathbb{E}_{\overline{P}^{k-1}_{z_m}}\left[\omega_0^k \left\{\overline{\mu}^{k+1} - \overline{\mu}^{k+1}_0\right\} + \omega_0^{k+1} \pi^k \mathbb{1}_{\bar{x}_k}(\bar{X}_k) \left\{\overline{\mu}^{k+2} - \mu^{k+1}\right\}\right] \\ &= \mathbb{E}_{\overline{P}^{k-1}_{z_m}}\left[\omega_0^k \left\{\overline{\mu}^{k+1} - \overline{\mu}^{k+1}_0\right\} + \pi^k \mathbb{1}_{\bar{x}_k}(\bar{X}_k) \left\{\overline{\mu}^{k+2} - \mu^{k+1}\right\}\right] \\ &= \mathbb{E}_{\overline{P}^{k-1}_{z_m}}\left[\omega_0^k \left\{\overline{\mu}^{k+1} - \overline{\mu}^{k+1}_0\right\} + \pi^k \mathbb{1}_{\bar{x}_k}(\bar{X}_k) \left\{\mu_0^{k+1} - \mu^{k+1}_k\right\}\right] \\ &= \mathbb{E}_{\overline{P}^{k-1}_{z_m}}\left[\omega_0^k \left\{\overline{\mu}^{k+1} - \overline{\mu}^{k+1}_0\right\} + \pi^k \mathbb{1}_{\bar{x}_k}(\bar{X}_k) \left\{\mu_0^{k+1} - \mu^{k+1}_k\right\right\}\right] \\ &= \mathbb{E}_{\overline{P}^{k-1}_{z_m}}\left[\pi^k \mathbb{1}_{\bar{x}_k}(\bar{X}_k) \left\{\mu^{k+1} - \mu^{k+1}_0\right\} + \pi^k \mathbb{1}_{\bar{x}_k}(\bar{X}_k) \left\{\mu_0^{k+1} - \mu^{k+1}_0\right\}\right] \\ &= \mathbb{E}_{\overline{P}^{k-1}_{z_m}}\left[\pi^k \mathbb{1}_{\bar{x}_k}(\bar{X}_k) \left\{\mu^{k+1} - \mu^{k+1}_0\right\} + \pi^k \mathbb{1}_{\bar{x}_k}(\bar{X}_k) \left\{\mu_0^{k+1} - \mu^{k+1}_0\right\}\right] \\ &= \mathbb{E}_{\overline{P}^{k-1}_{z_m}}\left[\omega_0^{k+1} \mathbb{1}_{\bar{x}_k}(\bar{X}_k) \left\{\mu^{k+1} - \mu^{k+1}_0\right\} + \left\{\pi_0^k - \pi^k\right\}\right] \\ &= \mathbb{E}_{\overline{P}^{k-1}_{z_m}}\left[\omega_0^{k+1} \mathbb{1}_{\bar{x}_k}(\bar{X}_k) \left\{\mu^{k+1} - \mu^{k+1}_0\right\}\right\} \\ &= \mathbb{E}_{\overline{P}^{k-1}_{z_m}}\left[\mathbb{1}_{\bar{x}_k}(\bar{X}_k) \left\{\mu^{k+1} - \mu^{k+1}_0\right\} + \left\{\pi_0^k - \pi^k\right\}\right] \\ &= \mathbb{E}_{\overline{P}^{k-1}_{z_m}}\left[\omega_0^{k+1} \mathbb{1}_{\bar{x}_k}(\bar{X}_k) \left\{\mu^{k+1} - \mu^{k+1}_0\right\}\right\} \\ &= \mathbb{E}_{\overline{P}^{k-1}_{z_m}}\left[\mathbb{1}_{\bar{x}_k}(\bar{X}_k) \left\{\mu^{k+1} - \mu^{k+1}_0\right\} + \left\{\pi_0^k - \pi^k\right\}\right] \\ &= \mathbb{E}_{\overline{P}^{k-1}_{z_m}}\left[\mathbb{1}_{\bar{x}_k}(\bar{X}_k) \left\{\mu^{k+1} - \mu^{k+1}_0\right\}\right\} \\ &= \mathbb{E}_{\overline{P}^{k-1}_{z_m}}\left[\mathbb{1}_{\bar{x}_k}(\bar{X$$

where

•  $\stackrel{10}{=}$  holds since

$$\begin{split} & \mathbb{E}_{P_{z_m}} \left[ \omega_0^{k+1}(\mathbf{B}^{(k+1)}) \pi^k(\mathbf{V}^{(k)}) \mathbb{1}_{\tilde{x}_k}(\tilde{X}_k) \mu_0^{k+2}(\mathbf{B}^{(k+1)}, x_{k+1}) \Big| \mathbf{V}^{(k-1)} \right] \\ &= \mathbb{E}_{P_{z_m}} \left[ \mathbb{E}_{P_{z_m}} \left[ \omega_0^{k+1}(\mathbf{B}^{(k+1)}) \pi^k(\mathbf{V}^{(k)}) \mathbb{1}_{\tilde{x}_k}(\tilde{X}_k) \mu_0^{k+2}(\mathbf{B}^{(k+1)}, \tilde{x}_{k+1}) \Big| \mathbf{V}^{(k)} \right] \Big| \mathbf{V}^{(k-1)} \right] \\ &= \mathbb{E}_{P_{z_m}} \left[ \pi^k(\mathbf{V}^{(k)}) \mathbb{1}_{\tilde{x}_k}(\tilde{X}_k) \mathbb{E}_{P_{z_m}} \left[ \omega_0^{k+1}(\mathbf{B}^{(k+1)}) \mu_0^{k+2}(\mathbf{B}^{(k+1)}, \tilde{x}_{k+1}) \Big| \mathbf{V}^{(k)} \right] \Big| \mathbf{V}^{(k-1)} \right] \\ &= \mathbb{E}_{P_{z_m}} \left[ \pi^k(\mathbf{V}^{(k)}) \mathbb{1}_{\tilde{x}_k}(\tilde{X}_k) \mathbb{E}_{P_{z_m}} \left[ \frac{P_{z_{k+1}}(\mathbf{B}^{(k+1)})}{P_{z_m}(\mathbf{B}^{(k+1)})} \mu_0^{k+2}(\mathbf{B}^{(k+1)}, \tilde{x}_{k+1}) \Big| \mathbf{V}^{(k)} \right] \Big| \mathbf{V}^{(k-1)} \right] \\ &= \mathbb{E}_{P_{z_m}} \left[ \pi^k(\mathbf{V}^{(k)}) \mathbb{1}_{\tilde{x}_k}(\tilde{X}_k) \mathbb{E}_{P_{z_{k+1}}} \left[ \mu_0^{k+2}(\mathbf{B}^{(k+1)}, \tilde{x}_{k+1}) \Big| \mathbf{V}^{(k)} \right] \Big| \mathbf{V}^{(k-1)} \right] \\ &= \mathbb{E}_{P_{z_m}} \left[ \pi^k(\mathbf{V}^{(k)}) \mathbb{1}_{\tilde{x}_k}(\tilde{X}_k) \mu_0^{k+1}(\mathbf{V}^{(k)}) \Big| \mathbf{V}^{(k-1)} \right]. \end{split}$$

•  $\stackrel{11}{=}$  holds since

$$\begin{split} & \mathbb{E}_{P_{z_m}} \left[ \omega_0^k(\mathbf{B}^{(k)}) \left\{ \mu^{k+1}(\mathbf{B}^{(k)}, \tilde{x}_k) - \mu_0^{k+1}(\mathbf{B}^{(k)}, \tilde{x}_k) \right\} \left| \mathbf{V}^{(k-1)} \right] \\ &= \mathbb{E}_{P_{z_m}} \left[ \omega_0^k(\mathbf{B}^{(k)}) \frac{\mathbb{1}_{\tilde{x}_k}(\tilde{X}_k)}{P_{z_m}(\tilde{X}_k | \mathbf{B}^{(k)})} \left\{ \mu^{k+1}(\mathbf{B}^{(k)}, \tilde{X}_k) - \mu_0^{k+1}(\mathbf{B}^{(k)}, \tilde{X}_k) \right\} \left| \mathbf{V}^{(k-1)} \right] \right] \\ &= \mathbb{E}_{P_{z_m}} \left[ \pi_0^k(\mathbf{V}^{(k)}) P_{z_m}(\tilde{X}_k | \mathbf{B}^{(k)}) \frac{\mathbb{1}_{\tilde{x}_k}(\tilde{X}_k)}{P_{z_m}(\tilde{X}_k | \mathbf{B}^{(k)})} \left\{ \mu^{k+1}(\mathbf{B}^{(k)}, \tilde{X}_k) - \mu_0^{k+1}(\mathbf{B}^{(k)}, \tilde{X}_k) \right\} \left| \mathbf{V}^{(k-1)} \right] \right] \\ &= \mathbb{E}_{P_{z_m}} \left[ \pi_0^k(\mathbf{V}^{(k)}) \mathbb{1}_{\tilde{x}_k}(\tilde{X}_k) \left\{ \mu^{k+1}(\mathbf{B}^{(k)}, \tilde{X}_k) - \mu_0^{k+1}(\mathbf{B}^{(k)}, \tilde{X}_k) \right\} \left| \mathbf{V}^{(k-1)} \right]. \end{split}$$

•  $\stackrel{12}{=}$  and  $\stackrel{13}{=}$  hold since

$$\begin{split} & \mathbb{E}_{P_{z_m}} \left[ \mathbb{1}_{\tilde{x}_k}(\tilde{X}_k) \left\{ \mu^{k+1}(\mathbf{V}^{(k)}) - \mu_0^{k+1}(\mathbf{V}^{(k)}) \right\} \left\{ \pi_0^k(\mathbf{V}^{(k)}) - \pi^k(\mathbf{V}^{(k)}) \right\} \left| \mathbf{V}^{(k-1)} \right] \\ &= \mathbb{E}_{P_{z_m}} \left[ \omega_0^{k+1}(\mathbf{V}^{(k)}) \mathbb{1}_{\tilde{x}_k}(\tilde{X}_k) \left\{ \mu^{k+1}(\mathbf{V}^{(k)}) - \mu_0^{k+1}(\mathbf{V}^{(k)}) \right\} \left\{ \pi_0^k(\mathbf{V}^{(k)}) - \pi^k(\mathbf{V}^{(k)}) \right\} \left| \mathbf{V}^{(k-1)} \right] \right] \\ &= \mathbb{E}_{P_{z_m}} \left[ \frac{P_{z_{k+1}}(\mathbf{V}^{(k)})}{P_{z_m}(\mathbf{V}^{(k)})} \mathbb{1}_{\tilde{x}_k}(\tilde{X}_k) \left\{ \mu^{k+1}(\mathbf{V}^{(k)}) - \mu_0^{k+1}(\mathbf{V}^{(k)}) \right\} \left\{ \pi_0^k(\mathbf{V}^{(k)}) - \pi^k(\mathbf{V}^{(k)}) \right\} \left| \mathbf{V}^{(k-1)} \right] \right] \\ &= \mathbb{E}_{P_{z_{k+1}}} \left[ \mathbb{1}_{\tilde{x}_k}(\tilde{X}_k) \left\{ \mu^{k+1}(\mathbf{V}^{(k)}) - \mu_0^{k+1}(\mathbf{V}^{(k)}) \right\} \left\{ \pi_0^k(\mathbf{V}^{(k)}) - \pi^k(\mathbf{V}^{(k)}) \right\} \left| \mathbf{V}^{(k-1)} \right], \end{split}$$

where the second equation hold since

$$\omega_0^{k+1}(\mathbf{V}^{(k)}) = \frac{P_{z_{k+1}}(\mathbf{V}^{(k)})}{P_{z_m}(\mathbf{V}^{(k)})} = 1$$

since  $Z_{k+1}, Z_m$  are non-descendants of  $\mathbf{V}^{(k)}$  so that  $P_{z_{k+1}}(\mathbf{V}^{(k)}) = P_{z_m}(\mathbf{V}^{(k)})$ .

- $\stackrel{14}{=}$  holds by the induction hypothesis.
- $\stackrel{15}{=}$  holds by Cauchy-Schwartz inequality.

Therefore, the induction hypothesis in Eq. (C.25) holds for all  $k = 1, 2, \dots, m-1$ . Therefore,

l.h.s. of Eq. (C.19) = 
$$\mathbb{E}_{P_{z_m}} \left[ Q_1 - \omega_0^1 \overline{\mu}_0^2 \right] = \sum_{i=2}^m O_{P_{z_i}} \left( \left\| \mu^i - \mu_0^i \right\| \left\| \pi^{i-1} - \pi_0^{i-1} \right\| \right),$$

where the second equation holds by plugging k = 1 into the verified hypothesis in Eq. (C.25). This completes the proof.

**Lemma C.11** (Error analysis of the DML estimator for AC-gMTI). Suppose Assumptions (2,11) hold. Let  $T^{dml}$  denote the estimator defined in Def. 12. Then,

$$T^{dml} - \mathbb{E}\left[Y|do(\mathbf{x})\right] = \sum_{i=1}^{m} R_i + \sum_{i=2}^{m} O_{P_{z_i}}\left(\left\|\mu^i - \mu_0^i\right\| \left\|\pi^{i-1} - \pi_0^{i-1}\right\|\right)\right)$$

where  $R_i$  for  $i = 1, 2, \dots, m$  are variables converging in mean-zero normal distribution at  $n_i^{-1/2}$  rates.

**Proof of Lemma C.11.** In the proof, we will use  $\tilde{X}_i := \overline{X}_i$  for  $i = 1, 2, \dots, m-2$ , and  $\tilde{X}_{m-1} := \{\overline{X}_{m-1}, \overline{X}_m, \overline{X}_{m+1}\}$ . Therefore,  $\tilde{X}_i$  partitions **X**. We will use  $B_i := \{A_i, \tilde{X}_{i-1}\}$  or all  $i = 1, 2, \dots, m$ , where  $X_0 := \emptyset$ . To simplify the notation, we sometimes simply denote  $\mu^i(\mathbf{B}^{(i-1)}, \tilde{X}_{i-1})$  as  $\mu^i$ ;  $\mu^i(\mathbf{B}^{(i-1)}, \tilde{X}_{i-1})$  as  $\overline{\mu^i}$ ; and  $\pi^i(\mathbf{V}^{(i)})$  as  $\pi^i$ .

Let  $T^{dml}(\{\pi^k\}_{k=1}^{m-1}, \{\mu^k\}_{k=2}^m)$  be a quantity defined in Eq. (C.18). We first note that

$$T^{dml}(\{\pi_0^k\}_{k=1}^{m-1}, \{\mu_0^k\}_{k=2}^m) = \mathbb{E}\left[Y|do(\mathbf{x})\right]$$

by Eq. (C.7). Then, by Lemma C.3,

$$T^{dml} - \mathbb{E}\left[Y|do(\mathbf{x})\right]$$

$$= T^{dml} - T^{dml}(\{\pi_{0}^{k}\}_{k=1}^{m-1}, \{\mu_{0}^{k}\}_{k=2}^{m})$$

$$= \sum_{i=2}^{m} \mathbb{E}_{P_{x_{i}}-D_{i}}\left[\pi_{0}^{(i-1)}\mathbb{1}_{\mathbf{x}^{(i-1)}}(\mathbf{X}^{(i-1)})\left\{\overline{\mu}_{0}^{i+1}-\mu_{0}^{i}\right\}\right] + \mathbb{E}_{P_{x_{1}}-D_{1}}\left[\overline{\mu}_{0}^{2}\right]$$

$$+ \sum_{i=2}^{m} \mathbb{E}_{P_{x_{i}}-D_{i}}\left[\pi_{0}^{(i-1)}\mathbb{1}_{\mathbf{x}^{(i-1)}}(\mathbf{X}^{(i-1)})\left\{\overline{\mu}_{0}^{i+1}-\mu_{0}^{i}\right\} - \pi^{(i-1)}\mathbb{1}_{\mathbf{x}^{(i-1)}}(\mathbf{X}^{(i-1)})\left\{\overline{\mu}_{0}^{i+1}-\mu^{i}\right\}\right] + \mathbb{E}_{P_{x_{1}}-D_{1}}\left[\overline{\mu}_{0}^{2}-\overline{\mu}_{0}^{2}\right]$$

$$(C.27)$$

$$+\sum_{i=2}^{m} \mathbb{E}_{P_{x_i}} \left[ \pi^{(i-1)} \mathbb{1}_{\mathbf{x}^{(i-1)}} (\mathbf{X}^{(i-1)}) \left\{ \overline{\mu}^{i+1} - \mu^i \right\} \right] + \mathbb{E}_{P_{x_1}} \left[ \overline{\mu}^2 - \overline{\mu}_0^2 \right].$$
(C.28)

We first note that

Eq. (C.14) = 
$$\sum_{i=1}^{m} o_{P_{x_i}}(n_i^{-1/2})$$

under Assumptions (2,8) by Lemma C.3.

Then,

Eq. (C.13) + Eq. (C.14) = 
$$\sum_{i=1}^{m} R_i$$
,

where  $R_i$  for  $i = 1, 2, \dots, m$  are variables converging in mean-zero normal distribution, by the central limit theorem and Slutsky's theorem.

Finally

Eq. (C.15) = 
$$T^{dml}(\{\pi^k\}_{k=1}^{m-1}, \{\mu^k\}_{k=2}^m) - \mathbb{E}[Y|do(\mathbf{x})]$$
  
=  $\sum_{i=k+1}^m O_{\overline{P}_{x_i}^{k-1}}(\|\mu^i - \mu_0^i\| \|\pi^{i-1} - \pi_0^{i-1}\|),$ 

where the second equation holds by Lemma C.6.

**Theorem 6** (Error analysis of the estimators for MTI). Under Assumptions (2,7,8,9) and AC-MTI in Def. 7, the errors of the estimators in Def. 9, denoted  $\epsilon^{est} := T^{est} - \mathbb{E}[Y|do(\mathbf{x})]$  for  $est \in \{reg, pw, dml\}$ , are:

$$\begin{split} \epsilon^{\text{reg}} &= R_1 + O_{P_{x_1}} \left( \| \mu^1 - \mu_0^1 \| \right), \\ \epsilon^{\text{pw}} &= R_m + O_{P_{x_m}} (\| \pi^{(m-1)} - \pi_0^{(m-1)} \|), \\ \epsilon^{\text{dml}} &= \sum_{i=1}^m R_i + \sum_{i=2}^m O_{P_{x_i}} \left( \| \mu^i - \mu_0^i \| \| \pi^{i-1} - \pi_0^{i-1} \| \right). \end{split}$$

where  $R_i$  is a random variable such that  $\sqrt{n_i}R_i$  converges in distribution to a zero-mean normal random variable, where  $n_i := |D_{x_i}|$  for  $i \in \{1, \dots, m\}$ .

*Proof of Theorem 6.* The proof is complete by Lemmas (C.4, C.5, C.7).

**Corollary 6** (Multiply robustness of the DML estimators (Corollary of Thm. 6)). Suppose Asumptions (2,7,8,9) and AC-MTI in Def. 7 hold. For  $i = 2, \dots, m-1$ , suppose either  $\pi^{i-1} = \pi_0^{i-1}$  or  $\mu^i = \mu_0^i$ . Then,  $T^{dml}$  in Def. 9 is an unbiased estimator of  $\mathbb{E}[Y|do(\mathbf{x})]$ .

**Proof of Corollary 6.** Let  $T^{dml}(\{\pi^k\}_{k=1}^{m-1}, \{\mu^k\}_{k=2}^m)$  be a quantity defined in Eq. (C.5). Let

$$T^{dml,i} \coloneqq \mathbb{E}_{D_i} \left[ \pi^{(i-1)}(\mathbf{V}^{(i-1)}) \mathbb{1}_{\mathbf{x}^{(i-1)}}(\mathbf{X}^{(i-1)}) \left\{ \mu^{i+1}(\mathbf{B}^{(i)}, x_i) - \mu^i(\mathbf{V}^{(i)}) \right\} \right], \ i = 2, \cdots, m$$

and

$$T^{dml,1} \coloneqq \mathbb{E}_{D_1}\left[\mu^2(B_1, x_1)\right].$$

Under the assumption that samples are i.i.d.,

$$\sum_{i=1}^{m} \mathbb{E}_{P_{x_i}} \left[ T^{dml,i} \right] = T^{dml}(\{\pi^k\}_{k=1}^{m-1}, \{\mu^k\}_{k=2}^m).$$

Then,

$$\sum_{i=1}^{m} \mathbb{E}_{P_{x_i}} \left[ T^{dml,i} \right] - \mathbb{E} \left[ Y | do(\mathbf{x}) \right]$$
  
=  $T^{dml} (\{ \pi^k \}_{k=1}^{m-1}, \{ \mu^k \}_{k=2}^m) - \mathbb{E} \left[ Y | do(\mathbf{x}) \right]$   
=  $\sum_{i=2}^{m} O_{P_{x_i}} \left( \left\| \mu^i - \mu_0^i \right\| \left\| \pi^{i-1} - \pi_0^{i-1} \right\| \right)$   
= 0,

where the third equation holds by Lemma C.6, and the last equation holds under the given condition.

**Theorem 8** (Error analysis of the AC-gMTI estimators). Under Assumptions (2,10,11) and AC-gMTI in Def. 10, the errors of the estimators in Def. 12, denoted  $\epsilon^{est} := T^{est} - \mathbb{E}[Y|do(\mathbf{x})]$  for  $est \in \{reg, pw, dml\}$ , are:

$$\begin{split} \epsilon^{reg} &= R_1 + O_{P_{z_1}}(\|\mu^1 - \mu_0^1\|), \\ \epsilon^{pw} &= R_m + O_{P_{z_m}}(\|\pi^{(m-1)} - \pi_0^{(m-1)}\|), \\ \epsilon^{dml} &= \sum_{i=1}^m R_i + \sum_{i=2}^m O_{P_{z_i}}(\|\mu^i - \mu_0^i\|\|\pi^{i-1} - \pi_0^{i-1}\|), \end{split}$$

where  $R_i$  is a variable such that  $\sqrt{n_i}R_i$  converges in distribution to a zero-mean normal random variable, where  $n_i \coloneqq |D_i|$  for  $i \in \{1, \dots, m\}$ .

*Proof of Theorem 8.* The proof is complete by Lemmas (C.8, C.9, C.11).

**Corollary 8** (Multiply robustness of the DML estimators (Corollary of Thm. 8)). Suppose Assumptions (2,10,11) and AC-gMTI in Def. 10 hold. For  $i = 2, \dots, m-1$ , suppose either  $\pi^{i-1} = \pi_0^{i-1}$  or  $\mu^i = \mu_0^i$ . Then,  $T^{dml}$  in Def. 12 is an unbiased estimator of  $\mathbb{E}[Y|do(\mathbf{x})]$ .

**Proof of Corollary 8.** Let  $T^{dml}(\{\pi^k\}_{k=1}^{m-1}, \{\mu^k\}_{k=2}^m)$  be a quantity defined in Eq. (C.18). Let

$$T^{dml,i} \coloneqq \mathbb{E}_{D_i} \left[ \pi^{(i-1)}(\mathbf{V}^{(i-1)}) \mathbb{1}_{\tilde{\mathbf{x}}^{(i-1)}}(\tilde{\mathbf{X}}^{(i-1)}) \left\{ \mu^{i+1}(\mathbf{B}^{(i)}, \tilde{x}_i) - \mu^i(\mathbf{V}^{(i)}) \right\} \right], \ i = 2, \cdots, m$$

and

$$T^{dml,1} \coloneqq \mathbb{E}_{D_1}\left[\mu^2(B_1, \tilde{x}_1)\right].$$

Under the assumption that samples are i.i.d.,

$$\sum_{i=1}^{m} \mathbb{E}_{P_{z_i}} \left[ T^{dml,i} \right] = T^{dml} \left( \{ \pi^k \}_{k=1}^{m-1}, \{ \mu^k \}_{k=2}^m \right)$$

Then,

$$\begin{split} &\sum_{i=1}^{m} \mathbb{E}_{P_{z_i}} \left[ T^{dml,i} \right] - \mathbb{E} \left[ Y | do(\mathbf{x}) \right] \\ &= T^{dml}(\{\pi^k\}_{k=1}^{m-1}, \{\mu^k\}_{k=2}^m) - \mathbb{E} \left[ Y | do(\mathbf{x}) \right] \\ &= \sum_{i=2}^{m} O_{P_{z_i}} \left( \left\| \mu^i - \mu_0^i \right\| \left\| \pi^{i-1} - \pi_0^{i-1} \right\| \right) \\ &= 0, \end{split}$$

where the third equation holds by Lemma C.10, and the last equation holds under the given condition.

## D. Project STAR: Estimating Joint Effects of Class Sizes to Academic Outcomes

We applied the proposed estimators to Project STAR dataset (Krueger & Whitmore, 2001; Schanzenbach, 2006). Project STAR is an experimental study investigating teacher/student ratios' impact on academic achievement for kindergarten through third-grade students. In the study, students were randomly assigned to three different class sizes: small-size classes, regular classes, and large-size classes. The objective was to evaluate how class size affects academic outcomes (Schanzenbach, 2006). In our analysis, we used the dataset introduced in the online complement of Stock et al. (2003).



Figure D.5: Example causal graphs for Section D. Nodes representing the treatment and outcome are marked in blue and red, respectively.

**Project STAR Dataset.** We denote the Project STAR dataset as D. The dataset D includes the following information: class size for kindergarten  $(X_1)$ , the academic outcome in kindergarten (W), class size for third grade  $(X_2)$ , the academic outcome in third grade (Y), and pre-treatment variables (C) including genders, age, ethnicity, qualification for free lunch, school types, and teacher's education levels. Since Project STAR is a longitudinal experimental study, the samples for variables  $\{C, X_1, W\}$  follow a distribution  $P_{\text{rand}(X_1)}(C, X_1, W)$ , and the samples for variables  $\{C, X_1, W, X_2, Y\}$  follow a distribution  $P_{\text{rand}(X_1, X_2)}(C, X_1, W, X_2, Y)$ .

Assumption on Dataset. We assume that the structural causal model  $\mathcal{M}$  generating the dataset D induces a causal graph depicted in Figure 4a. Specifically, since Project STAR is a longitudinal experimental study randomizing  $X_1$  and  $X_2$ , the submodel  $M_{x_1,x_2}$  for  $x_1, x_2 \in \mathfrak{D}_{X_1,X_2}$  generates the dataset D.

**Creation of Datasets from Marginal Experiments.** In this empirical study, we create two datasets from this dataset:  $D_1$  and  $D_2$ . The dataset  $D_1$  is a random subsample of D only including  $\{C, X_1, W\}$ . Then,  $D_1$  follows  $P_{\text{rand}(X_1)}(C, X_1, W)$ .

We now construct the dataset  $D_2$  following the marginal experimental distribution  $P_{\operatorname{rand}(X_2)}(C, X_1, W, X_2, Y)$  by introducing the confounding bias between  $X_1$  and W as follows. A specific procedure for introducing confounding bias from experimental studies follows an approach widely used in practice<sup>4</sup>, which is described below. Among attributes in C, we chose specific covariates  $C_{\text{bias}} \coloneqq \{\text{ethnicity}, \text{gender}, \text{free-lunch-eligibility}\}$ . Next, we assign probabilities for  $P_{\text{sample}}(x_1|c_{\text{bias}})$  for  $\forall x_1, c_{\text{bias}} \in \mathfrak{D}_{X_1, C_{\text{bias}}}$ . Then, we construct the dataset  $D_2$  as follows:  $D_2 \coloneqq \{\}$ , and for each samples in  $D \coloneqq \{C_{(i)}, X_{1,(i)}, W_{(i)}, X_{2,(i)}, Y_{(i)}\}_{i=1}^{|D|}$ , we repeat the following steps:

- 1. Generate the Bernouli random variable  $B_{(i)}$  with parameter  $P_{\text{sample}}(X_{1,(i)}|C_{\text{bias.}(i)})$ .
- 2. If  $B_{(i)} = 1$ , include  $\{C_{(i)}, X_{1,(i)}, W_{(i)}, X_{2,(i)}, Y_{(i)}\}$  in  $D_2$ .

Finally, we exclude the covariate 'ethnicity' from C in  $D_1$  and  $D_2$ . By doing so, we introduce unmeasured confounding bias between  $X_1$  and W in  $D_2$ . As a result,  $D_2$  follows a marginal experimental distribution  $P_{rand(X_2)}(C, X_1, W, X_2, Y)$ . In this empirical study, the construction of estimators solely relied on the datasets  $D_1$  and  $D_2$ , while the dataset D was exclusively leveraged to construct the ground-truth estimate. The following procedure outlines the specific steps for constructing the ground-truth estimate. Detailed procedures for creating the datasets  $D_1$  and  $D_2$  from D is provided in Appendix E.1.5.

**Goal.** In this empirical study, we aim to study the joint effect of the class size for kindergarten  $(X_1)$  and the third grade  $(X_2)$  on the third grade's academic outcome (Y); i.e.,  $\mathbb{E}[Y|do(x_1, x_2)]$ . Since D is a longitudi-

<sup>&</sup>lt;sup>4</sup>The following procedure introduces confounding bias in an experimental dataset by resampling the dataset with a probability depending on the treatment  $X_1$  and covariates C. The procedure has been used in prior research, such as (Hill, 2011; Louizos et al., 2017; Zhang & Bareinboim, 2019; Gentzel et al., 2021) for simulation purposes.

nal experimental dataset following  $P_{\text{rand}(X_1,X_2)}(C,X_1,W,X_2,Y)$ , the ground-truth  $\mathbb{E}[Y|do(x_1,x_2)]$  is estimated as  $\mathbb{E}_D[Y\mathbbm{1}_{x_1,x_2}(X_1,X_2)]/\mathbb{E}_D[\mathbbm{1}_{x_1,x_2}(X_1,X_2)].$ 

**Causal Effect Identification.** Identifying and estimating the causal effects  $\mathbb{E}[Y|do(x_1, x_2)]$  falls under Task TTI. To witness, we first recall that the datasets  $D_1$  and  $D_2$  consist of samples that follow the distributions  $P_{\text{rand}(X_1)}(C, X_1, W)$  and  $P_{\text{rand}(X_2)}(C, X_1, W, X_2, Y)$ , respectively. Furthermore, within each dataset, the samples  $D_{x_1}$  follow the distribution  $P_{x_1}(C, W)$  and the samples  $D_{x_2}$  follow the distribution  $P_{x_2}(C, X_1, W, Y)$ .

We first observe that  $\{C, W\}$  in the graph G (in Fig. D.5a) satisfies the AC-TTI in Def. 1 w.r.t  $\{(X_1, X_2), Y\}$ . Specifically,

- 1.  $(\{C, W\} \perp X_2 | X_1)_{G_{\overline{X_1}, X_2}};$  and
- 2.  $(Y \perp X_2 | C, W, X_1)_{G_{X_1} \overline{X_2}}$ .

Also, the positivity assumption in Assumption 1 is satisfied for  $D_1$  and  $D_2$ . Therefore, according to Theorem 1, the joint treatment effects  $\mathbb{E}[Y|do(x_1, x_2)]$  are identifiable and can be expressed as follows:

$$\mathbb{E}\left[Y|do(x_1, x_2)\right] = \mathbb{E}_{P_{x_1}}\left[\mathbb{E}_{P_{x_2}}\left[Y|C, W, x_1\right]\right].$$
(D.1)

**Causal Effect Estimation.** We define the nuisance as follows: For the fixed  $x_1, x_2 \in \mathfrak{D}_{X_1, X_2}$ ,

$$\mu_0(C, X_1, W) \coloneqq \mathbb{E}_{P_{x_2}}\left[Y|W, X_1, C\right],\tag{D.2}$$

$$\pi_0(C, X_1, W) \coloneqq \frac{P_{x_1}(W|C)}{P_{x_2}(W, X_1|C)}.$$
(D.3)

Then, besides Eq. (D.1), the causal effect can be expressed as follows:

Eq. (D.1) = 
$$\mathbb{E}_{P_{x_2}}[Y\pi_0(C, X_1, W)\mathbb{1}_{x_1}(X_1)], \text{ or },$$
 (D.4)

$$= \mathbb{E}_{P_{x_2}} \left[ \pi_0(C, X_1, W) \mathbb{1}_{x_1}(X_1) \{ Y - \mu_0(C, X_1, W) \} \right] + \mathbb{E}_{P_{x_1}} \left[ \mu_0(C, x_1, W) \right].$$
(D.5)

We then construct the regression-based, probability weighting-based, and double/debiased machine learning (DML)  $T^{\text{reg}}, T^{\text{pw}}, T^{\text{dml}}$  using the following procedure.

- 1. For each fixed  $x_i \in \mathfrak{D}_{X_i}$  and a sample set  $D_{x_i}$  for  $i \in \{1, 2\}$ , randomly split the sample as  $D_{x_i,t}$  and  $D_{x_i,e}$ .
- 2. Use  $\{D_{x_1,t}, D_{x_2,t}\}$  to train the model for learning nuisances in Eq. (D.2) and Eq. (D.3). Let  $\mu(C, X_1, W)$  and  $\pi(C, X_1, W)$  denote the learnt models. We use the XGBoost (Chen & Guestrin, 2016) to learn the model.
- 3. Then, each estimator is defined as follows:

$$T^{\mathsf{reg}} \coloneqq \mathbb{E}_{D_{\tau_1,e}}\left[\mu(C, x_1, W)\right] \tag{D.6}$$

$$T^{\mathsf{pw}} \coloneqq \mathbb{E}_{D_{x_2,e}}\left[\pi(C, X_1, W) \mathbb{1}_{x_1}(X_1)Y\right]$$
(D.7)

$$T^{\text{dml}} \coloneqq \mathbb{E}_{D_{x_2,e}} \left[ \pi(C, X_1, W) \mathbb{1}_{x_1}(X_1) \{ Y - \mu(C, X_1, W) \} \right] + \mathbb{E}_{D_{x_1,e}} \left[ \mu(C, x_1, W) \right].$$
(D.8)

With the following construction, the Assumption 2 is satisfied.

**Experimental Results.** As described in the Experimental Setup section (Sec. 5), we evaluated the AAE<sup>est</sup> of estimators  $T^{est}$  for est  $\in$  {reg, pw, dml} in Cases {1, 2, 3, 4}. The AAE plots for all cases can be seen in Fig. D.6. In this particular scenario, the sample size was not varied since the sample itself was externally given.

In Case 2, we introduced variation by adjusting the size of the converging noise  $\epsilon$ , which follows a normal distribution Normal $(n^{-\alpha}, n^{-2\alpha})$  for  $n \in \{200, 400, 600, 800, 1000\}$ . It was observed that the DML estimator  $T^{\text{dml}}$  outperformed the

	Case 1	Case 2	Case 3	Case 4
Fig. D.4a (Project STAR)		$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	0.5 Model 0.5 PW 0.4 0.3 0.2 . 0.1	

Figure D.6: AAE Plot for Fig. (D.5a) for Cases {1,2,3,4} depicted in the Experimental Setup in Sec. 5.

other two estimators by achieving fast convergence, as demonstrated in Theorem 2. For Cases  $\{3, 4\}$ , the DML estimator  $T^{\text{dml}}$  exhibited doubly robust properties, as illustrated in Corollary 2.

# **E.** Details of Experiments

As described in Sec. 5, we used the XGBoost (Chen & Guestrin, 2016) as a model for estimating nuisances  $\mu, \pi, \{\mu^i\}_{i=2}^m, \{\pi^i\}_{i=1}^m$ . We implemented the model using Python. In modeling nuisance using the XGBoost, we used the command xgboost.XGBClassifier(eval\_metric='logloss')<sup>5</sup> to use the XGBoost with the default parameter settings. In implementing the PW estimators  $T^{pw}$  and the DML estimators  $T^{dml}$ , we use the clipped weight by trimming samples yielding weights lower than 10 percentile or greater than 90 percentile (Crump et al., 2009). For Tasks (TTI,MTI), d is chosen to be 10. For Task gTTI, d is chosen to be 5. For Task gMTI, d is chosen to be 1.

We split the dataset as training and test samples with a 5:5 ratio. The training samples are used only for running parameters of the XGBoost models, and the test samples are used only for evaluating the trained XGBoost models.

#### **E.1.** Designs of Simulations

This section provides the structural causal models used for generating the dataset. Specifically, we provide a part of the code for generating the dataset.

#### E.1.1. TASK TTI

```
Generate Exogeneous Variables
Generate U_C1_W (Latent confounders between C1, W)
U_C1_W = np.random.normal(0, 1, size=(n,))
Generate U_C1_X1 (Latent confounders between C1, X1)
U_C1_X1 = np.random.normal(0, 1, size=(n,))
Generate U_X1_W (Latent confounders between X1, W)
U_X1_W = np.random.normal(0, 1, size=(n,))
Generate U_X1_X2 (Latent confounders between X1, X2)
U_X1_X2 = np.random.normal(0, 1, size=(n,))
Generate U_X2_Y (Latent confounders between X2, Y)
U_X2_Y = np.random.normal(0, 1, size=(n,))
Generate U_C2_X2 (Latent confounders between C2, X2)
```

<sup>&</sup>lt;sup>5</sup>Detailed parametrization of parameters including learning rates, maximum depth of the trees, etc. are explained in https: //xgboost.readthedocs.io/en/stable/python/python\_api.html#xgboost.XGBClassifier.

```
U_C2_X2 = np.random.normal(0, 1, size=(n,))
 # Generate U_C2_Y (Latent confounders between C2, Y)
U_C2_Y = np.random.normal(0, 1, size=(n,))
••• Generate Endogenous Variables
 # SCM for Covariates C1
def f_C1(n,d,U_C1_X1, U_C1_W):
             C1 = np.zeros((n,d))
             for idx in range(0,d):
                          C1[:,idx] = np.random.normal(0,1,size = (n,)) + U_C1_X1 + U_C1_W
             return(C1)
 # SCM for Treatment X1
def f_X1(n,d,C1, U_C1_X1, U_X1_W, U_X1_X2):
             coeff = np.repeat(1,d)
             X1_linfun = np.dot(C1,coeff) + U_C1_X1 + U_X1_W + U_X1_X2
             X1_param = 1/(1+np.exp(-X1_linfun))
             X1 = np.round(X1_param)
             return(X1)
 # SCM for Output W
def f_W(n, d, C1, X1, U_C1_W, U_X1_W):
             coeff1 = np.repeat(1,d)
             coeff2 = np.repeat(-1, d)
             \label{eq:w_linfun} \texttt{W_linfun} = \texttt{np.dot}(\texttt{C1}, \texttt{ coeff1}) + \texttt{np.dot}(\texttt{C1}, \texttt{ coeff2}) \\ \star \\ \texttt{X1} + \texttt{U_C1}_{\texttt{W}} + \texttt{U_X1}_{\texttt{W}} \\ \texttt{W_linfun} = \texttt{np.dot}(\texttt{C1}, \texttt{ coeff1}) \\ \star \\ \texttt{W_linfun} = \texttt{np.dot}(\texttt{C1}, \texttt{ coeff1
             W_param = 1 / (1 + np.exp(-W_linfun))
             W = np.round(W_param)
             return (W)
 # SCM for Covariates C2
def f_C2(n, d, C1,U_C2_X2, U_C2_Y):
             C2 = np.zeros((n, d))
             for idx in range(0, d):
                          C2[:, idx] = (2*C1[:,idx]-1) + U_C2_X2 + U_C2_Y
             return (C2)
 # SCM for Treatment X2
def f_X2(n, d, C2, X1, U_C2_X2, U_X2_Y, U_X1_X2):
             coeff1 = np.repeat(1, d)
             coeff2 = np.repeat(-1, d)
             X2_linfun = np.dot(C2, coeff1) + np.dot(C2, coeff2) *X1 \
                                            + U_C2_X2 + U_X2_Y + U_X1_X2
             X2_param = 1 / (1 + np.exp(-X2_linfun))
             X2 = np.round(X2_param)
             return (X2)
 # SCM for Y
def f_Y(n, d, C2, X2, W, U_C2_Y, U_X2_Y):
             coeff1 = np.repeat(1, d)
             coeff2 = np.repeat(2, d)
             coeff3 = np.repeat(-1, d)
             Y_linfun = np.dot(C2, coeff1) + np.dot(C2, coeff2) \star X2 \setminus
```

```
+ np.dot(C2, coeff3) * W + U_C2_Y + U_X2_Y
Y_param = 1 / (1 + np.exp(-Y_linfun))
Y = np.round(Y_param)
return (Y)
```

## E.1.2. TASK GTTI

```
···Generate Exogeneous Variables
# Generate U_X0_Z1 (Latent confounders between X0, Z1)
U_X0_Z1 = np.random.normal(0, 1, size=(n,))
## Generate U_X0_Z2 (Latent confounders between X0, Z2)
U_X0_Z2 = np.random.normal(0, 1, size=(n,))
## Generate U_Z1_W (Latent confounders between Z1, W)
U_Z1_W = np.random.normal(0, 1, size=(n,))
## Generate U_Z2_Y (Latent confounders between Z2, Y)
U_Z2_Y = np.random.normal(0, 1, size=(n,))
••• Generate Endogenous Variables
# SCM for Treatment C
def f_C(n,d):
   C = np.zeros((n, d))
    for idx in range(0, d):
        C[:, idx] = np.random.normal(0, 1, size=(n,))
    return (C)
# SCM for Treatment X0
def f_X0(n,d, U_X0_Z1 , U_X0_Z2):
    X0_linfun = U_X0_Z1 - U_X0_Z2 + 0.5 + np.random.normal(0, 1, size=(n,))
    X0_param = 1/(1+np.exp(-X0_linfun))
    X0 = np.round(X0_param)
    return(X0)
# SCM for Treatment Z1
def f_Z1(n, d, C, X0, U_X0_Z1, U_Z1_W):
    coeff1 = np.repeat(1, d)
    Z1_linfun = np.dot(C, coeff1) * (2 * X0−1) + U_X0_Z1 + U_Z1_W + X0 \
                + np.random.normal(0, 1, size=(n,))
    Z1_param = 1 / (1 + np.exp(-Z1_linfun))
    Z1 = np.round(Z1_param)
    return (Z1)
# SCM for W
def f_W(n,d,C,Z1,U_Z1_W):
    coeff1 = np.repeat(-0.5, d)
    W_linfun = np.dot(C, coeff1) * (2 * Z1-1) + U_Z1_W \setminus
             + np.random.normal(0, 1, size=(n,))
    W_param = 1 / (1 + np.exp(-W_linfun))
    W = np.round(W_param)
    return (W)
```

```
# SCM for Treatment Z2
def f_Z2(n, d, C, X0, Z1, U_X0_Z2, U_Z2_Y):
    coeff1 = np.repeat(-1, d)
    coeff2 = np.repeat(0.5, d)
    U = 0.5 \star (U_X0_Z2 + U_Z2_Y)
    Z2_linfun = np.dot(C,coeff1) * (2*X0-1) + np.dot(C,coeff2) * (2*Z1-1) + U \
                + np.random.normal(0, 1, size=(n,))
    Z2_param = 1 / (1 + np.exp(-Z2_linfun))
    Z2 = np.round(Z2_param)
    return (Z2)
# SCM for Y
def f_Y(n, d, C, X0, Z2, W, U_Z2_Y):
    coeff1 = np.repeat(-1, d)
    coeff3 = np.repeat(-0.5, d)
    coeff4 = np.repeat(2, d)
    U = 0.5 \star U_Z2_Y
    Y_{linfun} = np.dot(C, coeff1) * (2*X0-1) + np.dot(C, coeff3) * (2*Z2-1) 
              + np.dot(C, coeff4) * (2*W-1) + \
               U + np.random.normal(0, 1, size=(n,))
    Y_param = 1 / (1 + np.exp(-Y_linfun))
    Y = np.round(Y_param)
    return (Y)
```

## E.1.3. TASK MTI

```
··· Generate Exogeneous Variables ···
# Generate U_C1_W1 (Latent confounder between C1, W1)
U_C1_W1 = np.random.normal(0, 1, size=(n,))
## Generate U_C1_X1 (Latent confounder between C1, X1)
U_C1_X1 = np.random.normal(0, 1, size=(n,))
## Generate U_X1_W1 (Latent confounder between X1, W1)
U_X1_W1 = np.random.normal(0, 1, size=(n,))
## Generate U_C2_W2 (Latent confounder between C2, W2)
U_C2_W2 = np.random.normal(0, 1, size=(n,))
## Generate U_C2_X2 (Latent confounder between C2, X2)
U_C2_X2 = np.random.normal(0, 1, size=(n,))
## Generate U_X2_W2 (Latent confounder between X2, W2)
U_X2_W2 = np.random.normal(0, 1, size=(n,))
## Generate U_C3_Y (Latent confounder between C3, Y)
U_C3_Y = np.random.normal(0, 1, size=(n,))
## Generate U_C3_X3 (Latent confounder between C3, X3)
U_C3_X3 = np.random.normal(0, 1, size=(n,))
```

```
## Generate U_X3_Y (Latent confounder between X3, Y)
U_X3_Y = np.random.normal(0, 1, size=(n,))
··· Generate Endogenous Variables ···
# SCM for Covariates C1
def f_C1(n,d,U_C1_X1,U_C1_W1):
    C1 = np.zeros((n,d))
    for idx in range(0,d):
        C1[:,idx] = np.random.normal(0,1,size = (n,)) + U_C1_X1 + U_C1_W1
    return(C1)
# SCM for Treatment X1
def f_X1(n,d,C1, U_C1_X1, U_X1_W1):
    coeff = np.repeat(1,d)
    X1_linfun = np.dot(C1,coeff) + U_C1_X1 + U_X1_W1
    X1_param = 1/(1+np.exp(-X1_linfun))
    X1 = np.round(X1_param)
    return(X1)
# SCM for Output W1
def f_W1(n, d, C1, X1, U_C1_W1, U_X1_W1):
    coeff1 = np.repeat(1, d)
    coeff2 = np.repeat(-1, d)
    W1\_linfun = np.dot(C1, coeff1) + np.dot(C1, coeff2) * X1 + U\_C1\_W1 + U\_X1\_W1
    W1_param = 1 / (1 + np.exp(-W1_linfun))
    W1 = np.round(W1_param)
    return (W1)
# SCM for Covariates C2
def f_C2(n, d, C1,U_C2_X2, U_C2_W2):
    C2 = np.zeros((n, d))
    for idx in range(0, d):
        C2[:, idx] = (2*C1[:, idx]-1) + U_C2_X2 + U_C2_W2
    return (C2)
# SCM for Treatment X2
def f_X2(n, d, C2, U_C2_X2, U_X2_W2):
    coeff1 = np.repeat(1, d)
    X2\_linfun = np.dot(C2, coeff1) + U\_C2\_X2 + U\_X2\_W2
    X2_param = 1 / (1 + np.exp(-X2_linfun))
    X2 = np.round(X2_param)
    return (X2)
# SCM for Output W2
def f_W2(n, d, C2, X2, W1, U_C2_W2, U_X2_W2):
    coeff1 = np.repeat(1, d)
    coeff2 = np.repeat(2, d)
    coeff3 = np.repeat(-1, d)
    W2_linfun = np.dot(C2, coeff1) + np.dot(C2, coeff2) * X2 +\
                np.dot(C2, coeff3) * W1 + U_C2_W2 + U_X2_W2
    W2_param = 1 / (1 + np.exp(-W2_linfun))
    W2 = np.round(W2_param)
```

```
return (W2)
# SCM for Covariates C3
def f_C3(n, d, C2, U_C3_X3, U_C3_Y):
    C3 = np.zeros((n, d))
    for idx in range(0, d):
        C3[:, idx] = (2 * C2[:, idx] - 1) + U_C3_X3 + U_C3_Y
    return (C3)
# SCM for Treatment X3
def f_X3(n, d, C3, U_C3_X3, U_X3_Y):
    coeff1 = np.repeat(1, d)
    X3\_linfun = np.dot(C3, coeff1) + U_C3\_X3 + U_X3\_Y
    X3_param = 1 / (1 + np.exp(-X3_linfun))
    X3 = np.round(X3_param)
    return (X3)
# SCM for Output Y
def f_Y(n, d, C3, X3, W2, U_C3_Y, U_X3_Y):
    coeff1 = np.repeat(1, d)
    coeff2 = np.repeat(2, d)
    coeff3 = np.repeat(-1, d)
    Y_linfun = np.dot(C3, coeff1) + np.dot(C3, coeff2) * X3 +\
              np.dot(C3, coeff3) * W2 + U_C3_Y + U_X3_Y
    Y_param = 1 / (1 + np.exp(-Y_linfun))
    Y = np.round(Y_param)
    return (Y)
```

# E.1.4. TASK GMTI

```
Generate Exogeneous Variables
Generate U_X0_Z1 (Latent Confounders between X0, Z1)
U_X0_Z1 = np.random.normal(0, 1, size=(n,))
Generate U_X0_Z2 (Latent Confounders between X0, Z2)
U_X0_Z2 = np.random.normal(0, 1, size=(n,))
Generate U_X0_Z3 (Latent Confounders between X0, Z3)
U_X0_Z3 = np.random.normal(0, 1, size=(n,))
Generate U_Z1_W1 (Latent Confounders between Z1, W)
U_Z1_W1 = np.random.normal(0, 1, size=(n,))
Generate U_Z2_W2 (Latent Confounders between Z2, W2)
U_Z2_W2 = np.random.normal(0, 1, size=(n,))
Generate U_Z3_Y (Latent Confounders between Z3, Y)
U_Z3_Y = np.random.normal(0, 1, size=(n,))
Score Endogenous Variables
Score for Covariate C1
```

```
def f Cl(n,d):
    C1 = np.zeros((n, d))
    for idx in range(0, d):
        C1[:, idx] = np.random.normal(0, 1, size=(n,))
    return (C1)
# SCM for Treatment X0
def f_X0(n,d, U_X0_Z1 , U_X0_Z2):
    X0_linfun = U_X0_Z1 - U_X0_Z2 + 0.5 + np.random.normal(0, 1, size=(n,))
    X0_param = 1/(1+np.exp(-X0_linfun))
   X0 = np.round(X0_param)
    return(X0)
# SCM for Treatment Z1
def f_Z1(n, d, C1, X0, U_X0_Z1, U_Z1_W1):
    coeff1 = np.repeat(1, d)
    Z1_linfun = np.dot(C1, coeff1) * (2 * X0 - 1) + U_X0_Z1 + U_Z1_W1 \
               - X0 + np.random.normal(0, 1, size=(n,))
    Z1_param = 1 / (1 + np.exp(-Z1_linfun))
    Z1 = np.round(Z1_param)
    return (Z1)
# SCM for Outcome W1
def f_W1(n,d,C1,Z1,U_Z1_W1):
    coeff1 = np.repeat(-0.5, d)
    W1\_linfun = np.dot(C1, coeff1) * (2 * Z1-1) + U_Z1_W1 
                + np.random.normal(0, 1, size=(n,))
    W1_param = 1 / (1 + np.exp(-W1_linfun))
    W1 = np.round(W1_param)
    return (W1)
# SCM for Covariate C2
def f_C2(n, d):
    C2 = np.zeros((n, d))
    for idx in range(0, d):
        C2[:,idx] = np.random.normal(0, 1, size=(n,))
    return (C2)
# SCM for Treatment Z2
def f_Z2(n, d, C1, X0, Z1, C2, U_X0_Z2, U_Z2_W2):
    coeff1 = np.repeat(-1, d)
    coeff2 = np.repeat(0.5, d)
   U = 0.5 * (U_X0_Z2 + U_Z2_W2)
    Z2_linfun = np.dot(C1,coeff1)*(2*X0-1) + np.dot(C2,coeff2)*(2*Z1-1) \
             + U + np.random.normal(0, 1, size=(n,))
    Z2_param = 1 / (1 + np.exp(-Z2_linfun))
    Z2 = np.round(Z2_param)
    return (Z2)
# SCM for Outcome W2
def f_W2(n, d, C1, C2, Z2, W1, U_Z2_W2):
    coeff1 = np.repeat(-1, d)
    coeff2 = np.repeat(-0.5, d)
```

```
coeff3 = np.repeat(2, d)
    U = 0.5 * U_22_W2
    W2_linfun = np.dot(C1, coeff1) * (2*X0-1) + np.dot(C2, coeff2) * (2*Z2-1) \
               + np.dot(C1 + C2, coeff3) * (2*W1-1) + \setminus
               U + np.random.normal(0, 1, size=(n,))
    W2_param = 1 / (1 + np.exp(-W2_linfun))
    W2 = np.round(W2_param)
    return (W2)
# SCM for Treatment Z3
def f_Z3(n, d, C2, X0, Z2, U_X0_Z3, U_Z3_Y):
    coeff1 = np.repeat(-1, d)
    coeff2 = np.repeat(-0.5, d)
    coeff3 = np.repeat(2, d)
    U = 0.5 * (U_X0_Z3 + U_Z3_Y)
    Z3\_linfun = np.dot(C2, coeff1) + np.dot(C2, coeff2) * (2*X0-1) 
               + np.dot(C2, coeff3) * (2*Z2-1) + U*(2*X0-1)*(2*Z2-1) \
            +np.random.normal(0, 1, size=(n,))
    Z3_param = 1 / (1 + np.exp(-Z3_linfun))
    Z3 = np.round(Z3_param)
    return (Z3)
# SCM for Outcome Y
def f_Y(n, d, C2, X0, Z3, W2, U_Z3_Y):
    coeff1 = np.repeat(-1, d)
    coeff2 = np.repeat(-0.5, d)
    coeff3 = np.repeat(2, d)
    U = 0.5 \star U_Z3_Y
    Y_linfun = np.dot(C2, coeff1) \star (2 \star X0 - 1) \setminus
             + np.dot(C2, coeff2) * (2 * Z3 - 1) \
            + np.dot(C2, coeff3) * (2 * W2 - 1) \
            + U + np.random.normal(0, 1, size=(n,))
    Y_param = 1 / (1 + np.exp(-Y_linfun))
    Y = np.round(Y_param)
    return (Y)
```

E.1.5. DATA GENERATION FOR PROJECT STAR

We obtained the Project STAR dataset from the following R-package, https://rdrr.io/cran/AER/man/STAR. html. Then, we used the following code for constructing  $D_1$  and  $D_2$  datasets used in analyzing the Project STAR dataset.

```
def normalize(vec):
    veccopy = copy.copy(vec)
    maxval = np.max(vec)
    minval = np.min(vec)
    return (veccopy - minval) / (maxval - minval)
def preprocess_STAR_D1(STAR_D1):
```

Estimating Joint Treatment Effects by Combining Multiple Experiments

```
selected_columns = ['gender', 'birth', 'stark', 'readk', \
                        'mathk', 'lunchk', 'schoolk', 'degreek', 'experiencek']
   STAR_D1 = copy.copy(STAR_D1[selected_columns])
   STAR_D1 = STAR_D1.dropna()
    # Numericalize all columns
    ## Binarize the gender
    gender = np.array(STAR_D1['gender'])
   binarize_gender = [1 if item == 'male' else 0 for item in gender]
   STAR_D1.loc[:, 'gender'] = binarize_gender
    # ## one hot encoding of ethnicity
    # one_hot_encoded = pd.get_dummies(STAR_D1['ethnicity'], prefix='ethnicity')
    # STAR_D1 = STAR_D1.drop('ethnicity', axis=1)
    # STAR_D1 = pd.concat([STAR_D1, one_hot_encoded], axis=1)
    ## Numericalize the birth
   birth = np.array(STAR_D1['birth'])
   age_from_birth = [1988 - (int(year) + (int(quarter[-1]) - 1) / 4) \setminus
                    for year, quarter in [quarter.split(' ') for quarter in birth]]
   STAR_D1.loc[:, 'birth'] = age_from_birth
    ## Binarize the stark
   stark = np.array(STAR_D1['stark'])
   binarize_stark = [1 if item == 'small' else 0 for item in stark]
   STAR_D1.loc[:, 'stark'] = binarize_stark
    ## Make the resultk := readk + mathk
   STAR_D1.loc[:, 'resultk'] = STAR_D1['readk'] + STAR_D1['mathk']
   STAR_D1 = STAR_D1.drop('readk', axis=1)
   STAR_D1 = STAR_D1.drop('mathk', axis=1)
   STAR_D1.loc[:, 'resultk'] = normalize(STAR_D1['resultk'])
    ## Binarize the lunchk
   lunchk = np.array(STAR_D1['lunchk'])
   binarize_lunchk = [1 if item == 'small' else 0 for item in lunchk]
   STAR_D1.loc[:, 'lunchk'] = binarize_lunchk
    ## one hot encoding of schoolk
   one_hot_encoded = pd.get_dummies(STAR_D1['schoolk'], prefix='schoolk')
   STAR_D1 = STAR_D1.drop('schoolk', axis=1)
   STAR_D1 = pd.concat([STAR_D1, one_hot_encoded], axis=1)
    ## one hot encoding of degreek
   one_hot_encoded = pd.get_dummies(STAR_D1['degreek'], prefix='schoolk')
   STAR_D1 = STAR_D1.drop('degreek', axis=1)
   STAR_D1 = pd.concat([STAR_D1, one_hot_encoded], axis=1)
   return STAR_D1
def preprocess_STAR_D2(STAR_D2):
    selected_columns_before = ['gender', 'birth', 'stark', 'readk', \
                                'mathk', 'lunchk', 'schoolk', 'degreek', 'experiencek'] \
                                + ['read3', 'math3', 'star3'] + ['ethnicity']
```

Estimating Joint Treatment Effects by Combining Multiple Experiments

```
selected_columns = ['gender', 'birth', 'stark', 'readk', \
                    'mathk', 'lunchk', 'schoolk', 'degreek', 'experiencek'] \
                    + ['read3', 'math3', 'star3']
# Take the pre-k randomization
STAR_D2 = copy.copy(STAR_D2[selected_columns_before])
STAR_D2 = STAR_D2.dropna()
STAR_D2 = introduceConfoundingD2 (STAR_D2)
# Numericalize all columns
## Binarize the gender
gender = np.array(STAR_D2['gender'])
binarize_gender = [1 if item == 'male' else 0 for item in gender]
STAR_D2.loc[:, 'gender'] = binarize_gender
## one hot encoding of ethnicity
# one_hot_encoded = pd.get_dummies(STAR_D2['ethnicity'], prefix='ethnicity')
# STAR_D2 = STAR_D2.drop('ethnicity', axis=1)
# STAR_D2 = pd.concat([STAR_D2, one_hot_encoded], axis=1)
## Numericalize the birth
birth = np.array(STAR_D2['birth'])
age_from_birth = [1988 - (int(year) + (int(quarter[-1]) - 1) / 4)
               for year, quarter in [quarter.split(' ') for quarter in birth]]
STAR_D2.loc[:, 'birth'] = age_from_birth
## Binarize the stark
stark = np.array(STAR_D2['stark'])
binarize_stark = [1 if item == 'small' else 0 for item in stark]
STAR_D2.loc[:, 'stark'] = binarize_stark
## Make the resultk := readk + mathk
STAR_D2.loc[:, 'resultk'] = STAR_D2['readk'] + STAR_D2['mathk']
STAR_D2 = STAR_D2.drop('readk', axis=1)
STAR_D2 = STAR_D2.drop('mathk', axis=1)
STAR_D2.loc[:, 'resultk'] = normalize(STAR_D2['resultk'])
## Binarize the lunchk
lunchk = np.array(STAR_D2['lunchk'])
binarize_lunchk = [1 if item == 'small' else 0 for item in lunchk]
STAR_D2.loc[:, 'lunchk'] = binarize_lunchk
## one hot encoding of schoolk
one_hot_encoded = pd.get_dummies(STAR_D2['schoolk'], prefix='schoolk')
STAR_D2 = STAR_D2.drop('schoolk', axis=1)
STAR_D2 = pd.concat([STAR_D2, one_hot_encoded], axis=1)
## one hot encoding of degreek
one_hot_encoded = pd.get_dummies(STAR_D2['degreek'], prefix='schoolk')
STAR_D2 = STAR_D2.drop('degreek', axis=1)
STAR_D2 = pd.concat([STAR_D2, one_hot_encoded], axis=1)
## Binarize the star3
star3 = np.array(STAR_D2['star3'])
```

binarize\_star3 = [1 if item == 'small' else 0 for item in star3]

```
STAR_D2.loc[:, 'star3'] = binarize_star3
    ## Make the result3 := read3 + math3
   STAR_D2.loc[:, 'result3'] = STAR_D2['read3'] + STAR_D2['math3']
   STAR_D2 = STAR_D2.drop('read3', axis=1)
   STAR_D2 = STAR_D2.drop('math3', axis=1)
   STAR_D2.loc[:, 'result3'] = normalize(STAR_D2['result3'])
   return STAR_D2
def introduceConfoundingD2(STAR_D2):
    # ethnicity: cauc, afam, asian, hispanic, others
    # gender: female, male
   prob_treat_covariate_matrix = {('cauc', 'male', 'free'): 0.8,
                                    ('cauc', 'male', 'non-free'): 0.4,
                                    ('cauc', 'female', 'free'): 0.4,
                                    ('cauc', 'female', 'non-free'): 0.8,
                                    ('afam', 'male', 'free'): 0.2,
                                    ('afam', 'male', 'non-free'): 0.2,
                                    ('afam', 'female', 'free'): 0.9,
                                    ('afam', 'female', 'non-free'): 0.9,
                                    ('asian', 'male', 'free'): 0.56,
                                    ('asian', 'male', 'non-free'): 0.3,
                                    ('asian', 'female', 'free'): 0.7,
                                    ('asian', 'female', 'non-free'): 0.8,
                                    ('hispanic', 'male', 'free'): 0.7,
                                    ('hispanic', 'male', 'non-free'): 0.4,
                                    ('hispanic', 'female', 'free'): 0.6,
                                    ('hispanic', 'female', 'non-free'): 0.45,
                                    ('other', 'male', 'free'): 0.75,
                                    ('other', 'male', 'non-free'): 0.35,
                                    ('other', 'female', 'free'): 0.48,
                                    ('other', 'female', 'non-free'): 0.25,
                                    }
   def get_weights(row):
        if row['stark'] == 'small':
            return \
            prob_treat_covariate_matrix[row['ethnicity'], row['gender'], row['lunchk']]
        else:
            return \
            1- prob_treat_covariate_matrix[row['ethnicity'], row['gender'], row['lunchk']]
    # Apply the function to assign weights
   STAR_D2['ethnicity_weight'] = STAR_D2.apply(get_weights,axis=1)
    # Resample the DataFrame based on weights
    resampled_STAR_D2 = STAR_D2.sample(n=len(STAR_D2), replace=True, \
                        weights=STAR_D2['ethnicity_weight'], random_state=42)
   STAR_D2 = copy.copy(resampled_STAR_D2)
    return STAR_D2
```
```
def dataMatrixGen(seednum_train = 123):
    . . .
   Form the dataset D1 (pre-k randomization) from the selected columns
    ....
    # Read the CSV without any missing data
   STAR = pd.read_csv("STAR.csv")
    ## D1: In STAR, star3 is NaN
    ## D2: In STAR, star3 is not NaN
   STAR_D1 = copy.copy(STAR[STAR['star3'].isna()])
   STAR_D1 = preprocess_STAR_D1(STAR_D1)
   STAR_D2 = copy.copy(STAR.dropna(subset=['star3']))
   STAR_D2 = preprocess_STAR_D2(STAR_D2)
   X1_x1 = np.array(STAR_D1['stark'])
   W_x1 = np.array(STAR_D1['resultk'])
   D1_covariates = copy.copy(STAR_D1)
   D1_covariates = D1_covariates.drop('stark', axis=1)
   D1_covariates = D1_covariates.drop('resultk', axis=1)
   D1_covariates = STAR_D1[D1_covariates.columns]
   C_x1 = D1_covariates.values # confounders
   D1_mat = np.concatenate((C_x1, W_x1[:, np.newaxis], X1_x1[:, np.newaxis]), axis=1)
   \dim_C = C_{x1.shape[1]}
   X1_x2 = np.array(STAR_D2['stark'])
   X2_x2 = np.array(STAR_D2['star3'])
   Y_x2 = np.array(STAR_D2['result3'])
   W_x2 = np.array(STAR_D2['resultk'])
   D2 covariates = STAR D2[D1 covariates.columns]
   C_x2 = D2_covariates.values
   D2_mat = np.concatenate((C_x2, W_x2[:, np.newaxis], X1_x2[:, np.newaxis], \
                            X2_x2[:, np.newaxis], Y_x2[:, np.newaxis]), axis=1)
   return D1_mat, D2_mat, dim_C
def groundTruthGen(seednum=123):
    . . .
   Form the dataset D1 (pre-k randomization) from the selected columns
    ....
    # Read the CSV without any missing data
   STAR = pd.read_csv("STAR.csv").dropna()
   STAR = preprocess_STAR(STAR)
    ...
   Groundtruth data
    ...
   truth_data_00 = STAR[(STAR['stark'] == 0) & (STAR['star3'] == 0)]['result3'].tolist()
```

Estimating Joint Treatment Effects by Combining Multiple Experiments

```
truth_data_01 = STAR[(STAR['stark'] == 0) & (STAR['star3'] == 1)]['result3'].tolist()
   truth_data_10 = STAR[(STAR['stark'] == 1) & (STAR['star3'] == 0)]['result3'].tolist()
   truth_data_11 = STAR[(STAR['stark'] == 1) & (STAR['star3'] == 1)]['result3'].tolist()
   truth_data_list = [[truth_data_00, truth_data_01], [truth_data_10, truth_data_11]]
   truth_00 = np.mean(STAR[(STAR['stark'] == 0) & (STAR['star3'] == 0)]['result3'])
   truth_01 = np.mean(STAR[(STAR['stark'] == 0) & (STAR['star3'] == 1)]['result3'])
   truth_10 = np.mean(STAR[(STAR['stark'] == 1) & (STAR['star3'] == 0)]['result3'])
   truth_11 = np.mean(STAR[(STAR['stark'] == 1) & (STAR['star3'] == 1)]['result3'])
   truth_list = [[truth_00, truth_01], [truth_10, truth_11]]
   return truth_list
if __name__ == '__main__':
    seednum_train = 123
    . . .
    Form the dataset D1 (pre-k randomization) from the selected columns
    ...
    # Read the CSV without any missing data
   STAR = pd.read_csv("STAR.csv")
    ## D1: In STAR, star3 is NaN
    ## D2: In STAR, star3 is not NaN
   STAR_D1 = copy.copy(STAR[STAR['star3'].isna()])
   STAR_D1 = STAR.sample(n=len(STAR_D1), replace=True)
   STAR_D1 = preprocess_STAR_D1(STAR_D1)
   STAR_D2 = copy.copy(STAR.dropna(subset=['star3']))
   STAR_D2 = preprocess_STAR_D2 (STAR_D2)
```

## F. Discussion on Relaxation of Shared Covariates Assumptions

In this section, we will explore potential relaxations of the shared covariates assumptions, namely Assumptions (4,9).

## F.1. On Assumption 4.

Assumption 4 captures the scenario when  $C_1$  is a pre-treatment covariate for  $X_1$  and  $X_2$ . In this case, intervening on either  $X_1$  or  $X_2$  does not directly affect the distribution of  $C_1$ . Consequently, the distribution of  $C_1$  under treatment  $x_1 (P_{x_1}(C_1))$  is the same as the distribution of  $C_1$  under treatment  $x_2 (P_{x_2}(C_1))$ . Thus, Assumption 4 is satisfied. Under this assumption, the nuisance  $\pi_0$  was given as

$$\pi_0 \coloneqq \pi_0(W, X_1 | C_1) \coloneqq \frac{P_{x_1}(W | C_1)}{P_{x_2}(W, X_1 | C_1)}.$$
(F.1)

We note that this quantity can be easily estimated when W and  $X_1$  are low dimensional discrete variables.

We note that Assumption 4 may not hold when  $C_1$  is a pre-treatment variable for  $X_2$  but a post-treatment covariate for  $X_1$ . In such cases, the distribution of  $C_1$  may differ between the two treatment groups, violating Assumption 4. For example, consider the causal graph depicted in Fig. F.7a, where the AC-TTI criterion in Def. 1 is satisfied. In this graph, Assumption 4 is violated since  $C_1$  is a post-treatment covariate for  $X_1$  but a pre-treatment covariate for  $X_2$ , and therefore,



Figure F.7: An example causal graph such that the AC-TTI criterion in Def. 1 is violated while Assumption 4 is violated.

 $P_{x_2}(C_1) = P(C_1) \neq P_{x_1}(C_1).$ 

To relax Assumption 4, we re-define the nuisance  $\pi_0(C_1, X_1, W)$  in Def. 2 as follows:

$$\pi_0 \coloneqq \pi_0(C_1, X_1, W) \coloneqq \frac{P_{x_1}(W, C_1)}{P_{x_2}(W, X_1, C_1)}.$$
(F.2)

Let  $T^{pw}$  and  $T^{dml}$  denote PW and DML estimators for the AC-TTI functional in Def. 3 equipped with the re-defined nuisance  $\pi$  estimated for  $\pi_0$  in Eq. (F.2). Then, all error analysis results in Thm. 2 is preserved:

**Theorem S.1** (Error analysis of the AC-TTI estimators).  $T^{pw}$  and  $T^{dml}$  denote PW and DML estimators for the AC-TTI functional in Def. 3 equipped with the re-defined nuisance  $\pi$  estimated for  $\pi_0$  in Eq. (F.2). Under Assumptions (1,2,3) and AC-TTI in Def. 1, the error of the estimators in Def. 3, denoted  $\epsilon^{est} := T^{est} - \mathbb{E}[Y|do(x_1, x_2)]$  for  $est \in \{pw, dml\}$  are:

$$\epsilon^{pw} = R_2 + O_{P_{x_2}} \left( \|\pi - \pi_0\| \right),$$
  

$$\epsilon^{dml} = R_1 + R_2 + O_{P_{x_2}} \left( \|\pi - \pi_0\| \|\mu - \mu_0\| \right)$$

where  $R_i$  is a random variable such that  $\sqrt{n_i}R_i$  converges in distribution to a zero-mean normal random variable, where  $n_i := |D_{x_i}|$  for  $i \in \{1, 2\}$ .

**Proof of Theorem S.1**. For the error analysis of  $T^{pw}$  with the redefined nuisance  $\pi_0$  in Eq. (F.2), it suffices to show that

$$\mathbb{E}_{P_{x_2}}\left[\pi_0(W, C_1, X_1)\mathbb{1}_{x_1}(X_1)Y\right] = \mathbb{E}\left[Y|do(x_1, x_2)\right],\tag{F.3}$$

because the other proofs are the same as the proof for the original  $T^{pw}$ .

$$\begin{split} \mathbb{E}_{P_{x_2}} \left[ \pi_0(W, C_1, X_1) \mathbbm{1}_{x_1}(X_1) Y \right] &= \mathbb{E}_{P_{x_2}} \left[ \frac{P_{x_1}(W, C_1)}{P_{x_2}(W, X_1, C_1)} \mathbbm{1}_{x_1}(X_1) Y \right] \\ &= \mathbb{E}_{P_{x_2}} \left[ \frac{P_{x_1}(W, C_1)}{P_{x_2}(W, X_1, C_1)} \mathbbm{1}_{x_1}(X_1) \mu_0(W, X_1, C_1) \right] \\ &= \mathbb{E}_{P_{x_2}} \left[ \frac{P_{x_1}(W, C_1, X_1)}{P_{x_2}(W, X_1, C_1)} \mu_0(W, X_1, C_1) \right] \\ &= \mathbb{E}_{P_{x_1}} \left[ \mu_0(W, x_1, C_1) \right] \\ &= \mathbb{E} \left[ Y | do(x_1, x_2) \right]. \end{split}$$

For the error analysis of  $T^{dml}$  with the redefined nuisance  $\pi_0$  in Eq. (F.2), it suffices to show the following:

$$\mathbb{E}_{P_{x_1}} \left[ \mu(W, C_1, x_1) - \mu_0(W, C_1, x_1) \right] \\= \mathbb{E}_{P_{x_2}} \left[ \frac{P_{x_1}(W, C_1) \mathbb{1}_{x_1}(X_1)}{P_{x_2}(W, C_1, X_1)} \{ \mu(W, C_1, X_1) - \mu_0(W, C_1, X_1) \} \right]$$
(F.4)

$$= \mathbb{E}_{P_{x_2}} \left[ \pi_0(W, C_1, X_1) \mathbb{1}_{x_1}(X_1) \{ \mu(W, C_1, X_1) - \mu_0(W, C_1, X_1) \} \right],$$
(F.5)

because the other proofs are the same as the proof for the original  $T^{dml}$ .

In summary, the purpose of Assumption 4 is to simplify the estimation process by avoiding the need for joint density estimation in Eq. (F.2). Instead, it allows us to use the nuisance in Eq. (F.1), which is amenable to estimate.

## F.2. On Assumption 9.

We recall that the nuisance  $\pi_0^i$  was given as

$$\pi_0^i(\mathbf{W}^{(i)}, \mathbf{C}^{(i)}, \mathbf{X}^{(i)}) \coloneqq \frac{P_{x_i}(W_i | \mathbf{W}^{(i-1)}, \mathbf{C}^{(i)}, \mathbf{X}^{(i-1)})}{P_{x_m}(W_i, X_i | \mathbf{W}^{(i-1)}, \mathbf{C}^{(i)}, \mathbf{X}^{(i-1)})}.$$
(F.6)

We note that this quantity can be easily estimated when  $W_i$  and  $X_i$  are low dimensional discrete variables.

To relax Assumption 9, we re-define the nuisance  $\pi_0(C_1, X_1, W)$  in Def. 8 as follow: For  $i = 1, 2, \dots, m-1$ ,

$$\pi_0^i \coloneqq \pi_0^i(\mathbf{W}^{(i)}, \mathbf{C}^{(i)}, \mathbf{X}^{(i)}) \coloneqq \frac{P_{x_i}(W_i, C_i | \mathbf{W}^{(i-1)}, \mathbf{C}^{(i-1)}, \mathbf{X}^{(i-1)})}{P_{x_m}(W_i, C_i, X_i | \mathbf{W}^{(i-1)}, \mathbf{C}^{(i-1)}, \mathbf{X}^{(i-1)})}.$$
(F.7)

In this subsection, let  $T^{\text{pw}}$  and  $T^{\text{dml}}$  denote PW and DML estimators for the AC-MTI functional in Def. 9 equipped with the re-defined nuisance  $\pi^i$  estimated for  $\pi_0^i$  in Eq. (F.7). Then, all error analysis results in Thm. 6 is preserved:

**Theorem S.2** (Error analysis of AC-MTI estimators). Let  $T^{pw}$  and  $T^{dml}$  denote PW and DML estimators for the AC-MTI functional in Def. 9 equipped with the re-defined nuisance  $\pi^i$  estimated for  $\pi_0^i$  in Eq. (F.7). Under Assumptions (2,7,8) and AC-MTI in Def. 7, the errors of the estimators in Def. 9, denoted  $\epsilon^{est} := T^{est} - \mathbb{E}[Y|do(\mathbf{x})]$  for  $est \in \{pw, dml\}$ , are:

$$\epsilon^{pw} = R_m + O_{P_{x_m}}(\|\pi^{(m-1)} - \pi_0^{(m-1)}\|),$$
  
$$\epsilon^{dml} = \sum_{i=1}^m R_i + \sum_{i=2}^m O_{P_{x_i}}(\|\mu^i - \mu_0^i\|\|\pi^{i-1} - \pi_0^{i-1}\|)$$

where  $R_i$  is a random variable such that  $\sqrt{n_i}R_i$  converges in distribution to a zero-mean normal random variable, where  $n_i := |D_{x_i}|$  for  $i \in \{1, \dots, m\}$ .

**Proof of Theorem S.2.** For the error analysis of  $T^{pw}$  with the redefined nuisance  $\pi_0^j$  in Eq. (F.7), it suffices to show that

$$\mathbb{E}_{P_{x_i}}\left[\pi_0^{(m-1)}\mathbb{1}_{\mathbf{x}^{(m-1)}}(\mathbf{X}^{(m-1)})Y\right] = \mathbb{E}\left[Y|do(\mathbf{x})\right],\tag{F.8}$$

because the other remaining parts of the proof are the same as the one for the original  $T^{pw}$ . It can be witnessed as follows: By Eq. (C.4),

$$\mathbb{E}[Y|do(\mathbf{x})] = \mathbb{E}_{P_{x_i}}\left[\prod_{j=1}^{m-1} \frac{P_{x_j}(\mathbf{C}^{(j)}, \mathbf{W}^{(j)}, \mathbf{X}^{(j-1)})}{P_{x_{j+1}}(\mathbf{C}^{(j)}, \mathbf{W}^{(j)}, \mathbf{X}^{(j)})} \mathbb{1}_{\mathbf{x}^{(m-1)}}(\mathbf{X}^{(m-1)})Y\right].$$
(F.9)

Also,

$$\prod_{j=1}^{m-1} \frac{P_{x_j}(\mathbf{C}^{(j)}, \mathbf{W}^{(j)}, \mathbf{X}^{(j-1)})}{P_{x_{j+1}}(\mathbf{C}^{(j)}, \mathbf{W}^{(j)}, \mathbf{X}^{(j)})} = \frac{1}{P_{x_m}(\mathbf{C}^{(m-1)}, \mathbf{W}^{(m-1)}, \mathbf{X}^{(m-1)})} \prod_{j=1}^{m-1} P_{x_j}(W_j, C_j | \mathbf{W}^{(j-1)}, \mathbf{C}^{(j-1)}, \mathbf{X}^{(j-1)})$$
(F.10)

$$=\prod_{j=1}^{m-1} \frac{P_{x_j}(W_j, C_j | \mathbf{W}^{(j-1)}, \mathbf{C}^{(j-1)}, \mathbf{X}^{(j-1)})}{P_{x_m}(W_j, C_j, X_j | \mathbf{W}^{(j-1)}, \mathbf{C}^{(j-1)}, \mathbf{X}^{(j-1)})}$$
(F.11)

$$=\prod_{j=1}^{m-1} \pi_0^j(W_j, C_j, X_j).$$
(F.12)

Therefore, Eq. (F.8) is witnessed.

For the error analysis of  $T^{dml}$  with the redefined nuisance  $\pi_0$  in Eq. (F.7), it suffices to show that Eq. (C.9) holds with the redefined nuisance; i.e.,

$$\omega_0^i(\mathbf{C}^{(i)}, \mathbf{W}^{(i)}, \mathbf{X}^{(i-1)}) \coloneqq \frac{P_{x_i}(\mathbf{C}^{(i)}, \mathbf{W}^{(i)}, \mathbf{X}^{(i-1)})}{P_{x_m}(\mathbf{C}^{(i)}, \mathbf{W}^{(i)}, \mathbf{X}^{(i-1)})} = \pi_0^i(\mathbf{C}^{(i)}, \mathbf{W}^{(i)}, \mathbf{X}^{(i)}) P_{x_m}(X_i | \mathbf{C}^{(i)}, \mathbf{W}^{(i)}, \mathbf{X}^{(i-1)}), \quad (F.13)$$

because the other remaining parts of the proof are the same as the one for the original  $T^{dml}$ . It can be witnessed as follows:

$$\begin{split} \omega_{0}^{i}(\mathbf{C}^{(i)},\mathbf{W}^{(i)},\mathbf{X}^{(i-1)}) &= \frac{P_{x_{i}}(C_{i},W_{i}|\mathbf{C}^{(i-1)},\mathbf{W}^{(i-1)},\mathbf{X}^{(i-1)})}{P_{x_{m}}(C_{i},W_{i}|\mathbf{C}^{(i-1)},\mathbf{W}^{(i-1)},\mathbf{X}^{(i-1)})} \frac{P_{x_{i}}(\mathbf{C}^{(i-1)},\mathbf{W}^{(i-1)},\mathbf{X}^{(i-1)})}{P_{x_{m}}(\mathbf{C}^{(i-1)},\mathbf{W}^{(i-1)},\mathbf{X}^{(i-1)})} \\ &= \frac{P_{x_{i}}(C_{i},W_{i}|\mathbf{C}^{(i-1)},\mathbf{W}^{(i-1)},\mathbf{X}^{(i-1)})}{P_{x_{m}}(C_{i},W_{i}|\mathbf{C}^{(i-1)},\mathbf{C}^{(i-1)},\mathbf{X}^{(i-1)})} \\ &= \frac{P_{x_{i}}(W_{i},C_{i}|\mathbf{W}^{(i-1)},\mathbf{C}^{(i-1)},\mathbf{X}^{(i-1)})}{P_{x_{m}}(W_{i},C_{i},X_{i}|\mathbf{W}^{(i-1)},\mathbf{C}^{(i-1)},\mathbf{X}^{(i-1)})} P_{x_{m}}(X_{i}|\mathbf{W}^{(i)},\mathbf{C}^{(i)},\mathbf{X}^{(i-1)}) \\ &= \pi_{0}^{i}(\mathbf{C}^{(i)},\mathbf{W}^{(i)},\mathbf{X}^{(i)})P_{x_{m}}(X_{i}|\mathbf{C}^{(i)},\mathbf{W}^{(i)},\mathbf{X}^{(i-1)}). \end{split}$$

In summary, the purpose of Assumption 9 is to simplify the estimation process by avoiding the need for joint density estimation of  $C_i, W_i, X_i$  in Eq. (F.7). Instead, it allows us to use the nuisance in Eq. (F.6), which is more amenable to estimate.