

# Nonparametric Causal Decomposition of Group Disparities\*

Ang Yu<sup>†</sup>      Felix Elwert<sup>‡</sup>

April 26, 2023

## Abstract

We propose a nonparametric method for decomposing group disparities in terms of an intermediate treatment variable. Our decomposition contains four components capturing the contributions of group differences in baseline characteristics, treatment prevalence, average treatment effect, and selection into treatment. These components readily inform policy interventions on the treatment and its effect. Our main contribution is to neatly separate the roles of between-group and within-group effect heterogeneity in explaining group disparities, revealing a new lever for reducing disparities. This decomposition reformulates the classic Kitagawa-Blinder-Oaxaca decomposition in causal and model-free terms, supplements causal mediation analysis by targeting group disparities instead of group effects, and resolves conceptual difficulties of recent random equalization decompositions. We also provide a conditional decomposition that allows researchers to incorporate pre-treatment covariates in defining the estimands and the hypothetical interventions. Using efficient influence functions and cross-fitting, we propose nonparametric estimators that are  $\sqrt{n}$ -consistent, asymptotically normal, semiparametrically efficient, and doubly robust. We apply our framework to study the causal role of education in intergenerational income mobility.

---

\*We thank Paul Bauer, Eric Grodsky, Aleksei Opacic, Guanghui Pan, Chan Park, Ben Rosche, Jiwei Zhao, and especially Xiang Zhou for helpful comments and suggestions. Earlier versions of this paper has been presented at Causality in the Social Sciences Workshop at GESIS in 2021, Population Association of America annual meeting in 2022, and American Causal Inference Conference in 2022. We thank participants at these conferences for constructive discussions.

<sup>†</sup>Department of Sociology, University of Wisconsin-Madison. Email: ayu33@wisc.edu

<sup>‡</sup>Department of Sociology and Department of Biostatistics and Medical Informatics, University of Wisconsin-Madison

# 1 Introduction

Improving the explanations of inter-group disparities is an urgent task for the social and health sciences. For example, what is the contribution of educational attainment to racial disparities in adult health? How much of the gender wage gap can be explained by occupational sorting? And what is the role of college completion in the relationship between parental income and income during adulthood? The structure of such question is always the same. Researchers compare the outcomes of two groups whose membership is typically established early in life (e.g., race, gender, parental background, birth cohort, country) and ask how differences in outcomes can be attributed to some intermediate treatment variable that occurs between the establishment of group membership and the outcome.

Questions about explanations, contributions, and attribution are inherently causal. First, researchers often are interested in the *causes* of disparities and ask counterfactual questions about the impact of changing one or more intermediate treatment variables, which is necessary for informing real-world interventions that could alleviate the disparities. It is one thing to observe that racial differences in college completion coexist with wage gaps by college completion status; it is another to establish that differences in college completion are responsible for racial wages gaps, such that intervening to equalize college graduation rates would remediate wage inequality. Second, researchers often also posit counterfactual scenarios where the effects of a treatment variable in a group were set to different values than they are in reality. At a minimum, the discussion of effects also presumes causal identification.

Three prior approaches are concerned with questions of attribution and explanation. First, the Kitagawa-Blinder-Oaxaca (KBO) decomposition ([Kitagawa, 1955](#); [Blinder, 1973](#); [Oaxaca, 1973](#)) is a classic tool used for explaining group disparities and remains popular to this day (e.g., [Barber et al., 2015](#); [Mize, 2016](#); [Laurison and Friedman, 2016](#); [Storer et al., 2020](#); [Zajacova et al., 2021](#)). The KBO decomposition is based on group-specific linear regression models and decomposes the group disparity into an explained component and an unexplained component. The explained component reflects the differential averages of an intermediate variable by group, and the unexplained component captures differential slopes and intercepts. The KBO decomposition is thus defined in terms of regression coefficients and not nonparametrically formulated in causal terms. This renders the KBO decomposition inappropriate for explanation and attribution questions. In fact, sociologists and economists alike have criticized the KBO decomposition for failing to answer any counterfactual question by design ([Fortin et al., 2011](#),

p.13; [Lundberg et al., 2021](#), p.542).

Second, causal mediation analysis is also used to attribute the relationship between two variables to an intermediate variable ([VanderWeele, 2015](#)). Causal mediation estimands are clearly defined in counterfactual terms, however, they are not designed for explaining group disparities. This is because what they decompose are summary measures of treatment effects, such as the average treatment effect (ATE). However, the discussion of treatment effects is conceptually challenging when it comes to ascriptive variables such as race and gender, as it is not plausible to impose any intervention on these variables ([Rubin, 1974](#)). Some group variables, such as early-life socioeconomic status, can theoretically be intervened upon, but not for people who are already later in their life course. Indeed, in numerous applications in social and health sciences, descriptive group disparities, rather than average effects, is what researchers seek to understand and explain. Causal attribution of group disparities thereby requires analytical frameworks more tailored to it than mediation analysis.

Third, there have recently emerged variants of a causal decomposition of group disparities, which we call the random equalization decomposition ([VanderWeele and Robinson, 2014](#); [Jackson and VanderWeele, 2018](#); [Lundberg, 2022](#)). For this decomposition, the extent to which the group disparity can be attributed to the differential prevalence of an intermediate treatment variable is defined as the reduction in disparity brought about by a hypothetical intervention that randomly reassigns values of the treatment such that its prevalence is equalized across groups. The random equalization decomposition marks a important advancement towards establishing an appropriate framework for causal attribution of group disparities to an intermediate treatment variable. However, as we will show, the random equalization decomposition does not purely capture the contribution of differential prevalence of the treatment. Instead, it also incorporates the influence of differential patterns of selection into treatment across groups, which complicates its interpretation.

In this article, we develop a novel decomposition approach for group disparities in terms of an intermediate treatment variable. Our decompositions are defined counterfactually with respect to the treatment variable and descriptively with respect to the group variable. This approach is consistent with common practices in the social sciences and focuses on more modifiable factors in terms of policy intervention. Our unconditional decomposition contains four components, respectively capturing the contributions of baseline differences unrelated to the treatment, differential treatment prevalence, differential ATE, and differential selection into

the treatment. Compared with the KBO decomposition, our decomposition enables causal attribution and nonparametric estimation. It supplements causal mediation analysis by focusing on explaining group disparities instead of “effects” of the group variables. It also improves on the random equalization decomposition by providing components with unambiguous interpretations. The component in our decomposition that captures differential selection is a centerpiece in our contributions. First, this component reveals a novel policy lever for reducing disparities, and it did not appear in any of the previous decompositions. Second, by separating it from the component representing differential prevalence, we clarify the random equalization decomposition and provide a more clearly interpretable decomposition.

Our paper proceeds as follows. Section 2 introduces our unconditional decomposition and its interventional interpretations. We also formally relate our decomposition to the classic KBO decomposition, a causal mediation estimand, and the random equalization decomposition, through which we explicate the contributions of our framework. We also extend the unconditional decomposition to a conditional version, which allows researchers to take pre-treatment covariates into account in defining the estimands and the hypothetical interventions. In section 3, under the conditional ignorability assumption, we develop nonparametric estimators for the unconditional and conditional decompositions based on their efficient influence functions. These estimators can be implemented via data-adaptive methods such as machine learning (ML) models and accommodate high-dimensional confounders. We also derive the conditions under which the estimators are  $\sqrt{n}$ -consistent, asymptotically normal, and semi-parametrically efficient. In addition, they have double-robustness properties. In section 4, we apply our decompositions to study the causal role of college graduation in intergenerational income mobility. Section 5 concludes with possible future extensions.

## 2 Estimands

### 2.1 Unconditional decomposition

For each individual  $i$ , consider a binary treatment variable  $D_i \in \{0, 1\}$ . Let  $Y_{i0}$  and  $Y_{i1}$  be the potential outcomes (Rubin, 1974) of  $Y_i$  under the hypothetical intervention to set  $D_i = 0$  and  $D_i = 1$ , respectively. Let  $\tau_i := Y_{i1} - Y_{i0}$  denote the individual-level treatment effect. Henceforth, we largely suppress individual level subscripts on all variables to unburden notation. Suppose that the population contains two disjoint groups,  $g \in \{a, b\}$ , where  $a$  denotes the advantaged

group and  $b$  denotes the disadvantaged group. We use subscripts  $g$  to indicate group-specific quantities, for example,  $E_g(Y) := E(Y \mid G = g)$ .

We now only assume the stable unit treatment value assumption (SUTVA) ([Rubin, 1980](#)),

Assumption 1 (SUTVA).  $Y = DY_1 + (1 - D)Y_0$ .

Then, the observed outcomes disparity between group  $a$  and  $b$  can be decomposed into four components:

$$\begin{aligned}
& E_a(Y) - E_b(Y) \\
&= \underbrace{E_a(Y_0) - E_b(Y_0)}_{\text{baseline}} + \underbrace{E_b(\tau)[E_a(D) - E_b(D)]}_{\text{prevalence}} \\
&\quad + \underbrace{E_a(D)[E_a(\tau) - E_b(\tau)]}_{\text{effect}} + \underbrace{\text{Cov}_a(D, \tau) - \text{Cov}_b(D, \tau)}_{\text{selection}}. \tag{1}
\end{aligned}$$

The “baseline” component reflects the difference in average baseline potential outcome  $Y_0$  between groups, i.e., the difference in outcomes if nobody received treatment. The “prevalence” component indicates how much of the group disparity is due to differential prevalence of treatment. The “effect” component reflects unequal average treatment effects (ATE) across groups. Finally, the “selection” component captures the contribution of the differential selection into treatment based on the treatment effect,  $\tau$ , to the group disparity. For each group, if group members who would gain more from treatment are more likely to receive the treatment than those who would gain less, there will be a positive covariance between  $D$  and  $\tau$ . To our knowledge, none of prior decompositions of group disparities includes a selection component.<sup>1</sup> Both the effect and the selection components account for the relationship between heterogeneous effects and group disparities. But the “effect” component captures the contribution of *between*-group effect heterogeneity; and the “selection” component captures the contribution of *within*-group effect heterogeneity. In summary, a group will be more advantaged in the outcome if its members have a higher average baseline potential outcome, a higher prevalence of treatment (given a positive ATE), a higher ATE, or a higher level of selection into treatment based on treatment effect.

---

<sup>1</sup>[Zhou \(2022\)](#) recently developed a causal decomposition method for a special case of mediation analysis, which contains a similar covariance component capturing selection into a mediator.

## 2.2 Interventional interpretations

Our decomposition is formulated in counterfactual terms, hence is prescriptive for future interventions. The decomposition reveals three levers by which a policy may affect inter-group disparities through the treatment. First, policy makers could manipulate the prevalence of the treatment in each of the two groups. Second, they could manipulate who within each group receives treatment. For example, they could facilitate the matching of high-return individuals to treatment receipt in the disadvantaged group. Third, they might even be able to manipulate the average effect of the treatment for each group.<sup>2</sup> Specifically, each of our decomposition component quantifies the part of the outcome disparity that could be eliminated if the two groups were equalized with respect to that component while other characteristics were held unchanged.

To further fix ideas, we may more exactly pin down the equalization intervention implicated by each decomposition component. The prevalence component indicates the part of the disparity that would be eliminated by giving group  $b$  the treatment prevalence of group  $a$ , which can be seen by replacing  $E_b(D)$  with  $E_a(D)$  in equation (1). In the same vein, the effect component is the part of the disparity that would disappear if group  $a$  was given the ATE of group  $b$ . Finally, the disparity would be reduced by the amount of the selection component if the two groups were made to have the same selection level, regardless of the specific level at which their selection would be fixed. However, given that the prevalence, effect, and selection components are determined only by two variables, the treatment and the treatment effect, it is not always possible to intervene on a single component without changing the values of the other components. Nevertheless, we may alternatively conceive of a whole-package intervention. A whole-package intervention transfers the joint distribution of  $D$  and  $\tau$  from one group to the other. Our prevalence, effect, and selection components then add to the reduction in disparity of this whole-package intervention, and each of them can be interpreted as a distinct aspect of the intervention.

Even if one is not willing to consider interventions on treatment effects, the decomposition components can still be interpreted in terms of a two-step intervention. Under the randomized intervention notation proposed by [Didelez et al. \(2006\)](#),  $R(D \mid G = g')$  is defined to

---

<sup>2</sup>It is more unconventional to discuss interventions on treatment effects than on treatments. However, interventions on effects have appeared in both methodological ([Malinsky, 2018](#)) and empirical ([Brady et al., 2017](#)) literatures. Our notion of interventions on effects is a generalization of Malinsky’s (2018) interventions on “structural features”, as we allow for heterogeneous effects.

be a randomly drawn value of treatment from group  $g'$ . Then,  $E_g(Y_{R(D|G=g')})$  denotes the mean outcome for group  $g$  after each member of group  $g$  is counterfactually given a treatment value randomly drawn from group  $g'$ . When  $g = g'$ , the intervention amounts to a random permutation of treatment status within the group. Using its definition, we can rewrite the post-intervention mean as follows.

$$E_g(Y_{R(D|G=g')}) = E_g(Y_0) + E_{g'}(D) E_g(\tau). \quad (2)$$

Then it follows that our decomposition components can be re-written in the randomized intervention notation.

$$\begin{aligned} E_a(Y) - E_b(Y) - [E_a(Y_{R(D|G=a)}) - E_b(Y_{R(D|G=b)})] &= \text{selection} \\ E_b(Y_{R(D|G=a)}) - E_b(Y_{R(D|G=b)}) &= \text{prevalence} \\ E_a(Y_{R(D|G=a)}) - E_b(Y_{R(D|G=a)}) &= \text{baseline} + \text{effect} \end{aligned}$$

This points to a two-step intervention. In the first step, the pre-intervention disparity is  $E_a(Y) - E_b(Y)$ , and the post-intervention disparity is  $E_a(Y_{R(D|G=a)}) - E_b(Y_{R(D|G=b)})$ . Hence, this step internally randomizes the treatment in both groups without changing its prevalence, and the selection component represents the reduction in disparity resulting from this randomization intervention. In the second step, the pre-intervention disparity is  $E_a(Y_{R(D|G=a)}) - E_b(Y_{R(D|G=b)})$ , while the post-intervention disparity is  $E_a(Y_{R(D|G=a)}) - E_b(Y_{R(D|G=a)})$ . This step equalizes treatment prevalence by giving group  $b$  the treatment prevalence of group  $a$ . Therefore, the prevalence component is the reduction in disparity of an equalization intervention. Importantly, this equalization intervention is independent of the randomization intervention in the first step. Finally, the remaining disparity at the end of the two-step intervention is the sum of the baseline and the effect components.

## 2.3 Relation to the KBO decomposition

Our decomposition expands on the classic KBO decomposition in two ways. First, our decomposition is presented explicitly as a causal estimand, which can be estimated nonparametrically and interpreted interventionally. By contrast, the KBO decomposition is defined in terms of linear regressions and coincides with our decomposition only under a very strong assumption.

Second, our decomposition contains an additional selection component as a novel disparity-producing mechanism, which is assumed away in the KBO decomposition.

Consider a form of the KBO decomposition of the outcome disparity between groups  $a$  and  $b$  with respect to a treatment variable,  $D$ ,

$$\begin{aligned} & \bar{Y}_a - \bar{Y}_b \\ &= \underbrace{\hat{\alpha}_a - \hat{\alpha}_b}_{\text{intercept}} + \underbrace{\hat{\beta}_b[\bar{D}_a - \bar{D}_b]}_{\text{characteristic}} + \underbrace{\bar{D}_a[\hat{\beta}_a - \hat{\beta}_b]}_{\text{slope}}, \end{aligned}$$

which is based on estimates of the group-specific linear regressions:

$$Y = \alpha_g + \beta_g D + \epsilon.$$

Under a strong ignorability assumption,  $Y_d \perp\!\!\!\perp D \mid G = g, \forall d, g$ , in addition to assumption 1, the intercept, characteristic, and slope components in the KBO decomposition are unbiased estimators for our baseline, prevalence, and effect components, respectively. However, the strong ignorability assumption requires random assignment of treatment conditional only on group membership. It also directly rules out the selection component in group disparities. Hence, our decomposition is both more expansive and interventionally interpretable under far less stringent assumptions. Although components of the KBO decomposition have commonly been given counterfactual and interventional interpretations (e.g., [Jann, 2008](#)), our decomposition should be preferred for its more plausibly causal interpretability.

Before the current article, several studies have constructed or inferred causal estimands for the KBO decomposition. A methodological literature centers around interpreting the unexplained component (the sum of the intercept and the slope components) as an estimator for the treatment effect on the treated (ATT) ([Fortin et al., 2011](#); [Kline, 2011](#); [Yamaguchi, 2015](#)).<sup>3</sup> This literature addresses how an outcome difference associated with a treatment variable can be accounted for by *pre*-treatment confounders. By contrast, our decomposition focuses on a very different question, i.e., the extent to which an outcome difference between groups is explained by a *post*-group variable. [Huber \(2015\)](#) proposes an interpretation of the KBO components as causal mediation estimands, where the characteristic component is interpreted as the natural indirect effect ([Pearl, 2001](#)) of the group indicator via an intermediate variable

---

<sup>3</sup>Also see [Chernozhukov et al. \(2018b\)](#), which interprets a version of the unexplained component as an average partial effect.



and the unexplained component as the corresponding natural direct effect. The next section contrasts our decomposition with causal mediation analysis.

## 2.4 Relation to randomized interventional analogues of mediation estimands

Among the many causal mediation estimands, a four-way decomposition based on randomized interventional analogues (RIAs) of mediation estimands (VanderWeele, 2015, p.619-21; see also VanderWeele and Tchetgen Tchetgen, 2017; VanderWeele et al., 2014) is most akin to our decomposition. By showing how our decomposition is related to but differs from this decomposition, we highlight the unique advantages of our decomposition compared with all causal mediation estimands. A form of this decomposition can be written as

$$\begin{aligned}
& \underbrace{\mathbb{E}(Y_{G=a,R(D_{G=a})}) - \mathbb{E}(Y_{G=b,R(D_{G=b})})}_{\text{RIA of total effect}} \\
&= \underbrace{\mathbb{E}(Y_{G=a,D=0}) - \mathbb{E}(Y_{G=b,D=0})}_{\text{conditional direct effect (CDE)}} + \underbrace{\mathbb{E}(Y_{G=b,R(D_{G=a})}) - \mathbb{E}(Y_{G=b,R(D_{G=b})})}_{\text{RIA of pure indirect effect (PIE)}} \\
&+ \underbrace{\mathbb{E}(Y_{G=a,D=1} - Y_{G=a,D=0} - Y_{G=b,D=1} + Y_{G=b,D=0}) \mathbb{E}(D_{G=b})}_{\text{RIA of reference interaction effect (RIE)}} \\
&+ \underbrace{\mathbb{E}(Y_{G=a,D=1} - Y_{G=a,D=0} - Y_{G=b,D=1} + Y_{G=b,D=0}) (\mathbb{E}(D_{G=a}) - \mathbb{E}(D_{G=b}))}_{\text{RIA of mediated interaction effect (MIE)}},
\end{aligned}$$

where  $\mathbb{E}(Y_{G=g,R(D_{G=g'})})$  is the mean potential outcome of  $Y$  under assigning group  $g$  and a randomly drawn value of  $D$  from the population when assigned group  $g'$ , and  $\mathbb{E}(D_{G=g})$  is the the mean potential outcome of  $D$  under assigning group  $g$ .

We first note the connections between our decomposition and the RIA-based decomposition. Under unconditional ignorability of  $G$ , i.e.,  $Y_{G=g,D=d} \perp\!\!\!\perp G, \forall d, g$ , and  $D_{G=g} \perp\!\!\!\perp G, \forall g$ , and two SUTVA-type assumptions,  $\mathbb{E}(Y_{G=g,D=d} | G = g) = \mathbb{E}(Y_{D=d} | G = g)$  and  $\mathbb{E}(D_{G=g} | G = g) = \mathbb{E}(D | G = g), \forall d, g$ , the CDE equals our baseline component, the RIA of PIE equals our prevalence component, and the sum of RIAs of RIE and MIE equals our effect component. These connections are intuitive: both the CDE and the baseline component capture a group-based outcome difference when the intermediate treatment variable is held at 0; both the RIA of PIE and the prevalence component address the role of the treatment in the relationship between the group and the outcome; also, the RIA of RIE, the RIA of MIE, and the effect

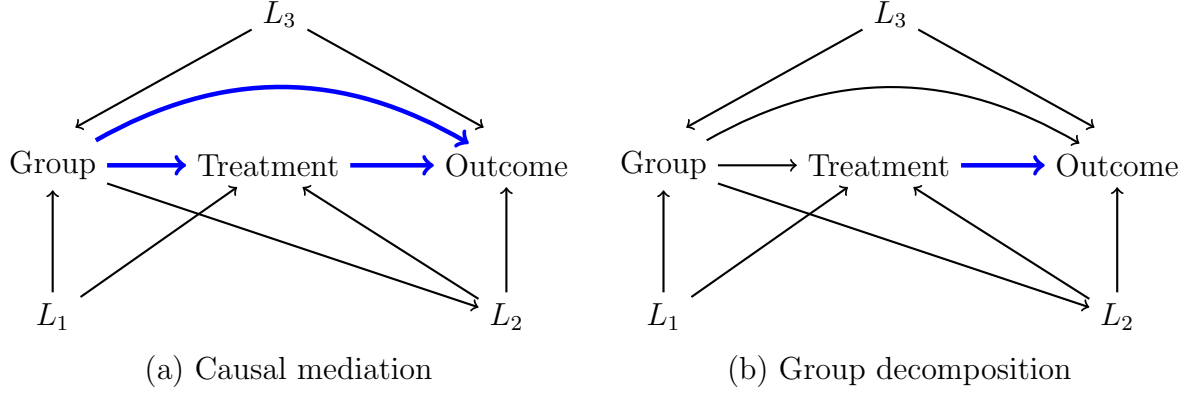


Figure 2. Causal graphs illustrating the difference in identification goals between causal mediation analysis and causal decomposition of group disparities. The highlighted edges are the paths each estimand seeks to identify. Causal mediation analysis aims to identify both the effects of group and treatment, hence requiring proper handling of all three types of confounders:  $L_1$ ,  $L_2$ , and  $L_3$ . Our causal decomposition only needs to identify the effect of treatment, so only treatment-outcome confounders, i.e.,  $L_2$ , are of concern.

component all, at least partly, reflect the heterogeneous effects of the treatment by group membership. Thus, as tools of causal explanation, our decomposition and causal mediation analysis are akin to each other.<sup>4</sup>

However, our decomposition and causal mediation analysis differ in two important ways, which makes our decomposition advantageous for studying group disparities. First, we decompose a descriptive group disparity rather than a total effect of group membership. This does not only cater to the goal of studying group disparities but also requires fewer identifying assumptions than causal mediation analysis. We illustrate the difference in identifying assumptions in Figure 1. Consequently, although the components of the RIA-based decomposition can be forcibly equated to the components of our decomposition, it would require the impossible scenario where the group membership is randomly assigned. Second, there is no decomposition in the causal mediation literature that contains a selection component. In this regard, our decomposition makes a novel contribution.

## 2.5 Relation to unconditional random equalization

In recent years, there have emerged two variants of what we call the unconditional random equalization decomposition, which decomposes the observed disparity into two components.

<sup>4</sup>These intuitions also carry over to VanderWeele’s (2014) four-way decomposition that is not based on RIAs. However, in order to establish a similar connection between the conventional four-way decomposition and our decomposition, a cross-world independence assumption needs to be additionally imposed.

One is the reduction in disparity that could be brought about by randomly reassigning the treatment such that the treatment prevalence is equalized across groups. The other is correspondingly a component of remaining disparity (VanderWeele and Robinson, 2014; Jackson and VanderWeele, 2018; Sudharsanan and Bijlsma, 2021; Lundberg, 2022). Like our approach, the random equalization decomposition seeks to identify the causal effect of the treatment, not that of group membership. However, the random equalization decomposition is a two-way decomposition that contains less information than our four-way decomposition. Moreover, the reduction in disparity in the random equalization decomposition combines the selection component and the prevalence component of our decomposition, impeding a clear interpretation.

The first variant of the random equalization decomposition is defined in terms of a random equalization intervention that randomly assigns treatment values of the advantaged group to the disadvantaged group. Using the randomized intervention notation introduced above, Jackson and VanderWeele (2018) decompose the observed group disparity into two components:

$$E_a(Y) - E_b(Y) = \underbrace{E_b(Y_{R(D|G=a)}) - E_b(Y)}_{\text{reduction in disparity}} + \underbrace{E_a(Y) - E_b(Y_{R(D|G=a)})}_{\text{remaining disparity}}.$$

The reduction in disparity is intended to quantify the contribution of the differential prevalence of treatment across groups. However, the random equalization intervention not only equalizes treatment prevalence across groups, but also randomizes the treatment within the disadvantaged group, making any selection into treatment disappear in that group. As a consequence, the reduction in disparity equals the combination of the prevalence component and the group  $b$  part of the selection component, i.e.,  $E_b(\tau)[E_a(D) - E_b(D)] - \text{Cov}_b(D, \tau)$ , which follows from equation (2). This leads to an interpretational difficulty that is reflected in an undesirable result: If group  $a$  is an exact replicate of group  $b$ , the random equalization decomposition would still contain a non-zero reduction in disparity, as long as  $\text{Cov}_b(D, \tau) \neq 0$ .

The second variant of the random equalization decomposition takes a somewhat different form (Lundberg, 2022), whose hypothetical intervention assigns each individual in both groups a treatment value randomly drawn from the pooled population. Hence, this random equalization intervention changes the treatment values in both groups instead of only in the disadvantaged group as in the first variant. The reduction in disparity of this variant mixes the prevalence component with the selection component, too. To show this, we first note that

the outcome disparity can also be decomposed as such:

$$\begin{aligned}
& E_a(Y) - E_b(Y) \\
&= E_a(Y_0) - E_b(Y_0) + E(\tau)[E_a(D) - E_b(D)] \\
&\quad + E(D)[E_a(\tau) - E_b(\tau)] + \text{Cov}_a(D, \tau) - \text{Cov}_b(D, \tau) \\
&\quad - [P_a - P_b][E_a(D) - E_b(D)][E_a(\tau) - E_b(\tau)],
\end{aligned}$$

where  $E(\tau)$  and  $E(D)$  are the overall ATE and treatment prevalence,  $P_a$  and  $P_b$  are the proportions of the population in group  $a$  and group  $b$ . And the remaining disparity in Lundberg’s (2022) unconditional decomposition equals  $E_a(Y_0) - E_b(Y_0) + E(D)[E_a(\tau) - E_b(\tau)]$ . It then follows that the reduction in disparity in this case contains the selection component,  $\text{Cov}_a(D, \tau) - \text{Cov}_b(D, \tau)$ , which does not involve the difference in treatment prevalence. The intuition is that in the random equalization intervention, both group  $a$  and group  $b$  would get random draws of  $D$ , thus selection into treatment is eliminated in both groups, making  $\text{Cov}_a(D, \tau) = \text{Cov}_b(D, \tau) = 0$ .

To conclude, the reduction in disparity in both variants of the random equalization decomposition is a mixture of the prevalent component and the selection component of group disparities. This is because its underlying intervention both equalizes and randomizes the treatment. In contrast, as explained in Section 2.1, the prevalence component in our decomposition is solely concerned with the disparity-reducing impact of an equalization intervention and avoids the influence of a randomization intervention.<sup>5</sup> If the goal is to capture the contribution of the unequal treatment prevalence to group disparities, one should prefer using the prevalence component in our decomposition.

## 2.6 Conditional decomposition

In this subsection, we develop a decomposition method whose corresponding interventions are conditional on a vector of variables  $\mathbf{Q}$ . The conditional decomposition is useful as it can shed light on the impact of normatively more desirable or realistically more feasible interventions, when the unconditional intervention appears less so (Jackson, 2021). For example, it may be

---

<sup>5</sup>This is also why the RIA of PIE can be translated to the prevalence component, but not to the reduction in disparity component in random equalization decompositions. Both the prevalence component and the RIA of PIE are a difference between two quantities that both contain randomized  $D$  values, so both estimands can represent the reduction in disparity resulting from a pure equalization intervention *after* a randomization intervention.

meaningful to consider equalizing college admissions conditional on test scores or equalizing drug prescriptions conditional on certain comorbidities. Assuming that the support of  $\mathbf{Q}$  is the same in the two groups and still under Assumption 1,

$$\begin{aligned}
E_a(Y) - E_b(Y) &= \underbrace{E_a(Y_0) - E_b(Y_0)}_{\text{baseline}} \\
&+ \underbrace{\int [E_a(D \mid \mathbf{Q} = \mathbf{q}) - E_b(D \mid \mathbf{Q} = \mathbf{q})] E_b(\tau \mid \mathbf{Q} = \mathbf{q}) f_b(\mathbf{q}) d\mathbf{q}}_{\text{conditional prevalence}} \\
&+ \underbrace{\int [E_a(\tau \mid \mathbf{Q} = \mathbf{q}) - E_b(\tau \mid \mathbf{Q} = \mathbf{q})] E_a(D \mid \mathbf{Q} = \mathbf{q}) f_a(\mathbf{q}) d\mathbf{q}}_{\text{conditional effect}} \\
&+ \underbrace{E_a[\text{Cov}_a(D, \tau \mid \mathbf{Q})] - E_b[\text{Cov}_b(D, \tau \mid \mathbf{Q})]}_{\text{conditional selection}} \\
&+ \underbrace{\int E_a(D \mid \mathbf{Q} = \mathbf{q}) E_b(\tau \mid \mathbf{Q} = \mathbf{q}) [f_a(\mathbf{q}) - f_b(\mathbf{q})] d\mathbf{q}}_{\mathbf{Q} \text{ distribution}}. \tag{3}
\end{aligned}$$

The conditional prevalence, effect, and selection components indicate the disparity-reducing impacts of equalizing the corresponding group characteristics within levels of  $\mathbf{Q}$ . Thus, the conditional decomposition allows researchers to incorporate prior information about the subjects when designing a re-distribution of treatments or their effects. As policies, the conditional equalization interventions are often more feasible but less influential than the marginal equalizations encoded in equation (1). In addition, the “ $\mathbf{Q}$  distribution” component captures the between- $\mathbf{Q}$  disparity, as opposed to within- $\mathbf{Q}$  disparities reflected in other components, but it should not be interpreted causally, as we do not make any identifying assumptions for  $\mathbf{Q}$ .

Similar to the unconditional case, we can also interventionally interpret the conditional decomposition using the randomized intervention notation. Let  $E_g(Y_{R(D|G=g', \mathbf{Q})})$  be the mean potential outcome of group  $g$  when its members were given treatment values randomly drawn from members of group  $g'$  who share the same  $\mathbf{Q}$  values with them. We can rewrite this mean potential outcome as follows.

$$E_g(Y_{R(D|G=g', \mathbf{Q})}) = E_g(Y_0) + \int E_g(\tau \mid \mathbf{Q} = \mathbf{q}) E_{g'}(D \mid \mathbf{Q} = \mathbf{q}) f_g(\mathbf{q}) d\mathbf{q}. \tag{4}$$

It then follows that components of the conditional decomposition can be written in the ran-

domized intervention notation:

$$\begin{aligned}
E_a(Y) - E_b(Y) - [E_a(Y_{R(D|G=a, \mathbf{Q})}) - E_b(Y_{R(D|G=b, \mathbf{Q})})] &= \text{conditional selection} \\
E_b(Y_{R(D|G=a, \mathbf{Q})}) - E_b(Y_{R(D|G=b, \mathbf{Q})}) &= \text{conditional prevalence} \\
E_a(Y_{R(D|G=a, \mathbf{Q})}) - E_b(Y_{R(D|G=a, \mathbf{Q})}) &= \text{conditional effect} + \text{conditional selection} \\
&\quad + \mathbf{Q} \text{ distribution.} \tag{5}
\end{aligned}$$

Therefore, a two-step intervention again underlies our conditional decomposition. In the first step, the treatment is randomized within each  $\mathbf{Q}$  level for members of both groups. The reduction in disparity resulting from this step is then the conditional selection component. The second step is an equalization intervention independent of the previous randomization step. For the second step, the pre-intervention disparity is  $E_a(Y_{R(D|G=a, \mathbf{Q})}) - E_b(Y_{R(D|G=b, \mathbf{Q})})$ , and the post-intervention disparity is  $E_a(Y_{R(D|G=a, \mathbf{Q})}) - E_b(Y_{R(D|G=a, \mathbf{Q})})$ . Hence, the conditional prevalence component is the reduction in disparity brought about by giving group  $b$  random draws of treatment from group  $a$  within  $\mathbf{Q}$  levels. What would remain after this two-step intervention is the sum of the conditional effect, the conditional selection, and the  $\mathbf{Q}$  distribution components.<sup>6</sup>

Similar to the unconditional random equalization decompositions, [Jackson \(2021\)](#) and [Lundberg \(2022\)](#) proposed two variants of what can be called the conditional random equalization decomposition (CRED), which decomposes outcome disparities into a reduction in disparity component and a remaining disparity component. The hypothetical intervention underlying the CRED is a random re-distribution of treatment values among people with the same  $\mathbf{Q}$  values such that the treatment prevalence would be equalized within  $\mathbf{Q}$  levels. Analogous to the unconditional case, we show that the reduction in disparity component in both variants of CRED conditionally mixes the prevalence component with the selection component, thus the conditional prevalence component in our decomposition should be preferred for its clearer interpretation.

In Jackson's (2021) variant of the CRED, the reduction in disparity is  $E_b(Y_{R(D|G=a, \mathbf{Q})}) -$

---

<sup>6</sup>For equation (4) to hold, we only require  $\text{supp}_g(\mathbf{Q}) \subseteq \text{supp}_{g'}(\mathbf{Q})$ , not necessarily  $\text{supp}_g(\mathbf{Q}) = \text{supp}_{g'}(\mathbf{Q})$  required for the decomposition (3) itself. This implies that the interventional interpretations for conditional selection and conditional prevalence in equation (5) hold as long as  $\text{supp}_b(\mathbf{Q}) \subseteq \text{supp}_a(\mathbf{Q})$ . Intuitively, at each level of  $\mathbf{Q}$  in group  $b$ , we must be able to find members of group  $a$  with the same  $\mathbf{Q}$  values in order to conduct the equalization intervention.

$E_b(Y)$ , which we can rewrite as

$$\int [E_a(D \mid \mathbf{Q} = \mathbf{q}) - E_b(D \mid \mathbf{Q} = \mathbf{q})] E_b(\tau \mid \mathbf{Q} = \mathbf{q}) f_b(\mathbf{q}) d\mathbf{q} - E_b[\text{Cov}_b(D, \tau \mid \mathbf{Q})]. \quad (6)$$

Therefore, the reduction in disparity is a mixture of the conditional prevalence component and part of the conditional selection component. Namely, it captures not only the differential treatment prevalence within  $\mathbf{Q}$  levels, but also the extent of conditional selection into treatment in group  $b$ . Intuitively, this is because the intervention of this CRED involves both a conditional equalization and a conditional randomization, and the randomization part eliminates selection into treatment conditional on  $\mathbf{Q}$  in group  $b$ .

A similar case can be shown for the reduction in disparity in the CRED of [Lundberg \(2022\)](#).

$$\begin{aligned} & E_a(Y) - E_b(Y) \\ & \{E_a[\Pr(D = 0 \mid \mathbf{Q})Y_0 + \Pr(D = 1 \mid \mathbf{Q})Y_1] - E_b[\Pr(D = 0 \mid \mathbf{Q})Y_0 + \Pr(D = 1 \mid \mathbf{Q})Y_1]\} \\ & = E_a[\text{Cov}_a(D, \tau \mid \mathbf{Q})] - E_b[\text{Cov}_b(D, \tau \mid \mathbf{Q})] + \int [E_a(D \mid \mathbf{Q} = \mathbf{q}) - E_b(D \mid \mathbf{Q} = \mathbf{q})] \\ & [E_a(\tau \mid \mathbf{Q} = \mathbf{q}) f_a(\mathbf{q}) \Pr(G = b \mid \mathbf{Q} = \mathbf{q}) + E_b(\tau \mid \mathbf{Q} = \mathbf{q}) f_b(\mathbf{q}) \Pr(G = a \mid \mathbf{Q} = \mathbf{q})] d\mathbf{q}, \quad (7) \end{aligned}$$

where the left hand side is the reduction in disparity in Lundberg's (2022) CRED. Hence, the reduction in disparity again involves both a term for the differential treatment prevalence within  $\mathbf{Q}$  levels and the conditional selection component. The CRED in this case corresponds to the combination of a conditional equalization and a conditional randomization in both groups. Consequently, both the differential treatment prevalence and the differential selection into treatment would be reduced to zero conditional on  $\mathbf{Q}$ . Distinct from the reduction in disparity in both variants of the CRED, the conditional prevalence component in our decomposition solely measures the impact of a conditional equalization intervention. By implication, only our conditional prevalence component can be interpreted as a pure measure of the contribution of the differential treatment prevalence within  $\mathbf{Q}$  levels to outcome disparities.

### 3 Estimation and Inference

We identify the components of our decompositions using two standard assumptions as follows.<sup>7</sup>

---

<sup>7</sup>Our unconditional decomposition can also be identified using the marginal treatment effect approach based on an instrumental variable ([Heckman and Vytlacil, 2005](#); [Zhou and Xie, 2020](#)), as identifying group-specific

Assumption 2 (Conditional ignorability).  $Y_d \perp\!\!\!\perp D \mid \mathbf{X} = \mathbf{x}, G = g, \forall d, \mathbf{x}, g;$

Assumption 3 (Overlap).  $0 < E(D \mid \mathbf{X} = \mathbf{x}, G = g) < 1, \forall \mathbf{x}, g.$

Without loss of generality, we assume  $\mathbf{Q} \in \mathbf{X}$ . We develop nonparametric and efficient estimators for components of our unconditional and conditional decompositions.<sup>8</sup> These estimators are “one-step” estimators based on the EIFs of the decomposition components, which remove the bias from naive substitution estimators (Bickel et al., 1998; Van der Vaart, 2000; Hines et al., 2022). The estimators contain some nuisance functions, which can be estimated using flexible ML methods coupled with cross-fitting. Under specified conditions, our estimators are  $\sqrt{n}$ -consistent, asymptotically normal, and semiparametrically efficient. Thus, we are able to construct asymptotically accurate Wald-type confidence intervals and hypothesis tests. Our estimators also have double robustness properties.

To introduce the EIFs, we define the following functions:

$$\begin{aligned}\mu(d, \mathbf{X}, g) &= E(Y \mid D, \mathbf{X}, G = g) \\ \pi(d, \mathbf{X}, g) &= \Pr(D = d \mid \mathbf{X}, G = g) \\ p_g &= \Pr(G = g) \\ p_g(\mathbf{Q}) &= \Pr(G = g \mid \mathbf{Q}).\end{aligned}$$

Below, we use  $\xi$  to denote the estimands (decomposition components),  $\phi$  to denote the EIFs, and  $\Psi$  to denote the resulting one-step estimators.

### 3.1 The unconditional decomposition

All components of the unconditional decomposition are simple linear combinations of the total disparity and two generic functions evaluated at appropriate values of  $d$ ,  $g$ , and  $g'$ :  $\xi_{dg} := E(Y_d \mid G = g)$  and  $\xi_{dgg'} := E(Y_d \mid G = g) E(D \mid G = g')$ . The relationship between components of the unconditional decomposition and the generic functions are as follows:

$$\text{Baseline} = \xi_{0a} - \xi_{0b}$$

---

ATTs and ATEs is sufficient (see Appendix E).

<sup>8</sup>Various regression, weighting, or matching estimators can also be used. In particular, under constant treatment effect, the prevalence component in the unconditional decomposition can be estimated using the traditional “product method” for the indirect effect of  $G$  through  $D$  on  $Y$  (Baron and Kenny, 1986). This again explicates the kinship between our decomposition and mediation analysis.



$$\text{Prevalence} = \xi_{1ba} - \xi_{0ba} - \xi_{1bb} + \xi_{0bb}$$

$$\text{Effect} = \xi_{1aa} - \xi_{0aa} - \xi_{1ba} + \xi_{0ba}$$

$$\text{Selection} = \text{Total} - \text{Baseline} - \text{Prevalence} - \text{Effect}.$$

Hence, the EIFs and one-step estimators for the decomposition components directly follow from those of  $\xi_{dg}$  and  $\xi_{dgg'}$ . In addition, these generic functions also provide a basis for estimating Jackson's (2018) version of the unconditional equalization estimand, since its reduction in disparity component can be represented as  $\xi_{0b} + \xi_{1ba} - \xi_{0ba} - E_b(Y)$ . Under assumptions 1, 2, and 3,  $\xi_{dg}$  and  $\xi_{dgg'}$  can be identified as the following quantities:

$$\xi_{dg} = E[\mu(d, \mathbf{X}, g) \mid G = g]$$

$$\xi_{dgg'} = E[\mu(d, \mathbf{X}, g) \mid G = g] E[D \mid G = g'].$$

These identification results then enable us to derive the EIFs for  $\xi_{dg}$  and  $\xi_{dgg'}$ .

**Theorem 1** (EIF). *Under assumptions 1, 2, and 3, the EIF of  $\xi_{dg}$  is*

$$\phi_{dg}(Y, \mathbf{X}) := \frac{\mathbb{1}(G = g)}{p_g} \left\{ \frac{\mathbb{1}(D = d)}{\pi(d, \mathbf{X}, g)} [Y - \mu(d, \mathbf{X}, g)] + \mu(d, \mathbf{X}, g) - \xi_{dg} \right\},$$

and the EIF of  $\xi_{dgg'}$  is

$$\begin{aligned} \phi_{dgg'}(Y, \mathbf{X}) := & \frac{\mathbb{1}(G = g)}{p_g} \left\{ \frac{\mathbb{1}(D = d)}{\pi(d, \mathbf{X}, g)} [Y - \mu(d, \mathbf{X}, g)] + \mu(d, \mathbf{X}, g) \right\} E(D \mid G = g') \\ & + \frac{\mathbb{1}(G = g')}{p_{g'}} E(Y_d \mid G = g) [D - E(D \mid G = g')] - \frac{\mathbb{1}(G = g)}{p_g} \xi_{dgg'}. \end{aligned}$$

In Appendix B, we derive the EIFs for the general case with survey weights for both unconditional and conditional decompositions.

It follows that the one-step estimators are

$$\begin{aligned} \Psi_{dg}(Y, \mathbf{X}) &:= \frac{1}{n} \sum \frac{\mathbb{1}(G = g)}{p_g} \left\{ \frac{\mathbb{1}(D = d)}{\pi(d, \mathbf{X}, g)} [Y - \mu(d, \mathbf{X}, g)] + \mu(d, \mathbf{X}, g) \right\} \\ \Psi_{dgg'}(Y, \mathbf{X}) &:= \frac{1}{n} \sum \frac{\mathbb{1}(G = g)}{p_g} \left\{ \frac{\mathbb{1}(D = d)}{\pi(d, \mathbf{X}, g)} [Y - \mu(d, \mathbf{X}, g)] + \mu(d, \mathbf{X}, g) \right\} E(D \mid G = g'). \end{aligned}$$

There are two nuisance functions,  $\pi(d, \mathbf{X}, g)$  and  $\mu(d, \mathbf{X}, g)$ . In order to relax the conditions required of nuisance function models (Kennedy, 2022; Chernozhukov et al., 2018a), we use cross-

fitting to estimate all nuisance functions. In particular, we first divide the sample randomly into two subsamples. Then we fit the nuisance functions using each subsample. Finally, we evaluate the fitted nuisance functions using data not in the subsample used to fit the functions. In practice, to improve the finite-sample performance of the estimator, we may stabilize the weights by dividing  $\frac{\mathbb{1}(D=d)}{\pi(d, \mathbf{X}, g)}$  by its sample average.

To study the asymptotic behavior of the one-step estimators for the unconditional decomposition, we invoke three additional assumptions,<sup>9</sup> which are consistent with assumptions required for the double ML estimator of ATE (Kennedy, 2022; Chernozhukov et al., 2018a). We let  $\|\cdot\|$  denote the  $L_2$ -norm.

Assumption 4a (Boundedness). With probability 1,  $\hat{\pi}(d, \mathbf{X}, g) \geq \eta$ ,  $\pi(d, \mathbf{X}, g) \geq \eta$ , and  $|Y - \hat{\mu}(d, \mathbf{X}, g)| \leq \zeta$ , for some  $\eta > 0$  and some  $\zeta < \infty$ ,  $\forall d, g$ .

Assumption 5a (Consistency).  $\|\hat{\mu}(d, \mathbf{X}, g) - \mu(d, \mathbf{X}, g)\| = o_p(1)$  and  $\|\hat{\pi}(d, \mathbf{X}, g) - \pi(d, \mathbf{X}, g)\| = o_p(1)$ ,  $\forall d, g$ .

Assumption 6a (Convergence rate).  $\|\hat{\pi}(d, \mathbf{X}, g) - \pi(d, \mathbf{X}, g)\| \|\hat{\mu}(d, \mathbf{X}, g) - \mu(d, \mathbf{X}, g)\| = o_p(n^{-1/2})$ ,  $\forall d, g$ .

**Theorem 2** (Asymptotic distribution). *Under assumptions 1, 2, 3, 4a, 5a, and 6a, the one step estimators for components of the unconditional decomposition are asymptotically normal and semiparametrically efficient, i.e.,  $\sqrt{n}(\Psi_{dg} - \xi_{dg}) \xrightarrow{d} \mathcal{N}(0, \sigma_{dg}^2)$ , and  $\sqrt{n}(\Psi_{dgg'} - \xi_{dgg'}) \xrightarrow{d} \mathcal{N}(0, \sigma_{dgg'}^2)$ , where  $\sigma_{dg}^2 = E[\phi_{dg}(Y, \mathbf{X})^2]$  and  $\sigma_{dgg'}^2 = E[\phi_{dgg'}(Y, \mathbf{X})^2]$  are the respective semiparametric efficiency bounds.*

Assumptions 4a through 6a are mild in the sense that they are often satisfied by popular ML models. The desirable asymptotic properties also carry over from  $\Psi_{dg}$  and  $\Psi_{dgg'}$  to the final estimators of the decomposition components. We consistently estimate  $\sigma_{dg}^2$  and  $\sigma_{dgg'}^2$  using the averages of squared estimated EIFs. The asymptotic distributions can then be used to construct hypothesis testing and confidence intervals. For example, a Wald-style  $1 - \alpha$  confidence interval for the prevalence component is

$$\Psi_{1ba} - \Psi_{1bb} - \Psi_{0ba} + \Psi_{0bb} \pm z_{1-\alpha/2} \cdot \frac{1}{n} \left[ \sum \left( \hat{\phi}_{1ba} - \hat{\phi}_{1bb} - \hat{\phi}_{0ba} + \hat{\phi}_{0bb} \right)^2 \right]^{1/2}.$$

Furthermore, our estimators also have a double robustness property.

**Theorem 3** (Double robustness). *Either consistent estimation of  $\mu(d, \mathbf{X}, g)$  or  $\pi(d, \mathbf{X}, g)$  for*

---

<sup>9</sup>Throughout this paper, the consistent estimation of  $p_g$  and  $E(D | G = g), \forall g$ , are left implicit.

all  $d$  and  $g$  is sufficient for the consistency of  $\Psi_{dg}$  and  $\Psi_{dgg'}$ .

Hence, our unconditional decomposition estimator only requires the same conditions as the classic augmented inverse probability weighting estimator for ATE (Robins et al., 1994; Hirano et al., 2003) to be consistent.

### 3.2 The conditional decomposition

Relative to the unconditional case, we only need to additionally consider one generic function:

$$\xi_{dgg'g''} := E[E(Y_d | \mathbf{Q}, g) E(D | \mathbf{Q}, g') | G = g''],$$

where  $(d, g, g', g'')$  is any combination of treatment status and group memberships out of 8 possible combinations. The relationship between components of the conditional decomposition and the generic functions are as follows:

$$\text{Baseline} = \xi_{0a} - \xi_{0b}$$

$$\text{Conditional Prevalence} = \xi_{1bab} - \xi_{0bab} - \xi_{1bbb} + \xi_{0bbb}$$

$$\text{Conditional Effect} = \xi_{1aaa} - \xi_{0aaa} - \xi_{1baa} + \xi_{0baa}$$

$$\mathbf{Q} \text{ Distribution} = \xi_{1baa} - \xi_{0baa} - \xi_{1bab} + \xi_{0bab}$$

$$\text{Conditional Selection} = \text{Total} - \text{Baseline}$$

$$- \text{Conditional Prevalence} - \text{Conditional Effect} - \mathbf{Q} \text{ Distribution}.$$

Hence, in this subsection, we focus on the inference of  $\xi_{dgg'g''}$ . The EIFs, one-step estimators, and their asymptotic distributions for components of the conditional decomposition will then straightforwardly follow. Moreover, we thereby also provide nonparametric inference for the reduction in disparity component in the CRED of Jackson (2021), which can be represented as  $\xi_{0b} + \xi_{1bab} - \xi_{0bab} - E_b(Y)$ . Under assumptions 1, 2, 3, we identify  $\xi_{dgg'g''}$  as

$$E\{E[E(Y | d, \mathbf{X}, g) | \mathbf{Q}, g] E(D | \mathbf{Q}, g') | G = g''\}.$$

**Theorem 4** (EIF). *Under assumptions 1, 2, 3, the EIF of  $\xi_{dgg'g''}$  is*

$$\phi_{dgg'g''}(Y, \mathbf{X}, \mathbf{Q})$$

$$\begin{aligned}
&= \frac{\mathbb{1}(G = g'')}{p_{g''}} [\mathbb{E}(Y_d \mid \mathbf{Q}, g) \mathbb{E}(D \mid \mathbf{Q}, g') - \xi_{dgg'g''}] \\
&\quad + \frac{\mathbb{1}(G = g)p_{g''}(\mathbf{Q})}{p_g(\mathbf{Q})p_{g''}} \left\{ \frac{\mathbb{1}(D = d)}{\pi(d, \mathbf{X}, g)} [Y - \mu(d, \mathbf{X}, g)] + \mu(d, \mathbf{X}, g) - \mathbb{E}(Y_d \mid \mathbf{Q}, g) \right\} \mathbb{E}(D \mid \mathbf{Q}, g') \\
&\quad + \frac{\mathbb{1}(G = g')p_{g''}(\mathbf{Q})}{p_{g'}(\mathbf{Q})p_{g''}} [D - \mathbb{E}(D \mid \mathbf{Q}, g')] \mathbb{E}(Y_d \mid \mathbf{Q}, g).
\end{aligned}$$

We again construct the one-step estimators based on the EIF.

$$\begin{aligned}
&\Psi_{dgg'g''}(Y, \mathbf{X}, \mathbf{Q}) \\
&= \frac{\mathbb{1}(G = g'')}{p_{g''}} \mathbb{E}(Y_d \mid \mathbf{Q}, g) \mathbb{E}(D \mid \mathbf{Q}, g') + \frac{\mathbb{1}(G = g')p_{g''}(\mathbf{Q})}{p_{g'}(\mathbf{Q})p_{g''}} [D - \mathbb{E}(D \mid \mathbf{Q}, g')] \mathbb{E}(Y_d \mid \mathbf{Q}, g) \\
&\quad + \frac{\mathbb{1}(G = g)p_{g''}(\mathbf{Q})}{p_g(\mathbf{Q})p_{g''}} \left\{ \frac{\mathbb{1}(D = d)}{\pi(d, \mathbf{X}, g)} [Y - \mu(d, \mathbf{X}, g)] + \mu(d, \mathbf{X}, g) - \mathbb{E}(Y_d \mid \mathbf{Q}, g) \right\} \mathbb{E}(D \mid \mathbf{Q}, g').
\end{aligned}$$

For this estimator, there are five nuisance functions:  $p_g(\mathbf{Q})$ ,  $\pi(d, \mathbf{X}, g)$ ,  $\mu(d, \mathbf{X}, g)$ ,  $\mathbb{E}(D \mid \mathbf{Q}, g)$ , and  $\mathbb{E}(Y_d \mid \mathbf{Q}, g)$ . We again cross-fit these nuisance functions using nonparametric methods. This estimator can be weight-stabilized by dividing  $\frac{\mathbb{1}(D=d)}{\pi(d, \mathbf{X}, g)}$ ,  $\frac{\mathbb{1}(G=g')p_{g''}(\mathbf{Q})}{p_{g'}(\mathbf{Q})p_{g''}}$ , and  $\frac{\mathbb{1}(G=g)p_{g''}(\mathbf{Q})}{p_g(\mathbf{Q})p_{g''}}$  by their respective sample averages. The estimation of  $\mathbb{E}(Y_d \mid \mathbf{Q}, g) = \mathbb{E}[\mu(d, \mathbf{X}, g) \mid \mathbf{Q}, g]$  requires a more involved procedure. We adopt a pseudo-outcome approach (e.g., [van der Laan, 2006](#); [Semenova and Chernozhukov, 2021](#)), where the pseudo outcome for each  $d$  is defined as the decentered EIF for  $\mathbb{E}(Y_d)$ , i.e.,

$$\delta_d(Y, \mathbf{X}, G) := \frac{\mathbb{1}(D = d)}{\pi(d, \mathbf{X}, G)} [Y - \mu(d, \mathbf{X}, G)] + \mu(d, \mathbf{X}, G).$$

We first estimate  $\pi(d, \mathbf{X}, G)$  and  $\mu(d, \mathbf{X}, G)$  within two random subsamples of the data *without* cross-fitting and obtain the estimated pseudo outcome  $\hat{\delta}_d(Y, \mathbf{X}, G)$ . Then, using  $\hat{\delta}_d(Y, \mathbf{X}, G)$  as a proxy of  $Y_d$ , we estimate  $\mathbb{E}[Y_d \mid \mathbf{Q}, g]$  *with* cross-fitting, i.e., we fit  $\mathbb{E}[\hat{\delta}_d(Y, \mathbf{X}, G) \mid \mathbf{Q}, g]$  in one subsample and plug in values of  $\mathbf{Q}$  in the other subsample. Using this procedure, we ensure that the fitting of  $\mathbb{E}(Y_d \mid \mathbf{Q}, g)$ , which relies on estimating the pseudo outcome, is done separately in each subsample.

The asymptotic theory of the one-step estimator for  $\xi_{dgg'g''}$  is less general than the theories for  $\xi_{dg}$  and  $\xi_{dgg'}$ , in that some conditions differ by the specific configuration of  $g, g'$  and  $g''$ . Note that  $g = g''$  for the conditional prevalence component, and  $g' = g''$  for the conditional effect component. For  $\sqrt{n}$ -consistency, asymptotic normality, and efficiency, we need three assumptions.

Assumption 4b (Boundedness). With probability 1,  $\hat{\pi}(d, \mathbf{X}, g) \geq \eta$ ,  $\pi(d, \mathbf{X}, g) \geq \eta$ ,  $\hat{p}_g(\mathbf{Q}) \geq \eta$ ,  $p_g(\mathbf{Q}) \geq \eta$ ,  $|Y - \hat{\mu}(d, \mathbf{X}, g)| \leq \zeta$ ,  $|Y - \mu(d, \mathbf{X}, g)| \leq \zeta$ ,  $|\mu(d, \mathbf{X}, g)| \leq \zeta$ ,  $|\mathbb{E}(Y_d | \mathbf{Q}, g)| \leq \zeta$ , and  $\left| \frac{\hat{p}_g(\mathbf{Q})}{\hat{p}_{g'}(\mathbf{Q})} \right| \leq \zeta$ , for some  $\eta > 0$  and  $\zeta < \infty$ ,  $\forall d, g, g'$ .

Assumption 5b (Consistency).  $\|\hat{\mu}(d, \mathbf{X}, g) - \mu(d, \mathbf{X}, g)\| = o_p(1)$ ,  $\|\hat{\pi}(d, \mathbf{X}, g) - \pi(d, \mathbf{X}, g)\| = o_p(1)$ ,  $\left\| \hat{\mathbb{E}}(Y_d | \mathbf{Q}, g) - \mathbb{E}(Y_d | \mathbf{Q}, g) \right\| = o_p(1)$ ,  $\left\| \hat{\mathbb{E}}(D | \mathbf{Q}, g) - \mathbb{E}(D | \mathbf{Q}, g) \right\| = o_p(1)$ , and  $\|\hat{p}_g(\mathbf{Q}) - p_g(\mathbf{Q})\| = o_p(1)$ ,  $\forall d, g$ .

Assumption 6b (Convergence rate).  $\|\hat{\pi}(d, \mathbf{X}, g) - \pi(d, \mathbf{X}, g)\| \|\hat{\mu}(d, \mathbf{X}, g) - \mu(d, \mathbf{X}, g)\| = o_p(n^{-1/2})$ ,  $\left\| \mathbb{1}(G = g) \mathbb{E}(Y_d | \mathbf{Q}, g) - \mathbb{1}(G = g') \frac{\hat{p}_g(\mathbf{Q})}{\hat{p}_{g'}(\mathbf{Q})} \hat{\mathbb{E}}(Y_d | \mathbf{Q}, g) \right\| \left\| \hat{\mathbb{E}}(D | \mathbf{Q}, g') - \mathbb{E}(D | \mathbf{Q}, g') \right\| = o_p(n^{-1/2})$ , and  $\left\| \mathbb{1}(G = g') \mathbb{E}(D | \mathbf{Q}, g') - \mathbb{1}(G = g) \frac{\hat{p}_{g'}(\mathbf{Q})}{\hat{p}_g(\mathbf{Q})} \hat{\mathbb{E}}(D | \mathbf{Q}, g') \right\| \left\| \hat{\mathbb{E}}(Y_d | \mathbf{Q}, g) - \mathbb{E}(Y_d | \mathbf{Q}, g) \right\| = o_p(n^{-1/2})$ ,  $\forall d, g, g'$ .

**Theorem 5** (Asymptotic distribution). *Under assumptions 1, 2, 3, 4b, 5b, 6b, the one step estimator for components of the conditional decomposition is asymptotically normal and semiparametrically efficient, i.e.,  $\sqrt{n}(\Psi_{dgg'g''} - \xi_{dgg'g''}) \rightarrow \mathcal{N}(0, \sigma_{dgg'g''}^2)$ , where  $\sigma_{dgg'g''}^2 = \mathbb{E}[\phi_{dgg'g''}(Y, \mathbf{X}, \mathbf{Q})^2]$  is the semiparametric efficiency bound.*

Some further insights can be gained about the second and third convergence rate assumptions. For a specific  $\xi_{dgg'g''}$ , first, when  $g = g' = g''$ ,  $\left\| \hat{\mathbb{E}}(Y_d | \mathbf{Q}, g) - \mathbb{E}(Y_d | \mathbf{Q}, g) \right\| \left\| \hat{\mathbb{E}}(D | \mathbf{Q}, g) - \mathbb{E}(D | \mathbf{Q}, g) \right\| = o_p(n^{-1/2})$  is sufficient. Hence, we obtain a form of rate double robustness with respect to  $\mathbb{E}(Y_d | \mathbf{Q}, g)$  and  $\mathbb{E}(D | \mathbf{Q}, g)$ . Second, when  $g = g'' \neq g'$ , the following set of conditions is sufficient: for a constant  $\zeta > 0$ ,  $\frac{\hat{p}_g(\mathbf{Q})}{\hat{p}_{g'}(\mathbf{Q})} \leq \zeta$  with probability 1,  $\left\| \hat{\mathbb{E}}(D | \mathbf{Q}, g') - \mathbb{E}(D | \mathbf{Q}, g') \right\| = o_p(n^{-1/2})$ , and  $\left\| \hat{\mathbb{E}}(Y_d | \mathbf{Q}, g) - \mathbb{E}(Y_d | \mathbf{Q}, g) \right\| = o_p(1)$ . Third, when  $g' = g'' \neq g$ , a sufficient set is: for a constant  $\zeta > 0$ ,  $\frac{\hat{p}_{g'}(\mathbf{Q})}{\hat{p}_g(\mathbf{Q})} \leq \zeta$  with probability 1,  $\left\| \hat{\mathbb{E}}(D | \mathbf{Q}, g') - \mathbb{E}(D | \mathbf{Q}, g') \right\| = o_p(1)$ , and  $\left\| \hat{\mathbb{E}}(Y_d | \mathbf{Q}, g) - \mathbb{E}(Y_d | \mathbf{Q}, g) \right\| = o_p(n^{-1/2})$ . Therefore, the assumption is weaker for the inference of the conditional prevalence component, which is also the case for Theorem 6 below.

**Theorem 6** (Double robustness). *We assume consistent estimation of  $\mathbb{E}(D | \mathbf{Q}, g), \forall g$ . For a  $\xi_{dgg'g''}$ , when  $g = g''$ , consistent estimation of either  $\mu(d, \mathbf{X}, g)$  or  $\pi(d, \mathbf{X}, g)$  is sufficient for the consistency of  $\Psi_{dgg'g''}$ ; when  $g \neq g''$ , if the estimation of  $\mathbb{E}(Y_d | \mathbf{Q}, g)$  is doubly robust with respect to  $\mu(d, \mathbf{X}, g)$  and  $\pi(d, \mathbf{X}, g)$ , then consistent estimation of either  $\mu(d, \mathbf{X}, g)$  or  $\pi(d, \mathbf{X}, g)$  is still sufficient for the consistency of  $\Psi_{dgg'g''}$ .*

Note that consistent estimation of  $\hat{p}_g(\mathbf{Q})$  is not required for the consistency of  $\Psi_{dgg'g''}$ . The double robustness condition for  $\mathbb{E}(Y_d | \mathbf{Q}, g)$  in the case of  $g \neq g''$  motivates our pseudo-outcome estimator of  $\mathbb{E}(Y_d | \mathbf{Q}, g)$ , which is built on both  $\mu(d, \mathbf{X}, g)$  and  $\pi(d, \mathbf{X}, g)$ .

Both the unconditional and conditional decompositions are implemented in the R ([R Core Team, 2023](#)) package `cdgd` ([Yu, 2023](#)), available from CRAN. The package offers various options for nuisance function estimation, including various ML methods, a parametric option with linear and logistic regressions, and an option that allows users to manually plug in nuisance function estimates.

## 4 Application

### 4.1 Overview

We apply our causal decomposition to investigate the role of college graduation in intergenerational income mobility in the United States. In this empirical setting, groups are defined by parental income, the outcome is adult income, and the treatment is college graduation. Previous research has touched upon the components of our unconditional decomposition to various degrees, but this is the first study to comprehensively evaluate the three distinct ways in which college graduation might contribute to the intergenerational transmission of income advantages.

The baseline component represents the part of the income attainment disparity that is unaccounted for by college graduation. Apart from college completion, parental income is associated with a variety of pre-college characteristics that have bearings on income attainment, such as academic achievement, cognitive skills, and noncognitive traits in adolescence ([Reardon, 2011](#); [Heckman et al., 2006](#); [Farkas, 2003](#)). Moreover, in adulthood, people from more privileged backgrounds may continue to benefit from various resources passed on within their families regardless of formal educational attainment. Any intervention imposed on the post-secondary credentialing process cannot eliminate channels of income persistence that are independent of college graduation.

In a way conceptually consistent with the prevalence component in our decomposition, social scientists have construed education as a mediator of the intergenerational reproduction of SES inequalities ([Blau and Duncan, 1967](#), chapter 4 & 5; [Featherman and Hauser, 1978](#), p.255-9; [Ishida et al., 1995](#); [Breen, 2010](#)). In particular, researchers have documented a large and increasing income gap in college graduation ([Ziol-Guest and Lee, 2016](#); [Bailey and Dynarski, 2011](#)). Moreover, using a simulation based on descriptive models, [Bloome et al. \(2018\)](#) concluded that rising educational inequality strengthened intergenerational income persistence

over time.

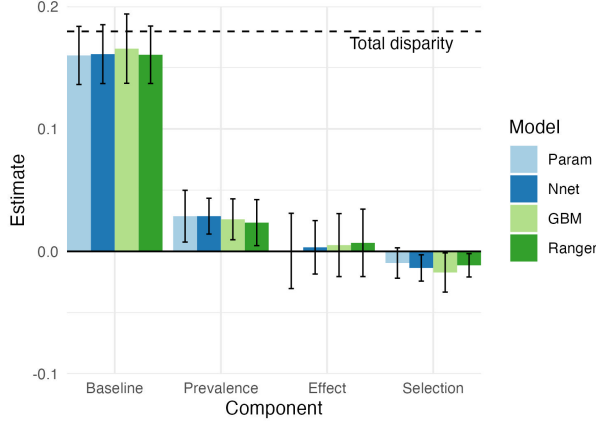
There is also an active literature on the heterogeneity of the college effect on SES attainment by SES origin, which underpins our effect component. Brand et al. (2021) and Cheng et al. (2021) found that the college premiums on wage and income outcomes are larger for people who are economically disadvantaged in adolescence. However, Zhou (2019) and Fiel (2020) found only insignificant heterogeneity in the effect of college completion on income by parental income. However, these works did not evaluate the extent to which income disparities can be attributed to groupwise differential effects of college.

Finally, the literature on college effects on later SES status has yielded mixed findings regarding whether selection on the treatment effect is positive or negative overall in the US. Brand and Xie (2010) and Brand et al. (2021) concluded that there is negative selection into college. However, Breen et al. (2015) warned that Brand and Xie’s (2010) result may be biased from unobserved confounding. The instrumental variable analysis of Heckman et al. (2018), on the other hand, reveals positive selection. This line of inquiry does not estimate group-specific selection patterns, thereby missing the linkage between group-differentiated selection and outcome disparities. In Appendix E, we clarify and synthesize the differences and connections between various concepts of “selection into college” in the social science literature.

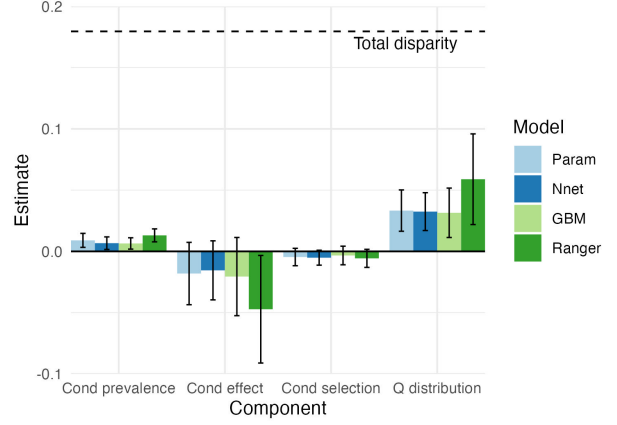
## 4.2 Data, variables and estimation details

We analyze the National Longitudinal Survey of Youth 1979 (NLSY79), which is a nationally representative dataset for a cohort of young men and women who were born between 1957 and 1964 and were living in the United States at the baseline survey in 1979. We limit the analysis to respondents who were between 14 to 17 years old in 1979 to ensure that income origin is measured before the respondents had a chance to graduate from college.

We contrast income-origin groups defined as the top 60% and bottom 40% of family incomes averaged over the first three waves of the survey (1979, 1980, and 1981), divided by the square root of the family size. The respondents were 14 to 20 years old during that time. The treatment variable is a binary indicator of whether the respondent graduated from college by the age of 31. We measure the outcome by the percentile rank of adult income, averaged over five survey waves between age 35 and 44, again divided by the square root of family size. For the conditional decomposition, we define  $Q$  to be the Armed Forces Qualification Test (AFQT) score measured in 1980, which reflects academic achievement in high school.



(c) Unconditional decomposition



(d) Conditional decomposition

Figure 2. Decomposition estimates. Cond=conditional, Para=parametric models, Nnet=neural networks, GBM=generalized boosting machine, Ranger=random forests. Error bars indicate 95% confidence intervals. Baseline component is omitted in panel (b) as the baseline component estimates for the conditional decomposition are, by definition, the same as those for the unconditional decomposition.

Covariates used for causal identification include gender, race, parental income percentile, parental education, parental presence, the number of siblings, urban residence, educational expectation, friends' educational expectation, AFQT score, age at baseline survey, the Rotter locus of control scale score, the Rosenberg esteem item response score, language spoken at home, urban/rural residence, presence of mother, school satisfaction, region of residence, and mother's working status.

We fit the nuisance functions using three alternative ML methods, gradient boosting machine (GBM), neural networks and random forests. For comparison, we also fit the nuisance functions using parametric regressions. The parametric models are set up as follows. For models predicting  $Y$ , we use linear regressions with all two-way interactions between the treatment and baseline variables including the group indicator, along with their main effects. For propensity score models and models predicting  $G$ , the covariates are entered linearly. For other models, where only  $G$  and  $Q$  are predictors, we include all main effects and two-way interactions between  $G$  and  $Q$ .<sup>10</sup>



### 4.3 Results

We illustrate estimates of the decomposition components in Figure 1, while the numeric estimates are presented in Appendix Table A2 and A3. To aid the interpretation of the decomposition components, in Appendix Table A1, we additionally report estimates of group-specific means of baseline potential outcome, treatment proportions, ATEs, and covariances between treatment and treatment effects, as well as group differences in these quantities.

To begin with, according to the estimated total disparity, people from lower income origin on average achieve 18 percentiles lower incomes in their 30s and 40s. This confirms the well-established pattern of intergenerational income persistence. Our decompositions reveal aspects of the formation of this persistence. Across models for the nuisance functions, the baseline component constitutes around 90% of the total disparity. This suggests that income advantages are transmitted between generations via many channels other than inequalities related to a college degree. The baseline component hence demarcates an upper limit for interventions aiming at post-secondary education.

Nevertheless, differential prevalence of college degrees across income origin groups plays a substantial and statistically significant role in explaining the income disparity in adulthood. In the unconditional decomposition, the prevalence component causally accounts for about 15% of the total disparity in adult income. Therefore, if an intervention could equalize college graduation rates across income origins, the outcome disparity would be reduced by 15%. What underlies the prevalence component is the vast inequality in college graduation by parental income: 29% in the higher-income group obtained a college degree by age 31 while only 9% in the lower-income group did so (see Appendix Table A1). The conditional prevalence component in the conditional decomposition remains statistically significant in all models, although it has much smaller magnitudes and only accounts for 3% to 7% of the total disparity. Therefore, equalizing chances of graduating college within levels of prior achievement would still make a difference, even small.

The contribution of the effect components in both unconditional and conditional decompositions is small in magnitude and generally not statistically significant. This is due to minimal between-group effect heterogeneity, either unconditionally or conditionally. In Appendix A1, we show that across models, the group-specific ATEs range from 12 to 15 percentiles and are all statistically significant. However, the groupwise difference is always insignificant. This is

---

<sup>10</sup>Code for this analysis is available at [https://github.com/ang-yu/causal\\_decomposition\\_case\\_study](https://github.com/ang-yu/causal_decomposition_case_study).

consistent with findings of [Zhou \(2019\)](#) and [Fiel \(2020\)](#), which are also based on NLSY79.

All three ML models are in agreement that the selection component in the unconditional decomposition is significantly negative, offsetting 48% to 65% of the outcome disparity caused by differential prevalence. Had the selection component been zero, the total disparity would be around 6% higher. In fact, selection into college is positive in the disadvantaged group and negative in the advantaged group, although not significantly. This echoes the view that obtaining a college degree is more of a rational decision based on expected returns and less of a cultural norm among the disadvantaged than among the advantaged ([Mare, 1980](#); [Hout, 2012](#)). On the other hand, the conditional selection component does not notably contribute to the disparity in income attainment.

Descriptively, the  $Q$  distribution component also plays an important role in the generation of the outcome disparity, because of a large income gap in academic achievement in high school. Finally, we also estimate the reduction in disparity components in the unconditional and conditional random equalization decompositions ([Jackson and VanderWeele, 2018](#); [Jackson, 2021](#)) and present the results in Appendix Table A2 and A3. In the unconditional case, relative to our prevalence component, the random equalization decomposition underestimates the extent to which differential college graduation rates contribute to the outcome disparity. This is because selection into college is positive in the lower-income group. By randomizing college degrees within the lower-income group, the intervention underlying the random equalization decomposition would be less powerful in reducing the outcome disparity than the pure equalization intervention corresponding to our prevalence component.

## 5 Discussion

We provide a new decomposition approach for quantifying the ways in which a treatment variable affects outcome disparities between ascriptive groups. Our decompositions improve on previous approaches of explaining group disparities and reveal novel policy levers for ameliorating disparities. Our decompositions are expressed as model-free estimands. We introduced nonparametric estimators that are efficient, asymptotically normal, and doubly robust. We applied our decompositions to study the role of college education in intergenerational income persistence. We highlight that differential prevalence of the college degree significantly contributes to the income disparity, even conditional on academic achievement in high school.

Conversely, differential selection into college graduation plays a significantly compensatory role and reduces the income disparity.

Our approach can be extended in multiple directions. First, it is possible to develop analogous decompositions for non-binary treatments. For multi-valued categorical treatments, the basic form of equation (1) can be retained, with prevalence, effect, and selection components replaced with sums of category-specific contrasts. For continuous treatments, the randomization and equalization interventions, as well as their conditional versions, remain well-defined and may serve as bases of decompositions. Second, our framework can be extended to accommodate multiple temporally-ordered treatments. In the case of two treatments,  $D_1$  and  $D_2$ , such extension can be based on the following decomposition of the outcome:

$$Y = Y(0, 0) + D_1[Y(1, 0) - Y(0, 0)] + D_2[Y(0, 1) - Y(0, 0)] + D_1 D_2[Y(1, 1) - Y(1, 0) - Y(0, 1) + Y(0, 0)],$$

where  $Y(d_1, d_2)$  denotes the potential outcome of  $Y$  under the assignment of  $D_1 = d_1$  and  $D_2 = d_2$ . Third, as an alternative to the double ML-style estimators in this paper, targeted learning (Van der Laan and Rose, 2011) could also be employed for estimation, which may have better finite-sample performance when the outcome is bounded.

## Appendices

### Appendix A: Proofs for Section 2

#### A.1. Equation (2)

Note that  $R(D \mid G = g')$  denotes a randomly drawn value of treatment  $D$  from group  $g'$ .

$$\begin{aligned} & E_g(Y_{R(D \mid G = g')}) \\ &= E_g(Y_1 \mid R(D \mid G = g') = 1) \Pr_g(R(D \mid G = g') = 1) + E_g(Y_0 \mid R(D \mid G = g') = 0) \Pr_g(R(D \mid G = g') = 0) \\ &= E_g(Y_1) E_{g'}(D) + E_g(Y_0)(1 - E_{g'}(D)) \\ &= E_g(Y_0) + E_{g'}(D) E_g(\tau). \end{aligned}$$

## A.2. Equivalence results in subsection 2.3

For the CDE,

$$\begin{aligned}
& E(Y_{G=a,D=0}) - E(Y_{G=b,D=0}) \\
&= E(Y_{G=a,D=0} \mid G = a) - E(Y_{G=b,D=0} \mid G = b) \\
&= E_a(Y_{D=0}) - E_b(Y_{D=0}),
\end{aligned}$$

where the first equality is by the unconditional ignorability of  $G$ , and the second equality is by consistency.

For the RIA of PIE,

$$\begin{aligned}
& E(Y_{G=b,R(D_{G=a})}) - E(Y_{G=b,R(D_{G=b})}) \\
&= E(Y_{G=b,D=1} \mid R(D_{G=a}) = 1) \Pr(R(D_{G=a}) = 1) + E(Y_{G=b,D=0} \mid R(D_{G=a}) = 0) \Pr(R(D_{G=a}) = 0) \\
&\quad - E(Y_{G=b,D=1} \mid R(D_{G=b}) = 1) \Pr(R(D_{G=b}) = 1) - E(Y_{G=b,D=0} \mid R(D_{G=b}) = 0) \Pr(R(D_{G=b}) = 0) \\
&= E(Y_{G=b,D=1}) \Pr(D_{G=a} = 1) + E(Y_{G=b,D=0}) \Pr(D_{G=a} = 0) \\
&\quad - E(Y_{G=b,D=1}) \Pr(D_{G=b} = 1) - E(Y_{G=b,D=0}) \Pr(D_{G=b} = 0) \\
&= E(Y_{G=b,D=1} \mid G = b) \Pr(D_{G=a} = 1 \mid G = a) + E(Y_{G=b,D=0} \mid G = b) \Pr(D_{G=a} = 0 \mid G = a) \\
&\quad - E(Y_{G=b,D=1} \mid G = b) \Pr(D_{G=b} = 1 \mid G = b) - E(Y_{G=b,D=0} \mid G = b) \Pr(D_{G=b} = 0 \mid G = b) \\
&= E(Y_{D=1} \mid G = b) \Pr(D = 1 \mid G = a) + E(Y_{D=0} \mid G = b) \Pr(D = 0 \mid G = a) \\
&\quad - E(Y_{D=1} \mid G = b) \Pr(D = 1 \mid G = b) - E(Y_{D=0} \mid G = b) \Pr(D = 0 \mid G = b) \\
&= E_b(\tau)[E_a(D) - E_b(D)],
\end{aligned}$$

where the third equality holds by the unconditional ignorability of  $G$ , and the fourth holds by the consistency assumptions.

For the sum of the RIA of RIE and the RIA of the MIE,

$$\begin{aligned}
& E(Y_{G=a,D=1} - Y_{G=a,D=0} - Y_{G=b,D=1} + Y_{G=b,D=0}) E(D_{G=a}) \\
&= [E(Y_{G=a,D=1} \mid G = a) - E(Y_{G=a,D=0} \mid G = a) - E(Y_{G=b,D=1} \mid G = b) + E(Y_{G=b,D=0} \mid G = b)] E(D_{G=a} \mid G = a) \\
&= [E(Y_{D=1} \mid G = a) - E(Y_{D=0} \mid G = a) - E(Y_{D=1} \mid G = b) + E(Y_{D=0} \mid G = b)] E(D \mid G = a) \\
&= [E_a(\tau) - E_b(\tau)] E_a(D),
\end{aligned}$$

where the first equality holds by unconditional ignorability of  $G$  and the second by consistency.

### A.3. Equation (3)

$$\begin{aligned}
E_a(Y) - E_b(Y) &= E_a(Y_0) + E_b(Y_0) \\
&\quad + E_a(D\tau) - E_b(D\tau) - E_a[E_a(D \mid \mathbf{Q}) E_a(\tau \mid \mathbf{Q})] + E_b[E_b(D \mid \mathbf{Q}) E_b(\tau \mid \mathbf{Q})] \\
&\quad + E_a[E_a(D \mid \mathbf{Q}) E_a(\tau \mid \mathbf{Q})] - E_b[E_b(D \mid \mathbf{Q}) E_b(\tau \mid \mathbf{Q})] \\
&= E_a(Y_0) + E_b(Y_0) + E_a[\text{Cov}_a(D, \tau \mid \mathbf{Q})] - E_b[\text{Cov}_b(D, \tau \mid \mathbf{Q})] \\
&\quad + \int E_a(D \mid \mathbf{Q} = \mathbf{q}) E_a(\tau \mid \mathbf{Q} = \mathbf{q}) f_a(\mathbf{q}) d\mathbf{q} - \int E_b(D \mid \mathbf{Q} = \mathbf{q}) E_b(\tau \mid \mathbf{Q} = \mathbf{q}) f_b(\mathbf{q}) d\mathbf{q} \\
&= E_a(Y_0) - E_b(Y_0) \\
&\quad + \int [E_a(D \mid \mathbf{Q} = \mathbf{q}) - E_b(D \mid \mathbf{Q} = \mathbf{q})] E_b(\tau \mid \mathbf{Q} = \mathbf{q}) f_b(\mathbf{q}) d\mathbf{q} \\
&\quad + \int E_a(D \mid \mathbf{Q} = \mathbf{q}) E_b(\tau \mid \mathbf{Q} = \mathbf{q}) [f_a(\mathbf{q}) - f_b(\mathbf{q})] d\mathbf{q} \\
&\quad + \int [E_a(\tau \mid \mathbf{Q} = \mathbf{q}) - E_b(\tau \mid \mathbf{Q} = \mathbf{q})] E_a(D \mid \mathbf{Q} = \mathbf{q}) f_a(\mathbf{q}) d\mathbf{q} \\
&\quad + E_a[\text{Cov}_a(D, \tau \mid \mathbf{Q})] - E_b[\text{Cov}_b(D, \tau \mid \mathbf{Q})].
\end{aligned}$$

Note that the last equality uses the assumption that  $\text{supp}_a(\mathbf{Q}) = \text{supp}_b(\mathbf{Q})$ .

### A.4. Equation (4)

$$\begin{aligned}
&E_g(Y_{R(D|G=g', \mathbf{Q})}) \\
&= \int E_g(Y_{R(D|G=g', \mathbf{Q})} \mid \mathbf{Q} = \mathbf{q}) f_g(\mathbf{q}) d\mathbf{q} \\
&= \int E_g(Y_1 \mid \mathbf{Q} = \mathbf{q}, R(D \mid G = g', \mathbf{Q} = \mathbf{q}) = 1) \Pr_g(R(D \mid G = g', \mathbf{Q} = \mathbf{q}) = 1 \mid \mathbf{Q} = \mathbf{q}) f_g(\mathbf{q}) d\mathbf{q} \\
&\quad + \int E_g(Y_0 \mid \mathbf{Q} = \mathbf{q}, R(D \mid G = g', \mathbf{Q} = \mathbf{q}) = 0) \Pr_g(R(D \mid G = g', \mathbf{Q} = \mathbf{q}) = 0 \mid \mathbf{Q} = \mathbf{q}) f_g(\mathbf{q}) d\mathbf{q} \\
&= \int [E_g(Y_1 \mid \mathbf{Q} = \mathbf{q}) E_{g'}(D \mid \mathbf{Q} = \mathbf{q}) + E_g(Y_0 \mid \mathbf{Q} = \mathbf{q}) (1 - E_{g'}(D \mid \mathbf{Q} = \mathbf{q}))] f_g(\mathbf{q}) d\mathbf{q} \\
&= E_g(Y_0) + \int E_g(\tau \mid \mathbf{Q} = \mathbf{q}) E_{g'}(D \mid \mathbf{Q} = \mathbf{q}) f_g(\mathbf{q}) d\mathbf{q}.
\end{aligned}$$

All expectations are taken over  $\mathbf{q} \in \text{supp}_g(\mathbf{Q})$ . For  $E_g(Y_{R(D|G=g', \mathbf{Q})})$  to be well-defined, we require  $\text{supp}_g(\mathbf{Q}) \subseteq \text{supp}_{g'}(\mathbf{Q})$ .

### A.5. Equation (5)

$$\begin{aligned}
& E_a(Y) - E_b(Y) - [E_a(Y_{R(D|G=a, \mathbf{Q})}) - E_b(Y_{R(D|G=b, \mathbf{Q})})] \\
&= E_a(Y_0) - E_b(Y_0) + E_a(D\tau) - E_b(D\tau) - [E_a(Y_{R(D|G=a, \mathbf{Q})}) - E_b(Y_{R(D|G=b, \mathbf{Q})})] \\
&= E_a[E_a(D\tau \mid \mathbf{Q} = \mathbf{q})] - E_a[E_a(\tau \mid \mathbf{Q} = \mathbf{q}) E_a(\tau \mid \mathbf{Q} = \mathbf{q})] \\
&\quad - \{E_b[E_b(D\tau \mid \mathbf{Q} = \mathbf{q})] - E_b[E_b(\tau \mid \mathbf{Q} = \mathbf{q}) E_b(\tau \mid \mathbf{Q} = \mathbf{q})]\} \\
&= E_a[\text{Cov}_a(D, \tau \mid \mathbf{Q} = \mathbf{q})] - E_b[\text{Cov}_b(D, \tau \mid \mathbf{Q} = \mathbf{q})].
\end{aligned}$$

Other results in Equation (5) follow directly from Equation (4).

### A.6. Equation (6)

$$\begin{aligned}
& E_b(Y_{R(D|G=a, \mathbf{Q})}) - E_b(Y) \\
&= E_b(Y_0) + \int E_b(\tau \mid \mathbf{Q} = \mathbf{q}) E_a(D \mid \mathbf{Q} = \mathbf{q}) f_b(\mathbf{q}) d\mathbf{q} - E_b(Y_0) - E_b(D\tau) \\
&= \int [E_a(D \mid \mathbf{Q} = \mathbf{q}) - E_b(D \mid \mathbf{Q} = \mathbf{q})] E_b(\tau \mid \mathbf{Q} = \mathbf{q}) f_b(\mathbf{q}) d\mathbf{q} \\
&\quad + E_b[E_b(D \mid \mathbf{Q}) E_b(\tau \mid \mathbf{Q})] - E_b[E_b(D\tau \mid \mathbf{Q})] \\
&= \int [E_a(D \mid \mathbf{Q} = \mathbf{q}) - E_b(D \mid \mathbf{Q} = \mathbf{q})] E_b(\tau \mid \mathbf{Q} = \mathbf{q}) f_b(\mathbf{q}) d\mathbf{q} - E_b[\text{Cov}_b(D, \tau \mid \mathbf{Q} = \mathbf{q})].
\end{aligned}$$

### A.7. Equation (7)

$$\begin{aligned}
& E_a(Y) - E_b(Y) - \\
& \{E_a[\text{Pr}(D = 0 \mid \mathbf{Q})Y_0 + \text{Pr}(D = 1 \mid \mathbf{Q})Y_1] - E_b[\text{Pr}(D = 0 \mid \mathbf{Q})Y_0 + \text{Pr}(D = 1 \mid \mathbf{Q})Y_1]\} \\
&= E_a(Y) - E_b(Y) - \{E_a(Y_0) + E_a[E(D \mid \mathbf{Q})\tau] - E_b(Y_0) - E_b[E(D \mid \mathbf{Q})\tau]\} \\
&= E_a(D\tau) - E_a[E(D \mid \mathbf{Q}) E_a(\tau \mid \mathbf{Q})] - \{E_b(D\tau) - E_b[E(D \mid \mathbf{Q}) E_b(\tau \mid \mathbf{Q})]\} \\
&= \int [E_a(D\tau \mid \mathbf{Q} = \mathbf{q}) - E(D \mid \mathbf{Q} = \mathbf{q}) E_a(\tau \mid \mathbf{Q} = \mathbf{q})] f_a(\mathbf{q}) d\mathbf{q} \\
&\quad - \int [E_b(D\tau \mid \mathbf{Q} = \mathbf{q}) - E(D \mid \mathbf{Q} = \mathbf{q}) E_b(\tau \mid \mathbf{Q} = \mathbf{q})] f_b(\mathbf{q}) d\mathbf{q} \\
&= E_a[\text{Cov}_a(D, \tau \mid \mathbf{Q})] - E_b[\text{Cov}_b(D, \tau \mid \mathbf{Q})] \\
&\quad + \int [E_a(D \mid \mathbf{Q} = \mathbf{q}) - E(D \mid \mathbf{Q} = \mathbf{q})] E_a(\tau \mid \mathbf{Q} = \mathbf{q}) f_a(\mathbf{q}) d\mathbf{q} \\
&\quad - \int [E_b(D \mid \mathbf{Q} = \mathbf{q}) - E(D \mid \mathbf{Q} = \mathbf{q})] E_b(\tau \mid \mathbf{Q} = \mathbf{q}) f_b(\mathbf{q}) d\mathbf{q} \\
&= E_a[\text{Cov}_a(D, \tau \mid \mathbf{Q})] - E_b[\text{Cov}_b(D, \tau \mid \mathbf{Q})]
\end{aligned}$$

$$\begin{aligned}
& + \int \Pr(G = b \mid \mathbf{Q} = \mathbf{q}) [E_a(D \mid \mathbf{Q} = \mathbf{q}) - E_b(D \mid \mathbf{Q} = \mathbf{q})] E_a(\tau \mid \mathbf{Q} = \mathbf{q}) f_a(\mathbf{q}) d\mathbf{q} \\
& - \int \Pr(G = a \mid \mathbf{Q} = \mathbf{q}) [E_b(D \mid \mathbf{Q} = \mathbf{q}) - E_a(D \mid \mathbf{Q} = \mathbf{q})] E_b(\tau \mid \mathbf{Q} = \mathbf{q}) f_b(\mathbf{q}) d\mathbf{q} \\
& = E_a[\text{Cov}_a(D, \tau \mid \mathbf{Q})] - E_b[\text{Cov}_b(D, \tau \mid \mathbf{Q})] + \int [E_a(D \mid \mathbf{Q} = \mathbf{q}) - E_b(D \mid \mathbf{Q} = \mathbf{q})] \\
& \quad [E_a(\tau \mid \mathbf{Q} = \mathbf{q}) f_a(\mathbf{q}) \Pr(G = b \mid \mathbf{Q} = \mathbf{q}) + E_b(\tau \mid \mathbf{Q} = \mathbf{q}) f_b(\mathbf{q}) \Pr(G = a \mid \mathbf{Q} = \mathbf{q})] d\mathbf{q}.
\end{aligned}$$

## Appendix B: Efficient Influence Functions

We use the Gateaux derivative approach to derive the EIFs (Ichimura and Newey, 2022), which results in more succinct derivation than the approach traditionally used in the semiparametric causal inference literature (e.g., Hahn, 1998). To further simplify the derivation, we leverage some practical rules of calculating Gateaux derivatives (Hines et al., 2022; Kennedy, 2022).

Let  $\mathbb{1}_{\tilde{o}}(o)$  be the point mass density at a single empirical observation,  $\tilde{o}$ . We also use  $\mathbb{1}(\cdot)$  as the indicator function. Let subscript  $\mathcal{P}_t$  indicate a regular parametric submodel indexed by  $t$ . When the true model is invoked, the subscript is omitted. By construction,  $f_{\mathcal{P}_t}(o) = t\mathbb{1}_{\tilde{o}}(o) + (1 - t)f(o)$ , i.e., the submodel is the true model perturbed in the direction of a single observation  $\tilde{o}$ . Under this construction, the EIF of an estimand,  $\xi$ , is the Gateaux derivative at the truth, i.e.,  $\phi(\xi) = \frac{\partial \xi_{\mathcal{P}_t}}{\partial t} \Big|_{t=0}$ . For an arbitrary function  $g(o)$ , we denote  $\frac{\partial g_{\mathcal{P}_t}(o)}{\partial t} \Big|_{t=0}$  as  $\partial g(o)$ .

We derive the EIFs for the general case of weighted estimands. Let  $w(\mathbf{X}, G)$  be the survey weight. Following Hirano et al. (2003), we assume the survey weight is a known function of the covariates. When no survey weights are needed,  $w(\mathbf{X}, G)$  reduces to 1 for every individual.

In this derivation, we also use the following definitions:

$$\begin{aligned}
h_g &:= E(w(\mathbf{X}, G) \mid G = g) \\
h_g(\mathbf{Q}) &:= E(w(\mathbf{X}, G) \mid \mathbf{Q}, G = g).
\end{aligned}$$

### B.1. EIFs for the unconditional decomposition

First of all, note that we only need to derive EIFs for two generic functions,  $\xi_{dg} := E\left(Y_d \frac{w(\mathbf{X}, g)}{h_g} \mid G = g\right)$  for an arbitrary group  $g$ ; and  $\xi_{dgg'} := E\left(Y_d \frac{w(\mathbf{X}, g)}{h_g} \mid G = g\right) E\left(D \frac{w(\mathbf{X}, g')}{h_{g'}} \mid G = g'\right)$  for two arbitrary groups  $g$  and  $g'$ , which may be the same group, and an arbitrary treatment status  $d$ . The EIFs for the decomposition components then follow from adding and subtracting these functions evaluated at appropriate  $g$ ,  $g'$ , and  $d$  values. Under conditional ignorability and SUTVA

assumptions, these estimands can be written as the following functionals:

$$\begin{aligned}\xi_{dg} &= \mathbb{E} \left[ \mu(d, \mathbf{X}, g) \frac{w(\mathbf{X}, g)}{h_g} \mid G = g \right] \\ \xi_{dgg'} &= \mathbb{E} \left[ \mu(d, \mathbf{X}, g) \frac{w(\mathbf{X}, g)}{h_g} \mid G = g \right] \mathbb{E} \left[ D \frac{w(\mathbf{X}, g')}{h_{g'}} \mid G = g' \right].\end{aligned}$$

We will also rely on the overlap assumption below, as  $\pi(d, \mathbf{X}, g)$  will appear in the denominator.

We start with  $\xi_{dg}$ .

$$\begin{aligned}\phi(\xi_{dg}) &= \partial \mathbb{E}_{\mathcal{P}_t} \left[ \mu_{\mathcal{P}_t}(d, \mathbf{X}, g) \frac{w(\mathbf{X}, g)}{\mathbb{E}_{\mathcal{P}_t}(w(\mathbf{X}, g) \mid G = g)} \mid G = g \right] \\ &= \frac{1}{h_g} \partial \mathbb{E}_{\mathcal{P}_t} [\mu_{\mathcal{P}_t}(d, \mathbf{X}, g) w(\mathbf{X}, g) \mid G = g] + \mathbb{E}[\mu(d, \mathbf{X}, g) w(\mathbf{X}, g) \mid G = g] \partial \frac{1}{\mathbb{E}_{\mathcal{P}_t}(w(\mathbf{X}, g) \mid G = g)} \\ &= \frac{1}{h_g} \frac{\mathbb{1}_{\tilde{g}}(g)}{p_g} \{ \mu(d, \tilde{\mathbf{x}}, g) w(\tilde{\mathbf{x}}, g) - \mathbb{E}[\mu(d, \mathbf{X}, g) w(\mathbf{X}, g) \mid G = g] \} \\ &\quad + \frac{1}{h_g} \mathbb{E}[\partial \mu_{\mathcal{P}_t}(d, \mathbf{X}, g) w(\mathbf{X}, g) \mid G = g] - \frac{1}{(h_g)^2} \frac{\mathbb{1}_{\tilde{g}}(g)}{p_g} [w(\tilde{\mathbf{x}}, g) - h_g] \mathbb{E}[\mu(d, \mathbf{X}, g) w(\mathbf{X}, g) \mid G = g] \\ &= \frac{\mathbb{1}_{\tilde{g}}(g)}{p_g} \frac{w(\tilde{\mathbf{x}}, g)}{h_g} [\mu(d, \tilde{\mathbf{x}}, g) - \xi_{dg}] + \frac{1}{h_g} \mathbb{E}[\partial \mu_{\mathcal{P}_t}(d, \mathbf{X}, g) w(\mathbf{X}, g) \mid G = g] \\ &= \frac{\mathbb{1}_{\tilde{g}}(g)}{p_g} \frac{w(\tilde{\mathbf{x}}, g)}{h_g} [\mu(d, \tilde{\mathbf{x}}, g) - \xi_{dg}] + \frac{1}{h_g} \mathbb{E} \left\{ \frac{\mathbb{1}_{\tilde{d}, \tilde{\mathbf{x}}, \tilde{g}}(d, \mathbf{X}, g)}{f(d, \mathbf{X}, g)} [\tilde{y} - \mu(d, \mathbf{X}, g)] w(\mathbf{X}, g) \mid G = g \right\} \\ &= \frac{\mathbb{1}_{\tilde{g}}(g)}{p_g} \frac{w(\tilde{\mathbf{x}}, g)}{h_g} [\mu(d, \tilde{\mathbf{x}}, g) - \xi_{dg}] + \frac{\mathbb{1}_{\tilde{g}}(g)}{p_g} \frac{w(\tilde{\mathbf{x}}, g)}{h_g} \frac{\mathbb{1}_{\tilde{d}}(d)}{\pi(d, \tilde{\mathbf{x}}, g)} [\tilde{y} - \mu(d, \tilde{\mathbf{x}}, g)] \\ &= \frac{\mathbb{1}(G = g)}{p_g} \frac{w(\mathbf{X}, g)}{h_g} \left\{ \frac{\mathbb{1}(D = d)}{\pi(d, \mathbf{X}, g)} [Y - \mu(d, \mathbf{X}, g)] + \mu(d, \mathbf{X}, g) - \xi_{dg} \right\}\end{aligned}$$

And without survey weights,  $\phi(\xi_{dg})$  simplifies to

$$\frac{\mathbb{1}(G = g)}{p_g} \left\{ \frac{\mathbb{1}(D = d)}{\pi(d, \mathbf{X}, g)} [Y - \mu(d, \mathbf{X}, g)] + \mu(d, \mathbf{X}, g) - \xi_{dg} \right\}.$$

Now, for  $\xi_{dgg'}$ ,

$$\phi(\xi_{dgg'}) = \phi(\xi_{dg}) \mathbb{E} \left( D \frac{w(\mathbf{X}, g')}{h_{g'}} \mid G = g' \right) + \xi_{dg} \phi \left[ \mathbb{E} \left( D \frac{w(\mathbf{X}, g')}{h_{g'}} \mid G = g' \right) \right].$$

Since

$$\begin{aligned}&\phi \left[ \mathbb{E} \left( D \frac{w(\mathbf{X}, g')}{h_{g'}} \mid G = g' \right) \right] \\ &= \partial \mathbb{E}_{\mathcal{P}_t} \left( D \frac{w(\mathbf{X}, g')}{\mathbb{E}_{\mathcal{P}_t}(w(\mathbf{X}, g') \mid G = g')} \mid G = g' \right)\end{aligned}$$



$$\begin{aligned}
&= \frac{1}{h_{g'}} \frac{\mathbb{1}_{\tilde{g}}(g')}{p_{g'}} \left[ \tilde{d}w(\tilde{\mathbf{x}}, g') - \mathbb{E}(Dw(\mathbf{X}, g') \mid G = g') \right] - \frac{1}{(h_{g'})^2} \frac{\mathbb{1}_{\tilde{g}}(g')}{p_{g'}} [w(\tilde{\mathbf{x}}, g') - h_{g'}] \mathbb{E}(Dw(\mathbf{X}, g') \mid G = g') \\
&= \frac{\mathbb{1}_{\tilde{g}}(g')}{p_{g'}} \frac{w(\tilde{\mathbf{x}}, g')}{h_{g'}} \left[ \tilde{d} - \mathbb{E} \left( D \frac{w(\mathbf{X}, g')}{h_{g'}} \mid G = g' \right) \right] \\
&= \frac{\mathbb{1}(G = g')}{p_{g'}} \frac{w(\mathbf{X}, g')}{h_{g'}} \left[ D - \mathbb{E} \left( D \frac{w(\mathbf{X}, g')}{h_{g'}} \mid G = g' \right) \right],
\end{aligned}$$

we obtain the EIF for  $\xi_{dgg'}$ ,

$$\begin{aligned}
\phi(\xi_{dgg'}) &= \frac{\mathbb{1}(G = g)}{p_g} \frac{w(\mathbf{X}, g)}{h_g} \left[ \frac{\mathbb{1}(D = d)}{\pi(d, \mathbf{X}, g)} (Y - \mu(d, \mathbf{X}, g)) + \mu(d, \mathbf{X}, g) \right] \mathbb{E} \left( D \frac{w(\mathbf{X}, g')}{h_{g'}} \mid G = g' \right) \\
&\quad + \frac{\mathbb{1}(G = g')}{p_{g'}} \frac{w(\mathbf{X}, g')}{h_{g'}} \mathbb{E} \left[ Y_d \frac{w(\mathbf{X}, g)}{h_g} \mid G = g \right] \left[ D - \mathbb{E} \left( D \frac{w(\mathbf{X}, g')}{h_{g'}} \mid G = g' \right) \right] \\
&\quad - \frac{\mathbb{1}(G = g)}{p_g} \frac{w(\mathbf{X}, g)}{h_g} \xi_{dgg'}.
\end{aligned}$$

Without survey weights,  $\phi(\xi_{dgg'})$  simplifies to

$$\begin{aligned}
&\frac{\mathbb{1}(G = g)}{p_g} \left[ \frac{\mathbb{1}(D = d)}{\pi(d, \mathbf{X}, g)} (Y - \mu(d, \mathbf{X}, g)) + \mu(d, \mathbf{X}, g) \right] \mathbb{E}(D \mid G = g') \\
&\quad + \frac{\mathbb{1}(G = g')}{p_{g'}} \mathbb{E}[Y_d \mid G = g] [D - \mathbb{E}(D \mid G = g')] - \frac{\mathbb{1}(G = g)}{p_g} \xi_{dgg'}.
\end{aligned}$$

Also note that the EIF for the total disparity,  $\mathbb{E} \left( Y \frac{w(\mathbf{X}, a)}{h_a} \mid G = a \right) - \mathbb{E} \left( Y \frac{w(\mathbf{X}, b)}{h_b} \mid G = b \right)$ ,

is

$$\begin{aligned}
&\phi(\text{Total}) \\
&= \frac{\mathbb{1}(G = a)}{p_a} \frac{w(\mathbf{X}, a)}{h_a} \left[ Y - \mathbb{E} \left( Y \frac{w(\mathbf{X}, a)}{h_a} \mid G = a \right) \right] - \frac{\mathbb{1}(G = b)}{p_b} \frac{w(\mathbf{X}, b)}{h_b} \left[ Y - \mathbb{E} \left( Y \frac{w(\mathbf{X}, b)}{h_b} \mid G = b \right) \right],
\end{aligned}$$

which, without survey weights, becomes

$$\frac{\mathbb{1}(G = a)}{p_a} [Y - \mathbb{E}(Y \mid G = a)] - \frac{\mathbb{1}(G = b)}{p_b} [Y - \mathbb{E}(Y \mid G = b)].$$

Finally, the EIFs for the unconditional decomposition components are

$$\phi(\text{Baseline}) = \phi(\xi_{0a}) - \phi(\xi_{0b})$$

$$\phi(\text{Prevalence}) = \phi(\xi_{1ba}) - \phi(\xi_{1bb}) - \phi(\xi_{0ba}) + \phi(\xi_{0bb})$$

$$\phi(\text{Effect}) = \phi(\xi_{1aa}) - \phi(\xi_{0aa}) - \phi(\xi_{1ba}) + \phi(\xi_{0ba})$$

$$\phi(\text{Selection}) = \phi(\text{Total}) - \phi(\text{Baseline}) - \phi(\text{Prevalence}) - \phi(\text{Effect}).$$

## B.2. EIFs for the conditional decomposition

Similarly to the unconditional case, we focus on the generic function

$$\xi_{dgg'g''} := \mathbb{E} \left[ \mathbb{E} \left( Y_d \frac{w(\mathbf{X}, g)}{h_g(\mathbf{Q})} \mid \mathbf{Q}, g \right) \mathbb{E} \left( D \frac{w(\mathbf{X}, g')}{h_{g'}(\mathbf{Q})} \mid \mathbf{Q}, g' \right) \frac{w(\mathbf{X}, g'')}{h_{g''}} \mid G = g'' \right],$$

where  $(g, g', g'')$  is an arbitrary combination of group memberships out of the 8 possible combinations. We still maintain the conditional ignorability, SUTVA, and overlap assumptions.

$$\begin{aligned} & \phi(\xi_{dgg'g''}) \\ &= \partial \mathbb{E}_{\mathcal{P}_t} \left[ \mathbb{E}_{\mathcal{P}_t} \left( Y_d \frac{w(\mathbf{X}, g)}{\mathbb{E}_{\mathcal{P}_t}(w(\mathbf{X}, g) \mid \mathbf{Q}, g)} \mid \mathbf{Q}, g \right) \mathbb{E}_{\mathcal{P}_t} \left( D \frac{w(\mathbf{X}, g')}{\mathbb{E}_{\mathcal{P}_t}(w(\mathbf{X}, g') \mid \mathbf{Q}, g')} \mid \mathbf{Q}, g' \right) \right. \\ & \quad \left. w(\mathbf{X}, g'') \mid G = g'' \right] \frac{1}{\mathbb{E}_{\mathcal{P}_t}(w(\mathbf{X}, g) \mid g'')} \\ &= \partial \mathbb{E}_{\mathcal{P}_t} \left[ \mathbb{E}_{\mathcal{P}_t} \left( Y_d \frac{w(\mathbf{X}, g)}{\mathbb{E}_{\mathcal{P}_t}(w(\mathbf{X}, g) \mid \mathbf{Q}, g)} \mid \mathbf{Q}, g \right) \mathbb{E}_{\mathcal{P}_t} \left( D \frac{w(\mathbf{X}, g')}{\mathbb{E}_{\mathcal{P}_t}(w(\mathbf{X}, g') \mid \mathbf{Q}, g')} \mid \mathbf{Q}, g' \right) \right. \\ & \quad \left. w(\mathbf{X}, g'') \mid G = g'' \right] \frac{1}{h_{g''}} \\ & \quad + \mathbb{E} \left[ \mathbb{E} \left( Y_d \frac{w(\mathbf{X}, g)}{\mathbb{E}(w(\mathbf{X}, g) \mid \mathbf{Q}, g)} \mid \mathbf{Q}, g \right) \mathbb{E} \left( D \frac{w(\mathbf{X}, g')}{\mathbb{E}(w(\mathbf{X}, g') \mid \mathbf{Q}, g')} \mid \mathbf{Q}, g' \right) w(\mathbf{X}, g'') \mid G = g'' \right] \\ & \quad \frac{1}{\partial \mathbb{E}_{\mathcal{P}_t}(w(\mathbf{X}, g) \mid g'')} \\ &= \frac{\mathbb{1}_{\tilde{g}(g'')}}{p_{g''}} \frac{w(\mathbf{X}, g'')}{h_{g''}} \mathbb{E} \left( Y_d \frac{w(\mathbf{X}, g)}{h_g(\mathbf{Q})} \mid \mathbf{Q}, g \right) \mathbb{E} \left( D \frac{w(\mathbf{X}, g')}{h_{g'}(\mathbf{Q})} \mid \mathbf{Q}, g' \right) - \frac{\mathbb{1}_{\tilde{g}(g'')}}{p_{g''}} \xi_{dgg'g''} \\ & \quad + \mathbb{E} \left[ \partial \mathbb{E}_{\mathcal{P}_t} \left( Y_d \frac{w(\mathbf{X}, g)}{\mathbb{E}_{\mathcal{P}_t}(w(\mathbf{X}, g) \mid \mathbf{Q}, g)} \mid \mathbf{Q}, g \right) \mathbb{E}_{\mathcal{P}_t} \left( D \frac{w(\mathbf{X}, g')}{\mathbb{E}_{\mathcal{P}_t}(w(\mathbf{X}, g') \mid \mathbf{Q}, g')} \mid \mathbf{Q}, g' \right) \right. \\ & \quad \left. w(\mathbf{X}, g'') \mid G = g'' \right] \frac{1}{h_{g''}} - \xi_{dgg'g''} h_{g''} \frac{1}{(h_{g''})^2} \frac{\mathbb{1}_{\tilde{g}(g'')}}{p_{g''}} [w(\mathbf{X}, g'') - h_{g''}] \\ &= \frac{\mathbb{1}(G = g'')}{p_{g''}} \frac{w(\mathbf{X}, g'')}{h_{g''}} \left[ \mathbb{E} \left( Y_d \frac{w(\mathbf{X}, g)}{h_g(\mathbf{Q})} \mid \mathbf{Q}, g \right) \mathbb{E} \left( D \frac{w(\mathbf{X}, g')}{h_{g'}(\mathbf{Q})} \mid \mathbf{Q}, g' \right) - \xi_{dgg'g''} \right] \\ & \quad + \mathbb{E} \left[ \partial \mathbb{E}_{\mathcal{P}_t} \left( Y_d \frac{w(\mathbf{X}, g)}{\mathbb{E}_{\mathcal{P}_t}(w(\mathbf{X}, g) \mid \mathbf{Q}, g)} \mid \mathbf{Q}, g \right) \mathbb{E}_{\mathcal{P}_t} \left( D \frac{w(\mathbf{X}, g')}{\mathbb{E}_{\mathcal{P}_t}(w(\mathbf{X}, g') \mid \mathbf{Q}, g')} \mid \mathbf{Q}, g' \right) \right. \\ & \quad \left. w(\mathbf{X}, g'') \mid G = g'' \right] \frac{1}{h_{g''}}. \end{aligned}$$

And

$$\begin{aligned} & \mathbb{E} \left[ \partial \mathbb{E}_{\mathcal{P}_t} \left( Y_d \frac{w(\mathbf{X}, g)}{\mathbb{E}_{\mathcal{P}_t}(w(\mathbf{X}, g) \mid \mathbf{Q}, g)} \mid \mathbf{Q}, g \right) \mathbb{E}_{\mathcal{P}_t} \left( D \frac{w(\mathbf{X}, g')}{\mathbb{E}_{\mathcal{P}_t}(w(\mathbf{X}, g') \mid \mathbf{Q}, g')} \mid \mathbf{Q}, g' \right) w(\mathbf{X}, g'') \mid G = g'' \right] \frac{1}{h_{g''}} \\ &= \mathbb{E} \left[ \partial \mathbb{E}_{\mathcal{P}_t} \left( Y_d \frac{w(\mathbf{X}, g)}{\mathbb{E}_{\mathcal{P}_t}(w(\mathbf{X}, g) \mid \mathbf{Q}, g)} \mid \mathbf{Q}, g \right) \mathbb{E} \left( D \frac{w(\mathbf{X}, g')}{h_{g'}(\mathbf{Q})} \mid \mathbf{Q}, g' \right) w(\mathbf{X}, g'') \mid G = g'' \right] \frac{1}{h_{g''}} \end{aligned}$$

$$\begin{aligned}
& + \mathbb{E} \left[ \mathbb{E} \left( Y_d \frac{w(\mathbf{X}, g)}{h_g(\mathbf{Q})} \mid \mathbf{Q}, g \right) \partial \mathbb{E}_{\mathcal{P}_t} \left( D \frac{w(\mathbf{X}, g')}{\mathbb{E}_{\mathcal{P}_t}(\mathbf{X}, g')} \mid \mathbf{Q}, g' \right) w(\mathbf{X}, g'') \mid G = g'' \right] \frac{1}{h_{g''}} \\
& = \frac{\mathbb{1}(G = g) p_{g''}(\mathbf{Q})}{p_g(\mathbf{Q}) p_{g''}} \frac{w(\mathbf{X}, g)}{h_g(\mathbf{Q})} \frac{w(\mathbf{X}, g'')}{h_{g''}} \left\{ \frac{\mathbb{1}(D = d)}{\pi(d, \mathbf{X}, g)} [Y - \mu(d, \mathbf{X}, g)] + \mu(d, \mathbf{X}, g) \right. \\
& \quad \left. - \mathbb{E} \left( Y_d \frac{w(\mathbf{X}, Q, g)}{h_g(\mathbf{Q})} \mid \mathbf{Q}, g \right) \right\} \mathbb{E} \left( D \frac{w(\mathbf{X}, Q, g')}{h_{g'}(\mathbf{Q})} \mid \mathbf{Q}, g' \right) \\
& \quad + \frac{\mathbb{1}(G = g') p_{g''}(\mathbf{Q})}{p_{g'}(\mathbf{Q}) p_{g''}} \frac{w(\mathbf{X}, g')}{h_{g'}(\mathbf{Q})} \frac{w(\mathbf{X}, g'')}{h_{g''}} \left[ D - \mathbb{E} \left( D \frac{w(\mathbf{X}, g')}{h_{g'}(\mathbf{Q})} \mid \mathbf{Q}, g' \right) \right] \mathbb{E} \left( Y_d \frac{w(\mathbf{X}, g)}{h_g(\mathbf{Q})} \mid \mathbf{Q}, g \right).
\end{aligned}$$

Hence,

$$\begin{aligned}
& \phi(\xi_{dgg'g''}) \\
& = \frac{\mathbb{1}(G = g'')}{p_{g''}} \frac{w(\mathbf{X}, g'')}{h_{g''}} \left[ \mathbb{E} \left( Y_d \frac{w(\mathbf{X}, g)}{h_g(\mathbf{Q})} \mid \mathbf{Q}, g \right) \mathbb{E} \left( D \frac{w(\mathbf{X}, g')}{h_{g'}(\mathbf{Q})} \mid \mathbf{Q}, g' \right) - \xi_{dgg'g''} \right] \\
& \quad + \frac{w(\mathbf{X}, g'')}{h_{g''}} \frac{\mathbb{1}(G = g) p_{g''}(\mathbf{Q})}{p_g(\mathbf{Q}) p_{g''}} \frac{w(\mathbf{X}, g)}{h_g(\mathbf{Q})} \left\{ \frac{\mathbb{1}(D = d)}{\pi(d, \mathbf{X}, g)} [Y - \mu(d, \mathbf{X}, g)] + \mu(d, \mathbf{X}, g) \right. \\
& \quad \left. - \mathbb{E} \left( Y_d \frac{w(\mathbf{X}, Q, g)}{h_g(\mathbf{Q})} \mid \mathbf{Q}, g \right) \right\} \mathbb{E} \left( D \frac{w(\mathbf{X}, Q, g')}{h_{g'}(\mathbf{Q})} \mid \mathbf{Q}, g' \right) \\
& \quad + \frac{w(\mathbf{X}, g'')}{h_{g''}} \frac{\mathbb{1}(G = g') p_{g''}(\mathbf{Q})}{p_{g'}(\mathbf{Q}) p_{g''}} \frac{w(\mathbf{X}, g')}{h_{g'}(\mathbf{Q})} \left[ D - \mathbb{E} \left( D \frac{w(\mathbf{X}, g')}{h_{g'}(\mathbf{Q})} \mid \mathbf{Q}, g' \right) \right] \mathbb{E} \left( Y_d \frac{w(\mathbf{X}, g)}{h_g(\mathbf{Q})} \mid \mathbf{Q}, g \right),
\end{aligned}$$

which, in the case of no survey weights, simplifies to

$$\begin{aligned}
& \frac{\mathbb{1}(G = g'')}{p_{g''}} [\mathbb{E}(Y_d \mid \mathbf{Q}, g) \mathbb{E}(D \mid \mathbf{Q}, g') - \xi_{dgg'g''}] \\
& \quad + \frac{\mathbb{1}(G = g) p_{g''}(\mathbf{Q})}{p_g(\mathbf{Q}) p_{g''}} \left\{ \frac{\mathbb{1}(D = d)}{\pi(d, \mathbf{X}, g)} [Y - \mu(d, \mathbf{X}, g)] + \mu(d, \mathbf{X}, g) - \mathbb{E}(Y_d \mid \mathbf{Q}, g) \right\} \mathbb{E}(D \mid \mathbf{Q}, g') \\
& \quad + \frac{\mathbb{1}(G = g') p_{g''}(\mathbf{Q})}{p_{g'}(\mathbf{Q}) p_{g''}} [D - \mathbb{E}(D \mid \mathbf{Q}, g')] \mathbb{E}(Y_d \mid \mathbf{Q}, g).
\end{aligned}$$

## Appendix C. Asymptotic inference

We follow the procedure of using the von Mises expansion to prove asymptotic properties of cross-fitting EIF-based one-step estimators (Hines et al., 2022; Kennedy, 2022; Fisher and Kennedy, 2021). In order for the cross-fitting one-step estimator to be  $\sqrt{n}$ -consistent, asymptotically normal, and semiparametrically efficient, we just need two conditions to hold. That is, both the empirical process term and the “remainder term” in the von Mises expansion are  $o_p(n^{-1/2})$ .

We use the notation  $\mathbb{P}(f(O)) := \int f(O) d\mathbb{P}(O)$ , and  $\mathbb{P}_n$  denotes the corresponding sample average. Also let  $\|\cdot\|$  denote the  $L_2$ -norm, such that  $\|f(O)\|^2 = \mathbb{P}(f(O)^2)$ . And  $\hat{\xi}$

is defined to be a substitution estimator for  $\xi$ . Formally, for all  $d, g$  and  $\xi$ , we need  $(\mathbb{P} - \mathbb{P}_n) \left[ \hat{\phi}(Y, d, \mathbf{X}, g) - \phi(Y, d, \mathbf{X}, g) \right] = o_p(n^{-1/2})$ , and  $\hat{\xi} + \mathbb{P}[\hat{\phi}(Y, d, \mathbf{X}, g)] - \xi = o_p(n^{-1/2})$ . In this appendix, we prove that the assumptions specified in the main text are sufficient for the unconditional and conditional decomposition to attain this convergence result. By cross-fitting, all of  $\hat{\mu}(d, \mathbf{X}, g)$ ,  $\hat{\pi}(d, \mathbf{X}, g)$ ,  $\hat{p}_g(\mathbf{Q})$ ,  $\hat{\mathbb{E}}(Y_d \mid \mathbf{Q}, g)$ , and  $\hat{\mathbb{E}}(D \mid \mathbf{Q}, g)$  are fitted using data not in the current subsample, which we implicitly condition on throughout.

### C.1. Inference for the unconditional decomposition

First, for  $\xi_{dg}$ , the remainder term is

$$\begin{aligned}
R_{2,dg} &= \hat{\xi}_{dg} + \mathbb{P} \left\{ \frac{\mathbb{1}(G=g)}{\hat{p}_g} \left[ \frac{\mathbb{1}(D=d)}{\hat{\pi}(d, \mathbf{X}, g)} (Y - \hat{\mu}(d, \mathbf{X}, g)) + \hat{\mu}(d, \mathbf{X}, g) - \hat{\xi}_{dg} \right] \right\} - \xi_{dg} \\
&= \hat{\xi}_{dg} + \mathbb{P} \left\{ \frac{\mathbb{1}(G=g)}{\hat{p}_g} \left[ \frac{\mathbb{1}(D=d)}{\hat{\pi}(d, \mathbf{X}, g)} (Y - \hat{\mu}(d, \mathbf{X}, g)) + \hat{\mu}(d, \mathbf{X}, g) - \mu(d, \mathbf{X}, g) \right] \right\} \\
&\quad + \mathbb{P} \left[ \frac{\mathbb{1}(G=g)}{\hat{p}_g} \mu(d, \mathbf{X}, G) \right] - \mathbb{P} \left[ \frac{\mathbb{1}(G=g)}{\hat{p}_g} \hat{\xi}_{dg} \right] - \xi_{dg} \\
&= \hat{\xi}_{dg} + \mathbb{P} \left\{ \frac{\mathbb{1}(G=g)}{\hat{p}_g} \left[ 1 - \frac{\pi(d, \mathbf{X}, g)}{\hat{\pi}(d, \mathbf{X}, g)} \right] (\hat{\mu}(d, \mathbf{X}, g) - \mu(d, \mathbf{X}, g)) \right\} + \frac{p_g}{\hat{p}_g} \xi_{dg} - \frac{p_g}{\hat{p}_g} \hat{\xi}_{dg} - \xi_{dg} \\
&= \left( 1 - \frac{p_g}{\hat{p}_g} \right) (\hat{\xi}_{dg} - \xi_{dg}) + \mathbb{P} \left\{ \frac{\mathbb{1}(G=g)}{\hat{p}_g} \left[ 1 - \frac{\pi(d, \mathbf{X}, g)}{\hat{\pi}(d, \mathbf{X}, g)} \right] [\hat{\mu}(d, \mathbf{X}, g) - \mu(d, \mathbf{X}, g)] \right\}.
\end{aligned}$$

The first term is a lower order term. For the second term,

$$\begin{aligned}
&\left| \mathbb{P} \left\{ \frac{\mathbb{1}(G=g)}{\hat{p}_g} \left[ 1 - \frac{\pi(d, \mathbf{X}, g)}{\hat{\pi}(d, \mathbf{X}, g)} \right] [\hat{\mu}(d, \mathbf{X}, g) - \mu(d, \mathbf{X}, g)] \right\} \right| \\
&\leq \frac{1}{\eta \hat{p}_g} |\mathbb{P} \{ [\hat{\pi}(d, \mathbf{X}, g) - \pi(d, \mathbf{X}, g)] [\hat{\mu}(d, \mathbf{X}, g) - \mu(d, \mathbf{X}, g)] \}| \\
&\leq \frac{1}{\eta \hat{p}_g} \|\hat{\pi}(d, \mathbf{X}, g) - \pi(d, \mathbf{X}, g)\| \|\hat{\mu}(d, \mathbf{X}, g) - \mu(d, \mathbf{X}, g)\| \\
&= o_p(n^{-1/2}),
\end{aligned}$$

where the second inequality uses the Cauchy–Schwarz inequality.

For the empirical process term,

$$\begin{aligned}
&\hat{\phi}_{dg}(Y, \mathbf{X}) - \phi_{dg}(Y, \mathbf{X}) \\
&= \frac{\mathbb{1}(G=g)}{\hat{p}_g} \left[ 1 - \frac{\mathbb{1}(D=d)}{\pi(d, \mathbf{X}, g)} \right] [\hat{\mu}(d, \mathbf{X}, g) - \mu(d, \mathbf{X}, g)] \\
&\quad + \frac{\mathbb{1}(G=g)}{\hat{p}_g} \frac{\mathbb{1}(D=d) [Y - \hat{\mu}(d, \mathbf{X}, g)]}{\hat{\pi}(d, \mathbf{X}, g) \pi(d, \mathbf{X}, g)} [\pi(d, \mathbf{X}, g) - \hat{\pi}(d, \mathbf{X}, g)]
\end{aligned}$$

$$\begin{aligned}
& + \mathbb{1}(G = g) \frac{p_g - \hat{p}_g}{\hat{p}_g p_g} \left\{ \frac{\mathbb{1}(D = d)}{\pi(d, \mathbf{X}, g)} [Y - \mu(d, \mathbf{X}, g)] + \mu(d, \mathbf{X}, g) \right\} \\
& + \frac{\mathbb{1}(G = g)}{\hat{p}_g} (\xi_{dg} - \hat{\xi}_{dg}) \\
& + \mathbb{1}(G = g) \frac{\hat{p}_g - p_g}{\hat{p}_g p_g} \xi_{dg}.
\end{aligned}$$

Note that  $(\mathbb{P}_n - \mathbb{P}) \left[ \frac{\mathbb{1}(G=g)}{\hat{p}_g} (\xi_{dg} - \hat{\xi}_{dg}) \right]$  is a lower order term. Then, using the Chebyshev's inequality argument commonly used in the double ML literature ([Chernozhukov et al., 2017](#)), the empirical process term is  $o_p(n^{-1/2})$  under stated conditions.

Second, for  $\xi_{dgg'}$ ,

$$\begin{aligned}
R_{2,dgg'} &= \hat{\xi}_{dgg'} + \mathbb{P} \left\{ \frac{\mathbb{1}(G = g)}{\hat{p}_g} \left[ \frac{\mathbb{1}(D = d)}{\hat{\pi}(d, \mathbf{X}, g)} (Y - \hat{\mu}(d, \mathbf{X}, g)) + \hat{\mu}(d, \mathbf{X}, g) \right] \hat{\mathbb{E}}(D | G = g') \right. \\
&\quad \left. + \frac{\mathbb{1}(G = g')}{\hat{p}_{g'}} \hat{\mathbb{E}}(Y_d | G = g) \left[ D - \hat{\mathbb{E}}(D | G = g') \right] - \frac{\mathbb{1}(G = g)}{\hat{p}_g} \hat{\xi}_{dgg'} \right\} - \xi_{dgg'} \\
&= \mathbb{P} \left\{ \frac{\mathbb{1}(G = g)}{\hat{p}_g} \left[ 1 - \frac{\pi(d, \mathbf{X}, g)}{\hat{\pi}(d, \mathbf{X}, g)} \right] [\hat{\mu}(d, \mathbf{X}, g) - \mu(d, \mathbf{X}, g)] \right\} \hat{\mathbb{E}}(D | G = g') \\
&\quad + \left[ \frac{p_g}{\hat{p}_g} \mathbb{E}(Y_d | G = g) - \frac{p_{g'}}{\hat{p}_{g'}} \hat{\mathbb{E}}(Y_d | G = g) \right] \left[ \hat{\mathbb{E}}(D | G = g') - \mathbb{E}(D | G = g') \right] \\
&\quad + \left( 1 - \frac{p_g}{\hat{p}_g} \right) (\hat{\xi}_{dgg'} - \xi_{dgg'}),
\end{aligned}$$

where, under stated conditions, the first term is  $o_p(n^{-1/2})$ , the second term is  $o_p(1)O_p(n^{-1/2}) = o_p(n^{-1/2})$ , and the last term is again a lower order term.

Also,

$$\begin{aligned}
& \hat{\phi}_{dgg'}(Y, \mathbf{X}) - \phi_{dgg'}(Y, \mathbf{X}) \\
&= \frac{\mathbb{1}(G = g)}{\hat{p}_g} \hat{\mathbb{E}}(D | G = g') \left[ 1 - \frac{\mathbb{1}(D = d)}{\pi(d, \mathbf{X}, g)} \right] [\hat{\mu}(d, \mathbf{X}, g) - \mu(d, \mathbf{X}, g)] \\
&\quad + \frac{\mathbb{1}(G = g)}{\hat{p}_g} \hat{\mathbb{E}}(D | G = g') \frac{\mathbb{1}(D = d) [Y - \hat{\mu}(d, \mathbf{X}, g)]}{\hat{\pi}(d, \mathbf{X}, g) \pi(d, \mathbf{X}, g)} [\pi(d, \mathbf{X}, g) - \hat{\pi}(d, \mathbf{X}, g)] \\
&\quad + \frac{\mathbb{1}(G = g)}{\hat{p}_g} \left[ \hat{\mathbb{E}}(D | G = g') - \mathbb{E}(D | G = g') \right] \left\{ \frac{\mathbb{1}(D = d)}{\pi(d, \mathbf{X}, g)} [Y - \mu(d, \mathbf{X}, g)] + \mu(d, \mathbf{X}, g) \right\} \\
&\quad + \mathbb{1}(G = g) \frac{p_g - \hat{p}_g}{\hat{p}_g p_g} \mathbb{E}(D | G = g') \left\{ \frac{\mathbb{1}(D = d)}{\pi(d, \mathbf{X}, g)} [Y - \mu(d, \mathbf{X}, g)] + \mu(d, \mathbf{X}, g) \right\} \\
&\quad + \frac{\mathbb{1}(G = g')}{\hat{p}_{g'}} \left[ \hat{\mathbb{E}}(Y_d | G = g) - \mathbb{E}(Y_d | G = g) \right] D \\
&\quad + \mathbb{1}(G = g') \frac{p_g - \hat{p}_g}{\hat{p}_g p_g} \mathbb{E}(Y_d | G = g) D
\end{aligned}$$

$$\begin{aligned}
& + \frac{\mathbb{1}(G = g')}{p_{g'}} \mathbb{E}(Y_d | G = g) \left[ \hat{\mathbb{E}}(D | G = g') - \mathbb{E}(D | G = g') \right] \\
& + \frac{\mathbb{1}(G = g')}{\hat{p}_{g'}} \left[ \hat{\mathbb{E}}(Y_d | G = g) - \mathbb{E}(Y_d | G = g) \right] \hat{\mathbb{E}}(D | G = g') \\
& + \mathbb{1}(G = g') \frac{p_{g'} - \hat{p}_{g'}}{\hat{p}_{g'} p_{g'}} \mathbb{E}(Y_d | G = g) \hat{\mathbb{E}}(D | G = g') \\
& + \mathbb{1}(G = g) \frac{\hat{p}_g - p_g}{\hat{p}_g p_g} \xi_{dgg'} \\
& + \frac{\mathbb{1}(G = g)}{\hat{p}_g} (\xi_{dgg'} - \hat{\xi}_{dgg'}).
\end{aligned}$$

Thus, the empirical process term,  $(\mathbb{P}_n - \mathbb{P}) \left[ \hat{\phi}_{dgg'}(Y, \mathbf{X}) - \phi_{dgg'}(Y, \mathbf{X}) \right]$ , is  $o_p(n^{-1/2})$  under stated conditions. Note that

$$\hat{\mathbb{E}}(Y_d | G = g) = \mathbb{P}_n \left\{ \frac{\mathbb{1}(D = d)}{\hat{\pi}(d, \mathbf{X}, g)} [Y - \hat{\mu}(d, \mathbf{X}, g)] + \hat{\mu}(d, \mathbf{X}, g) \right\},$$

hence consistent estimation of  $\pi(d, \mathbf{X}, g)$  and  $\mu(d, \mathbf{X}, g)$  makes sure that  $\hat{\mathbb{E}}(Y_d | G = g) - \mathbb{E}(Y_d | G = g) = o_p(1)$ .

## C.2. Inference for the conditional decomposition

For the components of our conditional decomposition, either  $g = g''$  or  $g' = g''$ . In what follows, we first show that the empirical process term,  $(\mathbb{P}_n - \mathbb{P}) \left[ \hat{\phi}_{dgg'g''}(Y, \mathbf{X}, \mathbf{Q}) - \phi_{dgg'g''}(Y, \mathbf{X}, \mathbf{Q}) \right]$ , is  $o_p(n^{-1/2})$ . Then, we show that the remainder term is also  $o_p(n^{-1/2})$  in both cases relevant to us, i.e., when  $g = g''$  and when  $g' = g''$ .

For the empirical process term,

$$\begin{aligned}
& \hat{\phi}_{dgg'g''}(Y, \mathbf{X}, \mathbf{Q}) - \phi_{dgg'g''}(Y, \mathbf{X}, \mathbf{Q}) \\
& = \frac{\mathbb{1}(G = g'')}{\hat{p}_{g''}} \left( \xi_{dgg'g''} - \hat{\xi}_{dgg'g''} \right) \\
& + \mathbb{1}(G = g'') \frac{\hat{p}_{g''} - p_{g''}}{\hat{p}_{g''} p_{g''}} \xi_{dgg'g''} \\
& + \mathbb{1}(G = g'') \frac{p_{g''} - \hat{p}_{g''}}{p_{g''} \hat{p}_{g''}} \mathbb{E}(Y_d | \mathbf{Q}, g) \mathbb{E}(D | \mathbf{Q}, g') \\
& + \frac{\mathbb{1}(G = g'')}{\hat{p}_{g''}} \mathbb{E}(Y_d | \mathbf{Q}, g) \left[ \hat{\mathbb{E}}(D | \mathbf{Q}, g') - \mathbb{E}(D | \mathbf{Q}, g') \right] \\
& + \frac{\mathbb{1}(G = g'')}{\hat{p}_{g''}} \hat{\mathbb{E}}(D | \mathbf{Q}, g') \left[ \hat{\mathbb{E}}(Y_d | \mathbf{Q}, g) - \mathbb{E}(Y_d | \mathbf{Q}, g) \right]
\end{aligned}$$

$$\begin{aligned}
& + \frac{\mathbb{1}(G = g) \hat{p}_{g''}(\mathbf{Q})}{\hat{p}_g(\mathbf{Q}) \hat{p}_{g''}} \mathbb{E}(Y_d | \mathbf{Q}, g) \left[ \mathbb{E}(D | \mathbf{Q}, g') - \hat{\mathbb{E}}(D | \mathbf{Q}, g') \right] \\
& + \frac{\mathbb{1}(G = g) \hat{p}_{g''}(\mathbf{Q})}{\hat{p}_g(\mathbf{Q}) \hat{p}_{g''}} \hat{\mathbb{E}}(D | \mathbf{Q}, g') \left[ \mathbb{E}(Y_d | \mathbf{Q}, g) - \hat{\mathbb{E}}(Y_d | \mathbf{Q}, g) \right] \\
& + \mathbb{1}(G = g) \frac{p_{g''}(\mathbf{Q})}{p_g(\mathbf{Q})} \frac{\hat{p}_{g''} - p_{g''}}{\hat{p}_{g''} p_{g''}} \mathbb{E}(Y_d | \mathbf{Q}, g) \mathbb{E}(D | \mathbf{Q}, g') \\
& + \frac{\mathbb{1}(G = g)}{\hat{p}_{g''}} \frac{1}{\hat{p}_g(\mathbf{Q})} [p_{g''}(\mathbf{Q}) - \hat{p}_{g''}(\mathbf{Q})] \mathbb{E}(Y_d | \mathbf{Q}, g) \mathbb{E}(D | \mathbf{Q}, g') \\
& + \frac{\mathbb{1}(G = g)}{\hat{p}_{g''}} \frac{\hat{p}_g(\mathbf{Q}) - p_g(\mathbf{Q})}{\hat{p}_g(\mathbf{Q}) p_g(\mathbf{Q})} p_{g''}(\mathbf{Q}) \mathbb{E}(Y_d | \mathbf{Q}, g) \mathbb{E}(D | \mathbf{Q}, g') \\
& + \frac{\mathbb{1}(G = g)}{\hat{p}_{g''}} \frac{\hat{p}_{g''}(\mathbf{Q})}{\hat{p}_g(\mathbf{Q})} \left[ 1 - \frac{\mathbb{1}(D = d)}{\pi(d, \mathbf{X}, g)} \right] [\hat{\mu}(d, \mathbf{X}, g) - \mu(d, \mathbf{X}, g)] \\
& + \frac{\mathbb{1}(G = g)}{\hat{p}_{g''}} \frac{\hat{p}_{g''}(\mathbf{Q})}{\hat{p}_g(\mathbf{Q})} \frac{\mathbb{1}(D = d)[Y - \hat{\mu}(d, \mathbf{X}, g)]}{\hat{\pi}(d, \mathbf{X}, g) \pi(d, \mathbf{X}, g)} [\pi(d, \mathbf{X}, g) - \hat{\pi}(d, \mathbf{X}, g)] \\
& + \mathbb{1}(G = g) \frac{p_{g''} - \hat{p}_{g''}}{p_{g''} \hat{p}_{g''}} \frac{p_{g''}(\mathbf{Q})}{p_g(\mathbf{Q})} \left\{ \frac{\mathbb{1}(D = d)}{\pi(d, \mathbf{X}, g)} [Y - \mu(d, \mathbf{X}, g)] + \mu(d, \mathbf{X}, g) \right\} \\
& + \frac{\mathbb{1}(G = g)}{\hat{p}_{g''}} \frac{p_g(\mathbf{Q}) - \hat{p}_g(\mathbf{Q})}{p_g(\mathbf{Q}) \hat{p}_g(\mathbf{Q})} p_{g''}(\mathbf{Q}) \left\{ \frac{\mathbb{1}(D = d)}{\pi(d, \mathbf{X}, g)} [Y - \mu(d, \mathbf{X}, g)] + \mu(d, \mathbf{X}, g) \right\} \\
& + \frac{\mathbb{1}(G = g') \hat{p}_{g''}(\mathbf{Q})}{\hat{p}_{g'}(\mathbf{Q}) \hat{p}_{g''}} D \left[ \hat{\mathbb{E}}(Y_d | \mathbf{Q}, g) - \mathbb{E}(Y_d | \mathbf{Q}, g) \right] \\
& + \mathbb{1}(G = g') \frac{p_{g''} - \hat{p}_{g''}}{p_{g''} \hat{p}_{g''}} \frac{p_{g''}(\mathbf{Q})}{p_{g'}(\mathbf{Q})} D \mathbb{E}(Y_d | \mathbf{Q}, g) \\
& + \frac{\mathbb{1}(G = g')}{\hat{p}_{g''}} \frac{1}{\hat{p}_{g'}(\mathbf{Q})} [\hat{p}_{g''}(\mathbf{Q}) - p_{g''}(\mathbf{Q})] D \mathbb{E}(Y_d | \mathbf{Q}, g) \\
& + \frac{\mathbb{1}(G = g')}{\hat{p}_{g''}} \frac{p_{g'}(\mathbf{Q}) - \hat{p}_{g'}(\mathbf{Q})}{p_{g'}(\mathbf{Q}) \hat{p}_{g'}(\mathbf{Q})} p_{g''}(\mathbf{Q}) D \mathbb{E}(Y_d | \mathbf{Q}, g) \\
& + \frac{\mathbb{1}(G = g') \hat{p}_{g''}(\mathbf{Q})}{\hat{p}_{g'}(\mathbf{Q}) \hat{p}_{g''}(\mathbf{Q})} \mathbb{E}(Y_d | \mathbf{Q}, g) \left[ \mathbb{E}(D | \mathbf{Q}, g') - \hat{\mathbb{E}}(D | \mathbf{Q}, g') \right] \\
& + \frac{\mathbb{1}(G = g') \hat{p}_{g''}(\mathbf{Q})}{\hat{p}_{g'}(\mathbf{Q}) \hat{p}_{g''}(\mathbf{Q})} \hat{\mathbb{E}}(D | \mathbf{Q}, g') \left[ \mathbb{E}(Y_d | \mathbf{Q}, g) - \hat{\mathbb{E}}(Y_d | \mathbf{Q}, g) \right] \\
& + \mathbb{1}(G = g') \frac{\hat{p}_{g''} - p_{g''}}{\hat{p}_{g''} p_{g''}} \frac{p_{g''}(\mathbf{Q})}{p_{g'}(\mathbf{Q})} \mathbb{E}(D | \mathbf{Q}, g') \mathbb{E}(Y_d | \mathbf{Q}, g) \\
& + \frac{\mathbb{1}(G = g')}{\hat{p}_{g''}} \frac{1}{\hat{p}_{g'}(\mathbf{Q})} [p_{g''}(\mathbf{Q}) - \hat{p}_{g''}(\mathbf{Q})] \mathbb{E}(D | \mathbf{Q}, g') \mathbb{E}(Y_d | \mathbf{Q}, g) \\
& + \frac{\mathbb{1}(G = g')}{\hat{p}_{g''}} \frac{\hat{p}_{g'}(\mathbf{Q}) - p_{g'}(\mathbf{Q})}{\hat{p}_{g'}(\mathbf{Q}) p_{g'}(\mathbf{Q})} p_{g''}(\mathbf{Q}) \mathbb{E}(D | \mathbf{Q}, g') \mathbb{E}(Y_d | \mathbf{Q}, g).
\end{aligned}$$

Using arguments similar to above, we can show that the empirical process term is indeed asymptotically negligible under stated conditions. Next, we turn to the remainder term.

### C.2.1 When $g = g''$

Note that for the conditional prevalence component, all  $\xi$  terms satisfy  $g = g''$ .

$$\begin{aligned}
& R_{2, dgg'g''} \\
&= \hat{\xi}_{dgg'g} - \xi_{dgg'g} \\
&\quad + \mathbb{P} \left\{ -\frac{\mathbb{1}(G=g)}{\hat{p}_g} \hat{\xi}_{dgg'g} \right\} \\
&\quad + \mathbb{P} \left\{ \frac{\mathbb{1}(G=g)}{\hat{p}_g} \left\{ \frac{\mathbb{1}(D=d)}{\hat{\pi}(d, \mathbf{X}, g)} [Y - \hat{\mu}(d, \mathbf{X}, g)] + \hat{\mu}(d, \mathbf{X}, g) \right\} \hat{\mathbb{E}}(D \mid \mathbf{Q}, g') \right\} \\
&\quad + \mathbb{P} \left\{ \frac{\mathbb{1}(G=g') \hat{p}_g(\mathbf{Q})}{\hat{p}_{g'}(\mathbf{Q}) \hat{p}_g} [D - \hat{\mathbb{E}}(D \mid \mathbf{Q}, g')] \hat{\mathbb{E}}(Y_d \mid \mathbf{Q}, g) \right\} \\
&= \left(1 - \frac{p_g}{\hat{p}_g}\right) \hat{\xi}_{dgg'g} - \xi_{dgg'g} \\
&\quad + \mathbb{P} \left\{ \frac{\mathbb{1}(G=g)}{\hat{p}_g} \left[1 - \frac{\pi(d, \mathbf{X}, g)}{\hat{\pi}(d, \mathbf{X}, g)}\right] [\hat{\mu}(d, \mathbf{X}, g) - \mu(d, \mathbf{X}, g)] \hat{\mathbb{E}}(D \mid \mathbf{Q}, g') \right\} \\
&\quad + \mathbb{P} \left\{ \frac{\mathbb{1}(G=g)}{\hat{p}_g} \mathbb{E}(Y_d \mid \mathbf{Q}, g) \hat{\mathbb{E}}(D \mid \mathbf{Q}, g') \right\} \\
&\quad + \mathbb{P} \left\{ \frac{\mathbb{1}(G=g') \hat{p}_g(\mathbf{Q})}{\hat{p}_{g'}(\mathbf{Q}) \hat{p}_g} [\mathbb{E}(D \mid \mathbf{Q}, g') - \hat{\mathbb{E}}(D \mid \mathbf{Q}, g')] \hat{\mathbb{E}}(Y_d \mid \mathbf{Q}, g) \right\} \\
&= \left(1 - \frac{p_g}{\hat{p}_g}\right) (\hat{\xi}_{dgg'g} - \xi_{dgg'g}) \\
&\quad + \mathbb{P} \left\{ \frac{\mathbb{1}(G=g)}{\hat{p}_g} \left[1 - \frac{\pi(d, \mathbf{X}, g)}{\hat{\pi}(d, \mathbf{X}, g)}\right] [\hat{\mu}(d, \mathbf{X}, g) - \mu(d, \mathbf{X}, g)] \hat{\mathbb{E}}(D \mid \mathbf{Q}, g') \right\} \\
&\quad + \mathbb{P} \left\{ \left[ \frac{\mathbb{1}(G=g)}{\hat{p}_g} \mathbb{E}(Y_d \mid \mathbf{Q}, g) - \frac{\mathbb{1}(G=g') \hat{p}_g(\mathbf{Q})}{\hat{p}_{g'}(\mathbf{Q}) \hat{p}_g} \hat{\mathbb{E}}(Y_d \mid \mathbf{Q}, g) \right] [\hat{\mathbb{E}}(D \mid \mathbf{Q}, g') - \mathbb{E}(D \mid \mathbf{Q}, g')] \right\}.
\end{aligned}$$

Then it follows from similar arguments as above that  $R_{2, dgg'g''} = o_p(n^{-1/2})$  under stated conditions.

To gain some intuition about the second convergence rate condition,

$$\left\| \mathbb{1}(G=g) \mathbb{E}(Y_d \mid \mathbf{Q}, g) - \mathbb{1}(G=g') \frac{\hat{p}_g(\mathbf{Q})}{\hat{p}_{g'}(\mathbf{Q})} \hat{\mathbb{E}}(Y_d \mid \mathbf{Q}, g) \right\| \left\| \hat{\mathbb{E}}(D \mid \mathbf{Q}, g') - \mathbb{E}(D \mid \mathbf{Q}, g') \right\| = o_p(n^{-1/2}),$$

note that if  $g = g' = g''$ ,

$$\begin{aligned}
& \mathbb{P} \left\{ \left[ \mathbb{1}(G=g) \mathbb{E}(Y_d \mid \mathbf{Q}, g) - \mathbb{1}(G=g') \frac{\hat{p}_g(\mathbf{Q})}{\hat{p}_{g'}(\mathbf{Q})} \hat{\mathbb{E}}(Y_d \mid \mathbf{Q}, g) \right] [\hat{\mathbb{E}}(D \mid \mathbf{Q}, g') - \mathbb{E}(D \mid \mathbf{Q}, g')] \right\} \\
&= \mathbb{P} \left\{ \mathbb{1}(G=g) [\mathbb{E}(Y_d \mid \mathbf{Q}, g) - \hat{\mathbb{E}}(Y_d \mid \mathbf{Q}, g)] [\hat{\mathbb{E}}(D \mid \mathbf{Q}, g') - \mathbb{E}(D \mid \mathbf{Q}, g')] \right\},
\end{aligned}$$



so  $\left\| \mathbb{E}(Y_d | \mathbf{Q}, g) - \hat{\mathbb{E}}(Y_d | \mathbf{Q}, g) \right\| \left\| \hat{\mathbb{E}}(D | \mathbf{Q}, g') - \mathbb{E}(D | \mathbf{Q}, g') \right\| = o_p(n^{-1/2})$  is sufficient to replace the second convergence rate condition.

If  $g = g'' \neq g'$ ,

$$\begin{aligned} & \mathbb{P} \left\{ \left[ \mathbb{1}(G = g) \mathbb{E}(Y_d | \mathbf{Q}, g) - \mathbb{1}(G = g') \frac{\hat{p}_g(\mathbf{Q})}{\hat{p}_{g'}(\mathbf{Q})} \hat{\mathbb{E}}(Y_d | \mathbf{Q}, g) \right] \left[ \hat{\mathbb{E}}(D | \mathbf{Q}, g') - \mathbb{E}(D | \mathbf{Q}, g') \right] \right\} \\ &= \mathbb{P} \left\{ \mathbb{1}(G = g') \frac{\hat{p}_g(\mathbf{Q})}{\hat{p}_{g'}(\mathbf{Q})} \left[ \mathbb{E}(Y_d | \mathbf{Q}, g) - \hat{\mathbb{E}}(Y_d | \mathbf{Q}, g) \right] \left[ \hat{\mathbb{E}}(D | \mathbf{Q}, g') - \mathbb{E}(D | \mathbf{Q}, g') \right] \right\} \\ &+ \mathbb{P} \left\{ \frac{\hat{p}_g(\mathbf{Q})}{\hat{p}_{g'}(\mathbf{Q})} [\mathbb{1}(G = g) - \mathbb{1}(G = g')] \mathbb{E}(Y_d | \mathbf{Q}, g) \left[ \hat{\mathbb{E}}(D | \mathbf{Q}, g') - \mathbb{E}(D | \mathbf{Q}, g') \right] \right\} \\ &+ \mathbb{P} \left\{ \mathbb{1}(G = g) \left[ 1 - \frac{\hat{p}_g(\mathbf{Q})}{\hat{p}_{g'}(\mathbf{Q})} \right] \mathbb{E}(Y_d | \mathbf{Q}, g) \left[ \hat{\mathbb{E}}(D | \mathbf{Q}, g') - \mathbb{E}(D | \mathbf{Q}, g') \right] \right\}, \end{aligned}$$

so the following conditions are sufficient to replace the second convergence rate condition: for a constant  $\zeta > 0$ ,  $\left| \frac{\hat{p}_g(\mathbf{Q})}{\hat{p}_{g'}(\mathbf{Q})} \right| \leq \zeta$  with probability 1,  $\left\| \hat{\mathbb{E}}(D | \mathbf{Q}, g') - \mathbb{E}(D | \mathbf{Q}, g') \right\| = o_p(n^{-1/2})$ , and  $\left\| \mathbb{E}(Y_d | \mathbf{Q}, g) - \hat{\mathbb{E}}(Y_d | \mathbf{Q}, g) \right\| = o_p(1)$ .

### C.2.2 When $g' = g''$

Note that all  $\xi$  terms satisfy  $g' = g''$  for the conditional effect component.

$$\begin{aligned} & R_{2, dgg'g''} \\ &= \hat{\xi}_{dgg'g'} - \xi_{dgg'g'} \\ &+ \mathbb{P} \left\{ -\frac{\mathbb{1}(G = g')}{\hat{p}_{g'}} \hat{\xi}_{dgg'g'} \right\} \\ &+ \mathbb{P} \left\{ \frac{\mathbb{1}(G = g) \hat{p}_{g'}(\mathbf{Q})}{\hat{p}_g(\mathbf{Q}) \hat{p}_{g'}} \left\{ \frac{\mathbb{1}(D = d)}{\hat{\pi}(d, \mathbf{X}, g)} [Y - \hat{\mu}(d, \mathbf{X}, g)] + \hat{\mu}(d, \mathbf{X}, g) - \hat{\mathbb{E}}(Y_d | \mathbf{Q}, g) \right\} \hat{\mathbb{E}}(D | \mathbf{Q}, g') \right\} \\ &+ \mathbb{P} \left\{ \frac{\mathbb{1}(G = g')}{\hat{p}_{g'}} D \cdot \hat{\mathbb{E}}(Y_d | \mathbf{Q}, g) \right\} \\ &= \left( 1 - \frac{p_g}{\hat{p}_g} \right) \left( \hat{\xi}_{dgg'g'} - \xi_{dgg'g'} \right) \\ &+ \mathbb{P} \left\{ \frac{\mathbb{1}(G = g) \hat{p}_{g'}(\mathbf{Q})}{\hat{p}_g(\mathbf{Q}) \hat{p}_{g'}} \left[ 1 - \frac{\pi(d, \mathbf{X}, g)}{\hat{\pi}(d, \mathbf{X}, g)} \right] [\hat{\mu}(d, \mathbf{X}, g) - \mu(d, \mathbf{X}, g)] \hat{\mathbb{E}}(D | \mathbf{Q}, g') \right\} \\ &+ \mathbb{P} \left\{ \left[ \frac{\mathbb{1}(G = g')}{\hat{p}_{g'}} \mathbb{E}(D | \mathbf{Q}, g') - \frac{\mathbb{1}(G = g) \hat{p}_{g'}(\mathbf{Q})}{\hat{p}_g(\mathbf{Q}) \hat{p}_{g'}} \hat{\mathbb{E}}(D | \mathbf{Q}, g') \right] \left[ \hat{\mathbb{E}}(Y_d | \mathbf{Q}, g) - \mathbb{E}(Y_d | \mathbf{Q}, g) \right] \right\}. \end{aligned}$$

Under stated conditions,  $R_{2, dgg'g''} = o_p(n^{-1/2})$ .

For the second convergence rate condition, i.e.,

$$\left\| \frac{\mathbb{1}(G = g')}{\hat{p}_{g'}} \mathbb{E}(D | \mathbf{Q}, g') - \frac{\mathbb{1}(G = g) \hat{p}_{g'}(\mathbf{Q})}{\hat{p}_g(\mathbf{Q}) \hat{p}_{g'}} \hat{\mathbb{E}}(D | \mathbf{Q}, g') \right\| \left\| \hat{\mathbb{E}}(Y_d | \mathbf{Q}, g) - \mathbb{E}(Y_d | \mathbf{Q}, g) \right\|,$$

note that when  $g' = g'' \neq g$ , the following conditions are sufficient: for a constant  $\zeta > 0$ ,  $\left| \frac{\hat{p}_{g'}(\mathbf{Q})}{\hat{p}_g(\mathbf{Q})} \right| \leq \zeta$  with probability 1,  $\left\| \hat{\mathbb{E}}(D \mid \mathbf{Q}, g') - \mathbb{E}(D \mid \mathbf{Q}, g') \right\| = o_p(1)$ , and  $\left\| \mathbb{E}(Y_d \mid \mathbf{Q}, g) - \hat{\mathbb{E}}(Y_d \mid \mathbf{Q}, g) \right\| = o_p(n^{-1/2})$ .

## Appendix D. Double robustness

The proof for Theorem 3, i.e., the double robustness of  $\Psi_{dg}$  and  $\Psi_{dgg'}$ , is omitted. It is similar to the proof of the double robustness of the AIPW estimator of ATE (Robins et al., 1994; Tsiatis, 2006).

Below, we prove the double robustness of  $\Psi_{dgg'g''}$ , which is stated in Theorem 6. We use a tilde to denote the probability limit of a nuisance estimator, i.e., for example,  $\hat{\mu}(d, \mathbf{X}, g)$  converges to  $\tilde{\mu}(d, \mathbf{X}, g)$ . Under consistent estimation of  $p_g$  and  $\mathbb{E}(D \mid \mathbf{Q}, g)$ , the one-step estimator  $\Psi_{dgg'g''}$  converges in probability to

$$\begin{aligned}
& \mathbb{E} \left\{ \frac{\mathbb{1}(G = g'')}{\tilde{p}_{g''}} \tilde{\mathbb{E}}(Y_d \mid \mathbf{Q}, g) \tilde{\mathbb{E}}(D \mid \mathbf{Q}, g') \right\} \\
& + \mathbb{E} \left\{ \left( 1 - \frac{p_{g''}}{\tilde{p}_{g''}} \right) \tilde{\xi}_{dgg'g''} \right\} \\
& + \mathbb{E} \left\{ \frac{\mathbb{1}(G = g) \tilde{p}_{g''}(\mathbf{Q})}{\tilde{p}_g(\mathbf{Q}) \tilde{p}_{g''}} \left\{ \frac{\mathbb{1}(D = d)}{\tilde{\pi}(d, \mathbf{X}, g)} [Y - \tilde{\mu}(d, \mathbf{X}, g)] + \tilde{\mu}(d, \mathbf{X}, g) - \tilde{\mathbb{E}}(Y_d \mid \mathbf{Q}, g) \right\} \tilde{\mathbb{E}}(D \mid \mathbf{Q}, g') \right\} \\
& + \mathbb{E} \left\{ \frac{\mathbb{1}(G = g') \tilde{p}_{g''}(\mathbf{Q})}{\tilde{p}_{g'}(\mathbf{Q}) \tilde{p}_{g''}} [D - \tilde{\mathbb{E}}(D \mid \mathbf{Q}, g')] \tilde{\mathbb{E}}(Y_d \mid \mathbf{Q}, g) \right\} \\
& = \mathbb{E} \left\{ \frac{\mathbb{1}(G = g'')}{p_{g''}} \tilde{\mathbb{E}}(Y_d \mid \mathbf{Q}, g) \mathbb{E}(D \mid \mathbf{Q}, g') \right\} \\
& + \mathbb{E} \left\{ \frac{\mathbb{1}(G = g) \tilde{p}_{g''}(\mathbf{Q})}{\tilde{p}_g(\mathbf{Q}) p_{g''}} \left\{ \frac{\mathbb{1}(D = d)}{\tilde{\pi}(d, \mathbf{X}, g)} [Y - \tilde{\mu}(d, \mathbf{X}, g)] + \tilde{\mu}(d, \mathbf{X}, g) - \tilde{\mathbb{E}}(Y_d \mid \mathbf{Q}, g) \right\} \mathbb{E}(D \mid \mathbf{Q}, g') \right\} \\
& + \mathbb{E} \left\{ \frac{\mathbb{1}(G = g') \tilde{p}_{g''}(\mathbf{Q})}{\tilde{p}_{g'}(\mathbf{Q}) p_{g''}} [D - \mathbb{E}(D \mid \mathbf{Q}, g')] \tilde{\mathbb{E}}(Y_d \mid \mathbf{Q}, g) \right\} \\
& = \mathbb{E} \left\{ \frac{\mathbb{1}(G = g'')}{p_{g''}} \tilde{\mathbb{E}}(Y_d \mid \mathbf{Q}, g) \mathbb{E}(D \mid \mathbf{Q}, g') \right\} \\
& + \mathbb{E} \left\{ \frac{\mathbb{1}(G = g) \tilde{p}_{g''}(\mathbf{Q})}{\tilde{p}_g(\mathbf{Q}) p_{g''}} \frac{\pi(d, \mathbf{X}, g)}{\tilde{\pi}(d, \mathbf{X}, g)} [\mu(d, \mathbf{X}, g) - \tilde{\mu}(d, \mathbf{X}, g)] \mathbb{E}(D \mid \mathbf{Q}, g') \right\} \\
& + \mathbb{E} \left\{ \frac{\mathbb{1}(G = g) \tilde{p}_{g''}(\mathbf{Q})}{\tilde{p}_g(\mathbf{Q}) p_{g''}} [\tilde{\mu}(d, \mathbf{X}, g) - \tilde{\mathbb{E}}(Y_d \mid \mathbf{Q}, g)] \mathbb{E}(D \mid \mathbf{Q}, g') \right\}
\end{aligned}$$

Now, note that when  $g = g''$ , this probability limit becomes

$$\mathbb{E} \left\{ \frac{\mathbb{1}(G = g)}{p_g} \left\{ \frac{\pi(d, \mathbf{X}, g)}{\tilde{\pi}(d, \mathbf{X}, g)} [\mu(d, \mathbf{X}, g) - \tilde{\mu}(d, \mathbf{X}, g)] + \tilde{\mu}(d, \mathbf{X}, g) \right\} \mathbb{E}(D \mid \mathbf{Q}, g') \right\}.$$

Hence, in this case, the double-robustness is achieved without requiring consistent estimation of  $p_g(\mathbf{Q})$  and  $E(Y_d | \mathbf{Q}, g)$  for any  $d$  and  $g$ . When  $g \neq g''$ , further assuming that  $\tilde{E}(Y_d | \mathbf{Q}, g) = E(Y_d | \mathbf{Q}, g)$  if either  $\tilde{\pi}(d, \mathbf{X}, g) = \pi(d, \mathbf{X}, g)$  or  $\tilde{\mu}(d, \mathbf{X}, g) = \mu(d, \mathbf{X}, g)$ , then the probability limit of  $\Psi_{dgg'g''}$  will still be  $\xi_{dgg'g''}$ . Note that consistent estimation of  $p_g(\mathbf{Q})$  is still not required in either case.

## Appendix E. Relation between selection concepts

Various concepts of ‘selection’, ‘sorting’, and ‘selectivity’ have appeared in the social science literature on the effects of education on later-life attainments. Below, we clarify the relationship between our selection concept and those appearing in prior works. In our framework, group-specific selectivity is defined as

$$\text{Cov}_g(D, \tau) = [E_g(\tau | D = 1) - E_g(\tau)] E_g(D).$$

First, our definition is closely related to the “sorting on gains” concept in [Heckman et al. \(2018\)](#), which is defined as the difference between ATT and ATE, i.e.,

$$E_g(\tau | D = 1) - E_g(\tau).$$

Since  $E_g(D)$  is always positive, the sign of our group-specific selectivity is always the same as that of sorting on gains in Heckman et al. (2018).

Second, akin to our framework, recent works in the “great equalizer” literature ([Zhou, 2019](#); [Fiel, 2020](#); [Karlson, 2019](#)) have discussed *differential* selection into college completion by parental income. Implicitly, these authors define group-specific selection as

$$E_g(Y | D = 1) - E_g(Y | D = 0) - E_g(\tau).$$

This implicit definition can be validated by rewriting the group difference in this selection term.

$$\begin{aligned} & E_a(Y | D = 1) - E_a(Y | D = 0) - E_a(\tau) - [E_b(Y | D = 1) - E_b(Y | D = 0) - E_b(\tau)] \\ &= \underbrace{E_a(Y | D = 1) - E_b(Y | D = 1) - [E_a(Y | D = 0) - E_b(Y | D = 0)]}_{\text{The descriptive test of the great equalizer thesis}} \end{aligned}$$

$$-\underbrace{\{E_a(Y_1) - E_b(Y_1) - [E_a(Y_0) - E_b(Y_0)]\}}_{\text{The causal test of the great equalizer thesis}}$$

If the descriptive test returns a smaller value than the causal test, these authors will conclude that members of the advantaged group are less selected than their disadvantaged peers in the transition to college completion.

Taking the difference between the selectivity term in the great equalizer literature and sorting on gains, we obtain  $E_g(Y_0 \mid D = 1) - E_g(Y_0 \mid D = 0)$ , which is selection on baseline outcome. Therefore, the selectivity term in the great equalizer literature captures selection on baseline outcome, on the top of selection on treatment effect captured in sorting on gains and our selectivity term (See [Morgan and Winship, 2014](#), p.58-9).

Third, a literature on the effect heterogeneity by propensity score ([Brand and Xie, 2010](#); [Xie et al., 2012](#); [Brand et al., 2021](#)) is closely aligned with our framework in its conceptualization of selection into treatment. Originally developed in the context of education effects, this framework has been applied in a wide array of topics. Recall that  $\pi(1, \mathbf{X}, g)$  denotes the propensity score, then under Assumption 2,

$$\begin{aligned} & \text{Cov}_g(D, \tau) \\ &= E_g[\text{Cov}_g(D, \tau \mid \mathbf{X})] + \text{Cov}_g[E_g(D \mid \mathbf{X}), E_g(\tau \mid \mathbf{X})] \\ &= \text{Cov}_g[\pi(1, \mathbf{X}, g), E_g(\tau \mid \mathbf{X})] \\ &= E_g[\pi(1, \mathbf{X}, g) E_g(\tau \mid \mathbf{X})] - E_g(D) E_g(\tau) \\ &= E_g\{E_g[\pi(1, \mathbf{X}, g) E_g(\tau \mid \mathbf{X}) \mid \pi(1, \mathbf{X}, g)]\} - E_g(D) E_g(\tau) \\ &= E_g\{\pi(1, \mathbf{X}, g) E_g[\tau \mid \pi(1, \mathbf{X}, g)]\} - E_g[\pi(1, \mathbf{X}, g)] E_g\{E_g[\tau \mid \pi(1, \mathbf{X}, g)]\} \\ &= \text{Cov}_g[\pi(1, \mathbf{X}, g), E_g(\tau \mid \pi(1, \mathbf{X}, g))]. \end{aligned}$$

Hence, under conditional ignorability, our selectivity term equals the covariance between the propensity score and the conditional treatment effect given the propensity score. Due to this relationship, the estimators we developed for the generic functions underlying our unconditional decomposition can also be used to nonparametrically estimate effect heterogeneity by propensity score.

## Appendix F. Supplemental tables

Table A1. Group-specific Estimates for the Unconditional Decomposition

	Top 60% income	Bottom 40% income	Top-Bottom
Treatment proportion	0.290 (0.269, 0.312)	0.086 (0.066, 0.107)	0.204 (0.174, 0.234)
Gradient Boosted Machine			
Baseline outcome	0.575 (0.553, 0.597)	0.410 (0.392, 0.427)	0.165 (0.137, 0.194)
ATE	0.146 (0.108, 0.184)	0.128 (0.048, 0.208)	0.018 (-0.071, 0.106)
Cov( $D, \tau$ )	-0.013 (-0.028, 0.001)	0.004 (-0.003, 0.011)	-0.017 (-0.033, -0.001)
Neural Networks			
Baseline outcome	0.569 (0.555, 0.584)	0.408 (0.389, 0.427)	0.161 (0.137, 0.185)
ATE	0.152 (0.118, 0.187)	0.141 (0.074, 0.208)	0.011 (-0.064, 0.087)
Cov( $D, \tau$ )	-0.009 (-0.018, 0.000)	0.004 (-0.002, 0.010)	-0.014 (-0.024, -0.003)
Random Forests			
Baseline outcome	0.570 (0.557, 0.584)	0.410 (0.390, 0.429)	0.161 (0.137, 0.184)
ATE	0.139 (0.109, 0.169)	0.115 (0.025, 0.205)	0.024 (-0.071, 0.119)
Cov( $D, \tau$ )	-0.006 (-0.013, 0.001)	0.005 (-0.001, 0.012)	-0.011 (-0.021, -0.002)
Parametric Regressions			
Baseline outcome	0.570 (0.556, 0.584)	0.410 (0.391, 0.429)	0.160 (0.136, 0.184)
ATE	0.142 (0.107, 0.177)	0.141 (0.041, 0.241)	0.001 (-0.105, 0.107)
Cov( $D, \tau$ )	-0.007 (-0.016, 0.003)	0.003 (-0.005, 0.011)	-0.009 (-0.022, 0.003)

Note: 95% confidence intervals are in the parentheses. Weight stabilization is used. For ML models, cross-fitting is used.

Table A2. Unconditional Decomposition Estimates

	GBM	Neural Networks	Random Forests	Parametric
Total	0.180 (0.158, 0.201)	0.180 (0.158, 0.201)	0.180 (0.158, 0.201)	0.180 (0.158, 0.201)
Baseline	0.165 (0.137, 0.194)	0.161 (0.137, 0.185)	0.161 (0.137, 0.184)	0.160 (0.136, 0.184)
Prevalence	0.026 (0.010, 0.043)	0.029 (0.014, 0.043)	0.023 (0.005, 0.042)	0.029 (0.008, 0.050)
Effect	0.005 (-0.021, 0.031)	0.003 (-0.018, 0.025)	0.007 (-0.021, 0.034)	0.000 (-0.030, 0.031)
Selection	-0.017 (-0.033, -0.001)	-0.014 (-0.024, -0.003)	-0.011 (-0.021, -0.002)	-0.009 (-0.022, 0.003)
Jackson reduction	0.022 (0.001, 0.043)	0.025 (0.007, 0.042)	0.018 (-0.006, 0.042)	0.026 (-0.001, 0.054)

Note: 95% confidence intervals are in the parentheses. Weight stabilization is used.  
For ML models, cross-fitting is used.

Table A3: Conditional Decomposition Estimates

	GBM	Neural Networks	Random Forests	Parametric
Total	0.180 (0.158,0.201)	0.180 (0.158,0.201)	0.180 (0.158,0.201)	0.180 (0.158,0.201)
Baseline	0.165 (0.137,0.194)	0.161 (0.137,0.185)	0.161 (0.137,0.184)	0.160 (0.136,0.184)
Conditional prevalence	0.006 (0.002,0.011)	0.007 (0.002,0.012)	0.013 (0.008,0.018)	0.009 (0.003,0.015)
Conditional effect	-0.021 (-0.052,0.011)	-0.015 (-0.040,0.009)	-0.047 (-0.091,-0.003)	-0.018 (-0.044,0.007)
Conditional selection	-0.003 (-0.011,0.004)	-0.005 (-0.011,0.001)	-0.006 (-0.013,0.002)	-0.005 (-0.012,0.003)
Q distribution	0.032 (0.011,0.052)	0.032 (0.017,0.048)	0.059 (0.022,0.096)	0.033 (0.016,0.050)
Conditional Jackson reduction	0.007 (-0.001,0.014)	0.006 (0.000,0.013)	0.011 (0.003,0.020)	0.007 (-0.003,0.016)

Note: 95% confidence intervals are in the parentheses. Weight stabilization is used.  
For ML models, cross-fitting is used.

## References

- Bailey, Martha and Susan Dynarski. 2011. "Gains and Gaps: Changing Inequality in U.S. College Entry and Completion." Technical Report w17633, National Bureau of Economic Research, Cambridge, MA.
- Barber, Jennifer S., Jennifer Eckerman Yarger, and Heather H. Gatny. 2015. "Black-White Differences in Attitudes Related to Pregnancy Among Young Women." *Demography* 52:751–786.
- Baron, Reuben M and David A Kenny. 1986. "The Moderator-Mediator Variable Distinction in Social Psychological Research: Conceptual, Strategic, and Statistical Considerations." *Journal of Personality and Social Psychology* 51.
- Bickel, Peter J, Chris A.J Klaassen, Ya'acov Ritov, and Jon A Wellner. 1998. *Efficient and Adaptive Estimation for Semiparametric Models*. New York, NY: Springer.
- Blau, Peter Michael and Otis Dudley Duncan. 1967. *The American Occupational Structure*. Free Press.
- Blinder, Alan S. 1973. "Wage Discrimination: Reduced Form and Structural Estimates." *The Journal of Human Resources* 8:436.
- Bloome, Deirdre, Shauna Dyer, and Xiang Zhou. 2018. "Educational Inequality, Educational Expansion, and Intergenerational Income Persistence in the United States." *American Sociological Review* 83:1215–1253.
- Brady, David, Ryan M. Finnigan, and Sabine Hübgen. 2017. "Rethinking the Risks of Poverty: A Framework for Analyzing Prevalences and Penalties." *American Journal of Sociology* 123:740–786.
- Brand, Jennie E. and Yu Xie. 2010. "Who Benefits Most from College?: Evidence for Negative Selection in Heterogeneous Economic Returns to Higher Education." *American Sociological Review* 75:273–302.
- Brand, Jennie E., Jiahui Xu, Bernard Koch, and Pablo Geraldo. 2021. "Uncovering Sociological Effect Heterogeneity Using Tree-Based Machine Learning." *Sociological Methodology* 51:189–223.
- Breen, R. 2010. "Educational Expansion and Social Mobility in the 20th Century." *Social Forces* 89:365–388.
- Breen, Richard, Seungsoo Choi, and Anders Holm. 2015. "Heterogeneous Causal Effects and Sample Selection Bias." *Sociological Science* 2:351–369.
- Cheng, Siwei, Jennie E. Brand, Xiang Zhou, Yu Xie, and Michael Hout. 2021. "Heterogeneous returns to college over the life course." *Science Advances* 7:eabg7641.
- Chernozhukov, Victor, Denis Chetverikov, Mert Demirer, Esther Duflo, Christian Hansen, and Whitney Newey. 2017. "Double/Debiased/Neyman Machine Learning of Treatment Effects." *American Economic Review* 107:261–265.

- Chernozhukov, Victor, Denis Chetverikov, Mert Demirer, Esther Duflo, Christian Hansen, Whitney Newey, and James Robins. 2018a. “Double/debiased machine learning for treatment and structural parameters.” *The Econometrics Journal* 21:C1–C68.
- Chernozhukov, Victor, Iván Fernández-Val, and Ye Luo. 2018b. “The Sorted Effects Method: Discovering Heterogeneous Effects Beyond Their Averages.” *Econometrica* 86:1911–1938.
- Didelez, Vanessa, A Philip Dawid, and Sara Geneletti. 2006. “Direct and Indirect Effects of Sequential Treatments.” In *Proceedings of the Twenty-Second Conference on Uncertainty in Artificial Intelligence*, edited by R Dechter and T. S Richardson, pp. 138–146, Arlington, Virginia. AUAI Press.
- Farkas, George. 2003. “Cognitive Skills and Noncognitive Traits and Behaviors in Stratification Processes.” *Annual Review of Sociology* 29:541–562.
- Featherman, David L. and Robert M. Hauser. 1978. *Opportunity and Change*. Studies in population. Academic Press.
- Fiel, Jeremy E. 2020. “Great Equalizer or Great Selector? Reconsidering Education as a Moderator of Intergenerational Transmissions.” *Sociology of Education* p. 003804072092788.
- Fisher, Aaron and Edward H. Kennedy. 2021. “Visually Communicating and Teaching Intuition for Influence Functions.” *The American Statistician* 75:162–172.
- Fortin, Nicole, Thomas Lemieux, and Sergio Firpo. 2011. “Decomposition Methods in Economics.” In *Handbook of Labor Economics*, volume 4, pp. 1–102. Elsevier.
- Hahn, Jinyong. 1998. “On the Role of the Propensity Score in Efficient Semiparametric Estimation of Average Treatment Effects.” *Econometrica* 66:315.
- Heckman, James J, John Eric Humphries, and Gregory Veramendi. 2018. “Returns to Education: The Causal Effects of Education on Earnings, Health, and Smoking.” *Journal of Political Economy* 126:50.
- Heckman, James J., Jora Stixrud, and Sergio Urzua. 2006. “The Effects of Cognitive and Noncognitive Abilities on Labor Market Outcomes and Social Behavior.” *Journal of Labor Economics* 24:411–482.
- Heckman, James J. and Edward Vytlacil. 2005. “Structural Equations, Treatment Effects, and Econometric Policy Evaluation1.” *Econometrica* 73:669–738.
- Hines, Oliver, Oliver Dukes, Karla Diaz-Ordaz, and Stijn Vansteelandt. 2022. “Demystifying Statistical Learning Based on Efficient Influence Functions.” *The American Statistician* pp. 1–13.
- Hirano, Keisuke, Guido W Imbens, and Geert Ridder. 2003. “Efficient Estimation of Average Treatment Effects Using the Estimated Propensity Score.” *Econometrica* 71:1161–1189.
- Hout, Michael. 2012. “Social and Economic Returns to College Education in the United States.” *Annual Review of Sociology* 38:379–400.
- Huber, Martin. 2015. “Causal Pitfalls in the Decomposition of Wage Gaps.” *Journal of Business & Economic Statistics* 33:179–191.



- Ichimura, Hidehiko and Whitney K. Newey. 2022. “The influence function of semiparametric estimators.” *Quantitative Economics* 13:29–61.
- Ishida, Hiroshi, Walter Muller, and John M. Ridge. 1995. “Class Origin, Class Destination, and Education: A Cross-National Study of Ten Industrial Nations.” *American Journal of Sociology* 101:145–193.
- Jackson, John W. 2021. “Meaningful Causal Decompositions in Health Equity Research: Definition, Identification, and Estimation Through a Weighting Framework.” *Epidemiology* 32:282–290.
- Jackson, John W. and Tyler J. VanderWeele. 2018. “Decomposition Analysis to Identify Intervention Targets for Reducing Disparities:.” *Epidemiology* 29:825–835.
- Jann, Ben. 2008. “The Blinder–Oaxaca Decomposition for Linear Regression Models.” *The Stata Journal: Promoting communications on statistics and Stata* 8:453–479.
- Karlson, Kristian Bernt. 2019. “College as equalizer? Testing the selectivity hypothesis.” *Social Science Research* 80:216–229.
- Kennedy, Edward H. 2022. “Semiparametric doubly robust targeted double machine learning: a review.” Number: arXiv:2203.06469 arXiv:2203.06469 [stat].
- Kitagawa, Evelyn M. 1955. “Components of a Difference Between Two Rates.” *Journal of the american statistical association* 50:1168–1194.
- Kline, Patrick. 2011. “Oaxaca-Blinder as a Reweighting Estimator.” *American Economic Review* 101:532–537.
- Laurison, Daniel and Sam Friedman. 2016. “The Class Pay Gap in Higher Professional and Managerial Occupations.” *American Sociological Review* 81:668–695.
- Lundberg, Ian. 2022. “The Gap-Closing Estimand: A Causal Approach to Study Interventions That Close Disparities Across Social Categories.” *Sociological Methods & Research* p. 004912412110557.
- Lundberg, Ian, Rebecca Johnson, and Brandon M Stewart. 2021. “What Is Your Estimand? Defining the Target Quantity Connects Statistical Evidence to Theory.” *American Sociological Review* 86:532–565.
- Malinsky, Daniel. 2018. “Intervening on structure.” *Synthese* 195:2295–2312.
- Mare, Robert D. 1980. “Social Background and School Continuation Decisions.” *Journal of the american statistical association* 75:295–305.
- Mize, Trenton D. 2016. “Sexual Orientation in the Labor Market.” *American Sociological Review* 81:1132–1160.
- Morgan, Stephen L. and Christopher Winship. 2014. *Counterfactuals and Causal Inference: Methods And Principles For Social Research*. Analytical Methods for Social Research. Cambridge, UK: Cambridge University Press, 2nd edition edition.
- Oaxaca, Ronald. 1973. “Male-Female Wage Differentials in Urban Labor Markets.” *International Economic Review* 14:693.

- Pearl, Judea. 2001. "Direct and Indirect Effects." In *Proceedings of the Seventeenth Conference on Uncertainty in Artificial Intelligence*, p. 411{20, San Francisco, CA. Morgan Kaufmann.
- R Core Team. 2023. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.
- Reardon, Sean F. 2011. "The Widening Academic Achievement Gap Between the Rich and the Poor: New Evidence and Possible Explanations." In *Whither Opportunity? Rising Inequality, Schools, and Children's Life Chances*, edited by Greg J. Duncan and Richard J. Murnane, pp. 91–115. Russell Sage Foundation.
- Robins, James M, Andrea Rotnitzky, and Lue Ping Zhao. 1994. "Estimation of Regression Coefficients When Some Regressors Are Not Always Observed." *Journal of the American statistical Association* 89:846–866.
- Rubin, Donald B. 1974. "Estimating causal effects of treatments in randomized and nonrandomized studies." *Journal of Educational Psychology* 66:688–701.
- Rubin, Donald B. 1980. "Randomization Analysis of Experimental Data: The Fisher Randomization Test Comment." *Journal of the American Statistical Association* 75:591–593.
- Semenova, Vira and Victor Chernozhukov. 2021. "Debiased machine learning of conditional average treatment effects and other causal functions." *The Econometrics Journal* 24:264–289.
- Storer, Adam, Daniel Schneider, and Kristen Harknett. 2020. "What Explains Racial/Ethnic Inequality in Job Quality in the Service Sector?" *American Sociological Review* 85:537–572.
- Sudharsanan, Nikkil and Maarten J Bijlsma. 2021. "Educational note: causal decomposition of population health differences using Monte Carlo integration and the g-formula." *International Journal of Epidemiology* p. dyab090.
- Tsiatis, Anastasios A. 2006. *Semiparametric Theory and Missing Data*. Springer Series in Statistics. New York, NY: Springer.
- van der Laan, Mark J. 2006. "Statistical Inference for Variable Importance." *The International Journal of Biostatistics* 2.
- Van der Laan, Mark J and Sherri Rose. 2011. *Targeted learning: causal inference for observational and experimental data*. Springer Series in Statistics. New York: Springer.
- Van der Vaart, Aad W. 2000. *Asymptotic Statistics*. Cambridge Series on Statistical and Probabilistic Mathematics. Cambridge University Press.
- VanderWeele, Tyler. 2015. *Explanation in causal inference: methods for mediation and interaction*. Oxford University Press.
- VanderWeele, Tyler J. 2014. "A Unification of Mediation and Interaction: A 4-Way Decomposition." *Epidemiology* 25:749–761.
- VanderWeele, Tyler J. and Whitney R. Robinson. 2014. "On the Causal Interpretation of Race in Regressions Adjusting for Confounding and Mediating Variables:." *Epidemiology* 25:473–484.

- VanderWeele, Tyler J. and Eric J. Tchetgen Tchetgen. 2017. “Mediation analysis with time varying exposures and mediators.” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 79:917–938.
- VanderWeele, Tyler J., Stijn Vansteelandt, and James M. Robins. 2014. “Effect Decomposition in the Presence of an Exposure-Induced Mediator-Outcome Confounder:.” *Epidemiology* 25:300–306.
- Xie, Yu, Jennie E. Brand, and Ben Jann. 2012. “Estimating Heterogeneous Treatment Effects with Observational Data.” *Sociological Methodology* 42:314–347.
- Yamaguchi, Kazuo. 2015. “Decomposition of Gender or Racial Inequality with Endogenous Intervening Covariates: An Extension of the DiNardo-Fortin-Lemieux Method.” *Sociological Methodology* 45:388–428.
- Yu, Ang. 2023. *cdgd: Causal Decomposition of Group Disparities*. R package version 0.3.0.
- Zajacova, Anna, Hanna Grol-Prokopczyk, and Zachary Zimmer. 2021. “Pain Trends Among American Adults, 2002–2018: Patterns, Disparities, and Correlates.” *Demography* 58:711–738.
- Zhou, Xiang. 2019. “Equalization or Selection? Reassessing the “Meritocratic Power” of a College Degree in Intergenerational Income Mobility.” *American Sociological Review* 84:459–485.
- Zhou, Xiang. 2022. “Attendance, Completion, and Heterogeneous Returns to College: A Causal Mediation Approach.” *Sociological Methods & Research* p. 004912412211138.
- Zhou, Xiang and Yu Xie. 2020. “Heterogeneous Treatment Effects in the Presence of Self-Selection: A Propensity Score Perspective.” *Sociological Methodology* 50:350–385.
- Ziol-Guest, Kathleen M. and Kenneth T. H. Lee. 2016. “Parent Income–Based Gaps in Schooling: Cross-Cohort Trends in the NLSYs and the PSID.” *AERA Open* 2:233285841664583.