

NONPARAMETRIC CAUSAL DECOMPOSITION OF GROUP DISPARITIES

BY ANG YU^{1,a}  AND FELIX ELWERT^{1,b} 

¹Department of Sociology, University of Wisconsin-Madison, ^aayu33@wisc.edu, ^belwert@wisc.edu

We propose a nonparametric causal framework that decomposes the contribution of a treatment variable to an outcome disparity between two groups. We decompose this contribution into group differences in 1) treatment prevalence, 2) average treatment effects, and 3) selection into treatment. This framework is particularly novel in revealing differential selection as a new mechanism for explaining and ameliorating disparities. It reformulates the classic Kitagawa-Blinder-Oaxaca decomposition in causal terms, supplements causal mediation analysis by explaining group disparities instead of group effects, and resolves a problem of recent random equalization decompositions with validly capturing differential treatment prevalence. We develop nonparametric estimators based on efficient influence functions that are \sqrt{n} -consistent, asymptotically normal, semiparametrically efficient, and multiply robust. We apply our framework to study the causal role of education in intergenerational income persistence. We find that both differential prevalence of and differential selection into college graduation significantly contribute to the disparity in income attainment between income origin groups.

1. Introduction. Social and health scientists often seek to decompose an outcome disparity between groups in terms of the contributions of an intermediate treatment variable. For example, how much, and in what ways, do racial differences in medical care contribute to racial disparities in health (Howe et al., 2014)? Does family formation contribute to the gender wage gap (Cha, Weeden and Schnabel, 2023)? And what are the roles of educational attainment in the relationship between socioeconomic origin and attainment (Ishida, Muller and Ridge, 1995)? The common structure of these questions is that they seek to quantify the mechanisms by which a treatment variable *causally* explains a *descriptive* group disparity.

Prior research has addressed questions of explanation and contribution using three approaches, none of which is fully appropriate for the task. First, the Kitagawa-Blinder-Oaxaca (KBO) decomposition (Kitagawa, 1955; Blinder, 1973; Oaxaca, 1973) is a classic tool that remains popular to this day. However, the KBO decomposition is defined in terms of regression coefficients and does not answer any causal question by design (Fortin, Lemieux and Firpo, 2011, p.13; Lundberg, Johnson and Stewart, 2021, p.542). Second, causal mediation analysis is also used to attribute the relationship between two variables to an intermediate variable (VanderWeele, 2015). Although causal mediation analysis is formulated as causal estimands, it decomposes effects of group variables rather than group disparities. Third, a recently developed approach, which we call random equalization decomposition (VanderWeele and Robinson, 2014; Jackson and VanderWeele, 2018; Lundberg, 2022), marks an important advancement towards an appropriate framework. However, as we will show, random equalization decomposition does not satisfy a null criterion for validly capturing the contribution of differential treatment prevalence to group disparities. Importantly, all three approaches neglect that outcome disparities can in part be explained by differential patterns of selection into treatment across groups.

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. Compared

Keywords and phrases: social inequality, mediation analysis, double machine learning, social mobility.

with the KBO decomposition, our decomposition enables causal attribution and interventional interpretation. It supplements causal mediation analysis by explaining group disparities instead of effects of group variables. It also improves on random equalization decomposition by providing a valid component of differential treatment prevalence.

Our framework reveals three distinct mechanisms through which an intermediate treatment variable can contribute to a group disparity in an outcome. First, groups may have differential prevalence of treatment. Second, the groupwise average treatment effects (ATEs) may differ across groups. Third, the patterns of selection into treatment based on treatment effects may also vary by group membership. All previous approaches are only limited to the first two mechanisms, which leads to the belief that they are the only mechanisms possible (Ward et al., 2019; Diderichsen, Hallqvist and Whitehead, 2019). Hence, our approach uniquely reveals differential selection into treatment as a source of group disparities and a novel policy lever for reducing disparities. A large number of research in the social sciences and public health that narrowly focused on differential prevalence and differential effects could be enhanced by this more comprehensive framework in explaining group disparities.

For our framework to inform policy interventions more broadly, we define two decompositions. First, the unconditional decomposition corresponds to marginal interventions on the treatment. For example, it can be used to answer how much of the racial disparities in income could be reduced if incarceration rates were equalized across racial groups. Second, the conditional decomposition corresponds to interventions that are conducted only within levels of certain pre-treatment covariates. For example, it may be desirable to equalize the receipt of a medical treatment conditional on existing comorbidities, and the conditional decomposition can be used to quantify the impact of such intervention.

For both unconditional and conditional decompositions, we develop nonparametric estimators for the decompositions under conditional ignorability of the treatment and using efficient influence functions (EIF). These estimators can be implemented via data-adaptive methods such as machine learning (ML) and accommodate high-dimensional confounders. We derive the conditions under which the estimators are \sqrt{n} -consistent, asymptotically normal, and semiparametrically efficient. The estimators are also doubly or quadruply robust to misspecification. The estimators are implemented in the R package *cdgd* (Yu, 2023), available from CRAN.

As an empirical application, we study the causal role of college graduation in intergenerational income mobility, which is defined as the disparity in income attainment between parental income groups. This study unifies and contributes to multiple sociology and economics literatures. Policy-wise, it provides insights on the ways in which policies on college education may alter the level of intergenerational persistence.

Our paper proceeds as follows. Section 2 introduces our decompositions and their interventional interpretations. We explicate the contributions of our framework by formally relating it to the KBO decomposition, causal mediation analysis, and the random equalization decomposition. In section 3, we introduce the estimators and their asymptotic theory. Section 4 presents the application. Section 5 concludes with future extensions.

2. Estimands. All proofs for Section 2 are collected in Appendix A.

2.1. Unconditional decomposition. We consider a binary treatment variable $D_i \in \{0, 1\}$ for each individual i . Let Y_i^0 and Y_i^1 be the potential outcomes (Rubin, 1974) of Y_i under the hypothetical intervention to set $D_i = 0$ and $D_i = 1$, respectively. Let $\tau_i := Y_i^1 - Y_i^0$ denote the individual-level treatment effect. For expositional convenience, we assume that higher values of Y_i are better in some sense. Suppose that the population contains two disjoint groups, $G_i = g \in \{a, b\}$, where a denotes the advantaged group and b denotes the disadvantaged

group. We use subscript g to indicate group-specific quantities, for example, $E_g(Y_i) := E(Y_i | g)$. Henceforth, we suppress subscript i to ease notation. In our empirical application, G is parental income groups, Y is adult income, and D is college graduation.

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

$$\text{ASSUMPTION 1 (SUTVA). } Y = DY^1 + (1 - D)Y^0.$$

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

$$\begin{aligned} (1) \quad & E_a(Y) - E_b(Y) \\ &= \underbrace{E_a(Y^0) - E_b(Y^0)}_{\text{baseline}} + \underbrace{E_b(\tau)[E_a(D) - E_b(D)]}_{\text{prevalence}} + \underbrace{E_a(D)[E_a(\tau) - E_b(\tau)]}_{\text{effect}} + \underbrace{\text{Cov}_a(D, \tau) - \text{Cov}_b(D, \tau)}_{\text{selection}}. \end{aligned}$$

First, the “baseline” component reflects the difference in mean baseline potential outcomes Y^0 between groups, i.e., the outcome disparity in the complete absence of treatment.¹ In our application, the baseline component is the part of the income disparity in adulthood that is not attributable to college graduation in any way. Second the “prevalence” component indicates how much of the group disparity is due to differential prevalence of treatment. For example, it indicates the extent to which the difference in college graduation rates across parental income groups contributes to the outcome disparity. Third, the “effect” component reflects the difference in average treatment effects (ATE) across groups. Thus, it reveals the contribution of group-differential average effects of college graduation to the adult income disparity.

Fourth, the “selection” component captures the contribution of differential selection into treatment based on the individual-level treatment effects. Selection into treatment within each group is captured by $\text{Cov}_g(D, \tau)$. This covariance is positive if group members who would benefit more from treatment are more likely to receive treatment. In our example, differential selection will increase the income disparity in adulthood if selection into college graduation is more positive in the higher parental income group than in the lower parental income group.

Both the effect component and the selection component account for the contribution of effect heterogeneity to group disparities. Whereas the effect component captures the contribution of *between*-group effect heterogeneity, the selection component captures the contribution of *within*-group effect heterogeneity. To our knowledge, no prior decomposition has captured the contribution of selection into treatment.

To further explicate our novel selection component, we provide two interpretations for the covariance between treatment and the treatment effect, $\text{Cov}(D, \tau)$. When the receipt of treatment is mainly based on self-selection (e.g., college graduation), the covariance may indicate the extent to which the choice to take up the treatment is rational with respect to returns from the treatment. This interpretation has been the focus of research in economics and sociology (Heckman and Vytlačil, 2005; Brand and Xie, 2010; Heckman, Humphries and Veramendi, 2018). On the other hand, if treatment assignment is mainly administered by external decision-makers (e.g., drug prescription), the covariance indicates how effectively treatments are assigned to individuals. In this case, the selection component could also be called an effectiveness component.

¹The baseline component is identical to the “counterfactual disparity measure” proposed by Naimi et al. (2016). In contrast to the two-way decomposition of Naimi et al. (2016), our four-way decomposition distinguishes more mechanisms.

2.2. *Interventional interpretation.* Our decomposition is formulated in counterfactual terms, and hence is prescriptive for future interventions. The decomposition reveals three policy levers for affecting group disparities. First, policy makers could manipulate the prevalence of the treatment in each of the two groups. Second, they could influence within-group selection into treatment based on group members' treatment effects. Third, they might even be able to manipulate the average treatment effect of each group.²

Different interventions will have different impacts on the prevalence, selection and effect components. Here, we explicate a two-step intervention that maps exactly onto the prevalence and selection components of our decomposition. We express this two-step intervention using randomized intervention notation (Didelez, Dawid and Geneletti, 2006; Geneletti, 2007), where $R(D | g)$ represents a randomly drawn value of treatment D from group g . Then, $E_g(Y^{R(D|g')})$ denotes the post-intervention mean potential outcome for group g after each member of group g has received a random draw of treatment D from group g' . When $g = g'$, the intervention amounts to a random redistribution of treatments within the group. Using its definition, we can rewrite the post-intervention mean potential outcome:

$$(2) \quad E_g(Y^{R(D|g')}) = E_g(Y^0) + E_{g'}(D) E_g(\tau).$$

It follows that our decomposition components can be re-written as follows:

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

This represents the two-step intervention. The first step internally randomizes the treatment in both groups without changing its prevalence. In this step, the pre-intervention disparity is $E_a(Y) - E_b(Y)$, and the post-intervention disparity is $E_a(Y^{R(D|a)}) - E_b(Y^{R(D|b)})$. Since randomizing treatment within each group sets $\text{Cov}_a(D, \tau) = \text{Cov}_b(D, \tau) = 0$ and removes differential selection between groups, the change of disparity resulting from this intervention equals our selection component.

The second step equalizes treatment prevalence without changing selection by giving members of group b random draws of treatment from group a . In this step, the pre-intervention disparity is $E_a(Y^{R(D|a)}) - E_b(Y^{R(D|b)})$, and the post-intervention disparity is $E_a(Y^{R(D|a)}) - E_b(Y^{R(D|a)})$. Therefore, the prevalence component is the change in disparity resulting from this equalization intervention.³ At the end of the two-step intervention, the remaining disparity is the sum of the baseline and the effect components. Using the fact that $\text{Cov}_g(D, \tau) = E_g(Y) - E_g(Y^{R(D|g)})$, we also present a visualization of our decomposition in Appendix B.

²It is more conventional to discuss interventions on treatments than interventions on effects, which cannot be represented using the potential outcome notation. However, the notion of interventions on effects or structural relations have appeared in multiple literatures (Malinsky, 2018; Diderichsen, Hallqvist and Whitehead, 2019; Brady, Finnigan and Hübgen, 2017; Canen and Song, 2023).

³This equalization intervention justifies the scaling factors on the prevalence and effect components of equation 1. Intuitively, the prevalence component is scaled by $E_b(\tau)$, because randomly changing treatment prevalence in group b affects the outcome disparity only to the extent that treatment has an effect in group b on average. The scaling factor on the effect component follows algebraically to complete the decomposition. Different interventions would lead to different scaling factors. For example, intervening to give group a the treatment prevalence of group b , the prevalence component would be scaled by $E_b(\tau)$.

Isolating the distinct contributions of the selection and the prevalence components is useful in practice and enables policy makers to implement only the first step of the two-step intervention. Such choice is appealing under at least two scenarios. First, it is possible that the selection component is positive but the prevalence component is negative. Second, intervening on prevalence may be normatively undesired or practically impossible due to budget constraints.

2.3. Comparison with the KBO decomposition. Disparities research in the social and biomedical sciences traditionally employs KBO decompositions. The form of the KBO decomposition most closely resembling our approach decomposes the outcome disparity between groups into four components with respect to treatment D and pre-treatment covariates \mathbf{X} :

$$E_a(Y) - E_b(Y) = \underbrace{\alpha_a - \alpha_b}_{\text{intercept}} + \underbrace{\beta_b[E_a(D) - E_b(D)]}_{\text{endowment}} + \underbrace{E_a(D)[\beta_a - \beta_b]}_{\text{slope}} + \underbrace{\gamma_b^\top[E_a(\mathbf{X}) - E_b(\mathbf{X})] + E_a(\mathbf{X})^\top(\gamma_a - \gamma_b)}_{\text{residual}},$$

where α_g , β_g , and γ_g are the coefficients from group-specific linear regressions:

$$Y = \alpha_g + \beta_g D + \gamma_g^\top \mathbf{X} + \epsilon.$$

This decomposition attains a causal interpretation under (i) the causal assumption of conditional ignorability of the treatment, $Y^d \perp\!\!\!\perp D \mid g, \mathbf{x}, \forall d, g, \mathbf{x}$, and (ii) the parametric assumption that the group-specific linear regressions are correctly specified. If and only if both assumptions are satisfied, the endowment and slope components in the KBO decomposition are equivalent to our prevalence and effect components, respectively; and the sum of the intercept and residual components equals our baseline component.⁴

Our decomposition differs from the KBO decomposition in three respects. First, our decomposition is inherently causal, because it is directly formulated as estimands in potential outcomes notation, whereas the KBO decomposition as an estimation framework requires additional assumptions to support a causal interpretation. Such assumptions are rarely stated in practice. Second, our decomposition is nonparametric, whereas KBO decompositions are model-based and hence rely on a particular functional form. Third, KBO decompositions do not contain a selection component, because the assumed functional form imposes effect homogeneity within each group. By contrast, our nonparametric decomposition does not impose effect homogeneity at any level and hence contains a selection component as a distinctive conceptual contribution.⁵

2.4. Comparison with causal mediation analysis. Causal mediation analysis (CMA) decomposes the total causal effect of an exposure into direct and indirect effects in terms of a mediating variable. Assuming familiarity with CMA concepts, we derive a three-way mediation decomposition that facilitates comparison with our approach. This mediation decomposition is an extension of prior mediation decompositions based on so-called randomized intervention analogues (RIA) (VanderWeele, 2015, p.619-21, VanderWeele and Tchetgen

⁴In practice, many KBO decompositions are farther removed from our approach. In particular, research often does not separate D from \mathbf{X} or heed the temporal order of variables.

⁵Prior work has offered alternative causal interpretations for KBO decompositions under various assumptions. For example, a prominent literature shows that KBO decompositions can estimate the treatment effect on the treated (ATT) (Fortin, Lemieux and Firpo, 2011; Kline, 2011; Yamaguchi, 2015). Similarly, Chernozhukov, Fernández-Val and Luo (2018) show that a KBO decomposition can estimate the partial treatment effect. And Huber (2015) discusses using a KBO decomposition for estimating natural indirect effects (Pearl, 2001). None of these interpretations accommodates a descriptive group variable.

Tchetgen, 2017; VanderWeele et al., 2014, VanderWeele, 2014). This RIA-based three-way mediation decomposition is:

$$\begin{aligned}
 (3) \quad & \underbrace{E(Y^{a,R(D^a)}) - E(Y^{b,R(D^b)})}_{\text{RTE}} \\
 &= \underbrace{E(Y^{a,0}) - E(Y^{b,0})}_{\text{CDE}} + \underbrace{E(Y^{b,R(D^a)}) - E(Y^{b,R(D^b)})}_{\text{RPIE}} + \underbrace{E(D^a) E(Y^{a,1} - Y^{a,0} - Y^{b,1} + Y^{b,0})}_{\text{RPAI}},
 \end{aligned}$$

where $Y^{g,R(D^{g'})}$ is the potential outcome of Y under the dual intervention of assigning both group g and a random value of D from the population when the population is assigned group g' ; D^g is the potential outcome of D when assigned group g ; Y^g is the potential outcome of Y when assigned group g , and $Y^{g,d}$ is the potential outcome of Y when assigned both group g and treatment d .

We first establish the connections between our decomposition and the RIA-based mediation decomposition by making two unconditional ignorability assumptions for G , i.e., $Y^{g,d} \perp\!\!\!\perp G, \forall d, g$, and $D^g \perp\!\!\!\perp G, \forall g$, and two SUTVA-type assumptions, $E_g(Y^{g,d}) = E_g(Y^d)$ and $E_g(D^g) = E_g(D), \forall d, g$. Then, the conditional direct effect (CDE) in the mediation decomposition equals our baseline component; the RIA of the pure indirect effect (RPIE) equals our prevalence component; and the randomized analogue of the portion attributable to interaction (RPAI) equals our effect component. These equivalences are intuitive: both the CDE and the baseline component capture a group-based outcome difference when the intermediate variable, D , is held at 0; both the RPIE and the prevalence component address the role of the prevalence of D in the relationship between the group and the outcome; finally, the RPAI and the effect component both reflect how the effect of D interacts with group membership.⁶

However, CMA also differs from our decomposition in three crucial ways. First, the identification assumptions of CMA are much stronger. Whereas, as we show below, our approach requires only ignorability of D , CMA requires ignorability of D and G (VanderWeele, 2015).

Second, CMA decomposes an entirely different quantity: a total effect of group membership, rather than the descriptive group disparity. In the social and health sciences, descriptive disparities between groups are often the object of interest in their own right and the focus of popular and policy concern, regardless of the causal impact of group membership (Jackson and VanderWeele, 2018; Lundberg, 2022). For example, income disparities in adulthood between the children of rich and poor parents are often viewed as concerning, regardless of whether these disparities originate from the causal effect of parental income or from confounding factors such as parents' education or race. Furthermore, some group variables may be immutable attributes such as race and gender, on which it is hard to define a hypothetical intervention (Rubin, 1974; Holland, 1986) so that causal mediation estimands may not even be well-defined (Jackson and VanderWeele, 2018).

Third, there is no selection component in the mediation decomposition (3), nor elsewhere in the CMA literature. In the case of RIA-based mediation decompositions, the selection component is ruled out by construction, as random assignment of D is a built-in feature of the RIA of the total effect (RTE). Interestingly, however, we can show that the difference between the total effect of G , $E(Y^a - Y^b)$, and the RTE can be written as

$$(4) \quad \text{Cov}(D^a, Y^{a,1} - Y^{a,0}) - \text{Cov}(D^b, Y^{b,1} - Y^{b,0}).$$

⁶These intuitions also carry over to VanderWeele's (2014) three-way decomposition that is not based on RIAs. However, establishing a similar connection between VanderWeele's three-way mediation decomposition and our decomposition would require a cross-world independence assumption that is not needed in our approach.

Under maintained assumptions, equation (4) equals our selection component. To our knowledge, this difference has not previously appeared as an estimand of interest in the CMA literature.

2.5. Comparison with unconditional random equalization decomposition. In recent years, there have emerged two variants of random equalization decomposition, each of which further has an unconditional version and a conditional version (VanderWeele and Robinson, 2014; Jackson and VanderWeele, 2018; Sudharsanan and Bijlsma, 2021; Lundberg, 2022; Park et al., 2023). In this section, we compare our unconditional decomposition with unconditional random equalization decomposition (URED). URED decomposes the observed disparity into two components. One is the change 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 the remaining disparity.

Similar to our approach, URED is associated with an intervention on treatment, not on group membership. However, URED is a two-way decomposition that contains less information than our four-way decomposition. Moreover, we argue that it is useful to have a decomposition capable of capturing the contribution of differential treatment prevalence. Although the change in disparity may appear to represent this contribution, we show that it does not satisfy a natural null criterion. The null criterion states that a valid measure of the contribution of differential prevalence should be zero when there is no group difference in treatment prevalence, i.e., $E_a(D) = E_b(D)$.

The first variant of URED 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:

$$E_a(Y) - E_b(Y) = \underbrace{E_b\left(Y^{R(D|a)}\right) - E_b(Y)}_{\text{change in disparity}} + \underbrace{E_a(Y) - E_b\left(Y^{R(D|a)}\right)}_{\text{remaining disparity}}.$$

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 change 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 a violation of the null criterion, because when $E_a(D) = E_b(D)$, the change in disparity is clearly still nonzero, as long as $\text{Cov}_b(D, \tau) \neq 0$.

The second variant of URED (Lundberg, 2022) takes a somewhat different form, 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 change in disparity of this variant is $E_a(Y) - E_b(Y) - [E_a(Y^{R(\bar{D})}) - E_b(Y^{R(\bar{D})})]$, where $Y^{R(\bar{D})}$ is the potential outcome under the assignment of a D value randomly drawn from the pooled population. We can rewrite the change in disparity as

$$(5) \quad E(\tau)[E_a(D) - E_b(D)] - [p_a - p_b][E_a(D) - E_b(D)][E_a(\tau) - E_b(\tau)] + \text{Cov}_a(D, \tau) - \text{Cov}_b(D, \tau),$$

where $E(\tau)$ is the overall ATE, and $p_g = \Pr(G = g)$. Therefore, the change in disparity of this variant mixes the prevalence component with the selection component, too. Intuitively, in this random equalization intervention, both group a and group b would get random draws of D , thus selection into treatment is eliminated in both groups. Consequently, the change in disparity again does not satisfy the null criterion, as long as $\text{Cov}_a(D, \tau) \neq \text{Cov}_b(D, \tau)$.

To conclude, the change in disparity in both variants of URED reflects not only differential prevalence but also the selection component of group disparities. This is because URED's underlying interventions both equalize and randomize the treatment. In contrast, as explained in Section 2.2, 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. Consequently, by the null criterion, only our prevalence component is a valid measure of the contribution of unequal treatment prevalence to group disparities. If obtaining such measure is the goal, our prevalence component should be preferred.

2.6. Conditional decomposition. In this subsection, we extend our decomposition to a conditional version whose corresponding interventions are conditional on a vector of pre-treatment covariates \mathbf{Q} . The conditional decomposition is useful as it can shed light on the impact of interventions that are normatively more desirable or realistically more feasible, when the unconditional intervention appears less so (Jackson, 2021). For example, it may be meaningful to consider equalizing college admissions conditional on test scores or equalizing medical treatments conditional on comorbidities. In order to introduce the conditional decomposition, we need a population common support assumption on \mathbf{Q} .

ASSUMPTION 2 (Common support). $\text{supp}_a(\mathbf{Q}) = \text{supp}_b(\mathbf{Q})$.

Under Assumptions 1 and 2, we have

$$\begin{aligned}
 (6) \quad E_a(Y) - E_b(Y) &= \underbrace{E_a(Y^0) - E_b(Y^0)}_{\text{baseline}} + \underbrace{\int [E_a(D | \mathbf{q}) - E_b(D | \mathbf{q})] E_b(\tau | \mathbf{q}) f_b(\mathbf{q}) d\mathbf{q}}_{\text{conditional prevalence}} \\
 &\quad + \underbrace{\int [E_a(\tau | \mathbf{q}) - E_b(\tau | \mathbf{q})] E_a(D | \mathbf{q}) f_a(\mathbf{q}) d\mathbf{q}}_{\text{conditional effect}} + \underbrace{E_a[\text{Cov}_a(D, \tau | \mathbf{Q})] - E_b[\text{Cov}_b(D, \tau | \mathbf{Q})]}_{\text{conditional selection}} \\
 &\quad + \underbrace{\int E_a(D | \mathbf{q}) E_b(\tau | \mathbf{q}) [f_a(\mathbf{q}) - f_b(\mathbf{q})] d\mathbf{q}}_{\mathbf{Q}\text{-distribution}}.
 \end{aligned}$$

The conditional prevalence, effect, and selection components are analogous to their unconditional counterparts but defined within levels of \mathbf{Q} . The sum of the baseline component and the new \mathbf{Q} -distribution component represents the part of the outcome disparity that is not attributable to the conditional contribution of the treatment.⁷

Similar to the unconditional case, we can also interventionally interpret the conditional decomposition using the randomized intervention notation. Let $E_g(Y^{R(D|g', \mathbf{Q})})$ 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.

$$(7) \quad E_g(Y^{R(D|g', \mathbf{Q})}) = E_g(Y^0) + \int E_g(\tau | \mathbf{q}) E_{g'}(D | \mathbf{q}) f_g(\mathbf{q}) d\mathbf{q}.$$

⁷Define the sum of the prevalence, effect, and selection components in equation (1) as the unconditional contribution of treatment to the outcome disparity, and the sum of the corresponding conditional terms given \mathbf{Q} in equation (6) as the conditional contribution of treatment. Then, the \mathbf{Q} -distribution component equals the difference between the unconditional and the conditional contributions of treatment. Clearly, if $(D, \tau) \perp\!\!\!\perp \mathbf{Q}$, the unconditional and conditional decompositions coincide, and the \mathbf{Q} -distribution component is zero.

It then follows that components of the conditional decomposition can be written in the randomized intervention notation:

$$\begin{aligned}
 E_a(Y) - E_b(Y) - \left[E_a \left(Y^{R(D|a, \mathbf{Q})} \right) - E_b \left(Y^{R(D|b, \mathbf{Q})} \right) \right] &= \text{conditional selection} \\
 E_b \left(Y^{R(D|a, \mathbf{Q})} \right) - E_b \left(Y^{R(D|b, \mathbf{Q})} \right) &= \text{conditional prevalence} \\
 E_a \left(Y^{R(D|a, \mathbf{Q})} \right) - E_b \left(Y^{R(D|a, \mathbf{Q})} \right) &= \text{baseline} + \text{conditional effect} \\
 (8) \qquad \qquad \qquad &+ \mathbf{Q}\text{-distribution.}
 \end{aligned}$$

Therefore, a two-step intervention again underlies our conditional decomposition. The first step is a randomization of treatment within groups within \mathbf{Q} levels, where the pre-treatment disparity is $E_a(Y) - E_b(Y)$ and post-intervention disparity is $E_a \left(Y^{R(D|a, \mathbf{Q})} \right) - E_b \left(Y^{R(D|b, \mathbf{Q})} \right)$. The resulting change in disparity is the conditional selection component. The second step is a conditional equalization within \mathbf{Q} levels given the prior conditional randomization. In this step, the pre-intervention disparity is $E_a \left(Y^{R(D|a, \mathbf{Q})} \right) - E_b \left(Y^{R(D|b, \mathbf{Q})} \right)$, the post-intervention disparity is $E_a \left(Y^{R(D|a, \mathbf{Q})} \right) - E_b \left(Y^{R(D|a, \mathbf{Q})} \right)$, and the change in disparity is the conditional prevalence component. What would remain after this two-step intervention is the sum of the baseline, conditional effect, and the \mathbf{Q} -distribution components.⁸ In our application, we choose \mathbf{Q} to be academic achievement in high school. Thus, the conditional decomposition informs interventions on college graduation that are conducted only within levels of prior achievement, preserving the positive relationship between prior achievement and degree attainment.

Jackson (2021) and Lundberg (2022) proposed two variants of conditional random equalization decomposition (CRED), which again decomposes group disparities into a change in disparity component and a remaining disparity component. The hypothetical intervention underlying CRED is a random reassignment of treatment values for people with the same \mathbf{Q} values such that treatment prevalence would be equalized within \mathbf{Q} levels. Analogous to the unconditional case, we show that the change in disparity component in both variants of CRED does not satisfy a conditional null criterion. The conditional null criterion states that a valid measure of the contribution of differential prevalence within levels of \mathbf{Q} should be zero when there is no differential prevalence within all levels of \mathbf{Q} , i.e., $E_a(D | \mathbf{q}) = E_b(D | \mathbf{q}), \forall \mathbf{q}$.

In Jackson's (2021) variant of CRED, the change in disparity is $E_b \left(Y^{R(D|a, \mathbf{Q})} \right) - E_b(Y)$, which we can rewrite as

$$(9) \qquad \int [E_a(D | \mathbf{q}) - E_b(D | \mathbf{q})] E_b(\tau | \mathbf{q}) f_b(\mathbf{q}) d\mathbf{q} - E_b[\text{Cov}_b(D, \tau | \mathbf{Q})].$$

Therefore, the change in disparity is the sum of the conditional prevalence component and the group b part of the conditional selection component. Namely, it captures not only 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 . Consequently, the change in disparity does not satisfy the conditional null criterion.

⁸For equation (7) to hold, we only require $\text{supp}_g(\mathbf{Q}) \subseteq \text{supp}_{g'}(\mathbf{Q})$. This implies that the two-step intervention is well-defined 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. Note that $\text{supp}_b(\mathbf{Q}) \subseteq \text{supp}_a(\mathbf{Q})$ is a weaker condition than Assumption 2. However, under the weaker condition, only the pre- and post-intervention disparities are well-defined, not all components in equation (6).

A similar case can be shown for the change in disparity in the CRED of [Lundberg \(2022\)](#), which is $E_a(Y) - E_b(Y) - [E_a(Y^{R(D|Q)}) - E_b(Y^{R(D|Q)})]$, where $Y^{R(D|Q)}$ is the potential outcome of an individual when they were given a random draw of D from those in the pooled population who share the same Q value with them. We can rewrite this change in disparity as

$$(10) \quad E_a[\text{Cov}_a(D, \tau | Q)] - E_b[\text{Cov}_b(D, \tau | Q)] + \int [E_a(D | q) - E_b(D | q)][E_a(\tau | q)f_a(q)p_b(q) + E_b(\tau | q)f_b(q)p_a(q)]dq,$$

where $p_g(q) = \Pr(G = g | q)$. Hence, the change in disparity involves both a term for differential treatment prevalence within Q levels and the conditional selection component. The intervention in this case corresponds to the combination of a conditional equalization and a conditional randomization in both groups. Consequently, both differential treatment prevalence and differential selection into treatment would be reduced to zero conditional on Q . Therefore, the change in disparity again does not satisfy the conditional null criterion. Distinct from the change in disparity in both variants of CRED, the conditional prevalence component in our decomposition solely measures the impact of a conditional equalization intervention. By the conditional null criterion, only our conditional prevalence component can be interpreted as a valid measure of the contribution of differential treatment prevalence within Q levels to group disparities.

3. Identification, Estimation, and Inference. We identify our unconditional and conditional decompositions using the standard assumptions of conditional ignorability and overlap. Without loss of generality, we let $Q \subseteq X$.

ASSUMPTION 3 (Conditional ignorability). $Y^d \perp\!\!\!\perp D | X = x, G = g, \forall d, x, g$.

ASSUMPTION 4 (Overlap). $0 < E(D | X = x, G = g) < 1, \forall x, g$.

We develop nonparametric and efficient estimators for our decompositions. 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 or quadruple robustness properties.

To unburden notation, we define the following functions of the observed data: $\mu(d, X, g) = E(Y | d, X, g)$, $\pi(d, X, g) = \Pr(D = d | X, g)$, and $\omega(d, Q, g) = E[\mu(d, X, g) | Q, g]$. Also recall that $p_g = \Pr(G = g)$, and $p_g(Q) = \Pr(G = g | Q)$. We use circumflexes to denote estimated quantities.

3.1. Unconditional decomposition. All components of the unconditional decomposition can be expressed as simple linear combinations of the total disparity and two generic functions of the potential outcomes evaluated at appropriate values of d , g , and g' : $\xi_{dg} := E(Y^d | g)$ and $\xi_{dgg'} := E(Y^d | g) E(D | g')$. The relationships between the components of the unconditional decomposition and the generic functions are as follows:

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

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

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

$$\text{Selection} = E_a(Y) - E_b(Y) - \xi_{0a} + \xi_{0b} - \xi_{1aa} + \xi_{0aa} + \xi_{1bb} - \xi_{0bb}.$$

Hence, the EIFs and one-step estimators for the decomposition components directly follow from those for ξ_{dg} and $\xi_{dgg'}$.^{9,10} Under Assumptions 1, 3, and 4, ξ_{dg} and $\xi_{dgg'}$ can be identified as the following quantities:

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

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

These identification results then enable the derivation of the EIFs for ξ_{dg} and $\xi_{dgg'}$.

THEOREM 1 (EIF, unconditional decomposition). *Under Assumptions 1, 3, and 4, the EIF of ξ_{dg} is*

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

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

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

*Proof is in Appendix C, where we also derive the general EIF with survey weights.*¹¹

We use the EIFs as estimating equations, i.e., set their sample averages to zero and solve for ξ_{dg} and $\xi_{dgg'}$, respectively. The one-step estimators of ξ_{dg} and $\xi_{dgg'}$ thus are

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

Each estimator contains two nuisance functions, $\pi(d, \mathbf{X}, g)$ and $\mu(d, \mathbf{X}, g)$. The estimators are consistent as long as either one of the two nuisance functions is consistently estimated.

⁹These generic functions could also provide a basis for the estimation of Jackson and VanderWeele's (2018) version of the URED, since its change in disparity component can be represented as $\xi_{0b} + \xi_{1ba} - \xi_{0ba} - E_b(Y)$.

¹⁰Note that although the selection component is defined in terms of individual-level treatment effects in equation 1, its estimation does not require estimation of individualized treatment effects.

¹¹Related EIFs have appeared in prior work. Park and Kang (2023) give the EIF for $\xi_{1g} - \xi_{0g}$. The EIF for $\xi_{0,g} + \xi_{1gg'} - \xi_{0gg'}$ coincides with the EIF for $E_g(Y^{g,R(D^{g'})})$ in Díaz et al. (2021) for binary D when omitting confounders for G . This is because $E_g(Y^{g,R(D^{g'})}) = \xi_{0,g} + \xi_{1gg'} - \xi_{0gg'}$ under the assumptions of Section 2.4. Neither of these prior works accommodate survey weights.

THEOREM 2 (Double robustness in consistency, unconditional decomposition). *Under Assumptions 1, 3, and 4, either consistent estimation of $\mu(d, \mathbf{X}, g)$ or of $\pi(d, \mathbf{X}, g)$ is sufficient for the consistency of $\hat{\xi}_{dg}$ and $\hat{\xi}_{dgg'}$.*¹²

The nuisance functions $\pi(d, \mathbf{X}, g)$ and $\mu(d, \mathbf{X}, g)$ can be estimated using various approaches. To avoid imposing parametric assumptions, we recommend using flexible ML models with cross-fitting (Kennedy, 2022; Chernozhukov et al., 2018). To improve the finite-sample performance of the estimator, we may stabilize the weight, $\mathbb{1}(D = d)/\hat{\pi}(d, \mathbf{X}, g)$, by dividing it by its sample average

To study the asymptotic distribution of the one-step estimators for the unconditional decomposition, we invoke three additional assumptions about the nuisance functions, which are the same as the assumptions required for the double ML estimator of the ATE (Kennedy, 2022; Chernozhukov et al., 2018). The consistent estimation of p_g and $E(D | g)$, $\forall g$, is left implicit. We let $\|\cdot\|$ denote the L_2 -norm.

ASSUMPTION 5a (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 5b (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 5c (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 3 (Asymptotic distributions, unconditional decomposition). *Under Assumptions 1, 3, 4, 5a, 5b, and 5c, the cross-fitted one-step estimators for $\hat{\xi}_{dg}$ and $\hat{\xi}_{dgg'}$ are \sqrt{n} -consistent, asymptotically normal, and semiparametrically efficient, i.e., $\sqrt{n}(\hat{\xi}_{dg} - \xi_{dg}) \xrightarrow{d} \mathcal{N}(0, \sigma_{dg}^2)$, and $\sqrt{n}(\hat{\xi}_{dgg'} - \xi_{dgg'}) \xrightarrow{d} \mathcal{N}(0, \sigma_{dgg'}^2)$, where $\sigma_{dg}^2 := E[\phi_{dg}(Y, D, \mathbf{X}, G)^2]$ and $\sigma_{dgg'}^2 := E[\phi_{dgg'}(Y, D, \mathbf{X}, G)^2]$ are the respective semiparametric efficiency bounds. Proof is in Appendix D.*

We consistently estimate σ_{dg}^2 and $\sigma_{dgg'}^2$ using the averages of the squared estimated EIFs. The asymptotic distributions can then be used to construct hypothesis tests and confidence intervals. Since the unconditional decomposition components are simple additive functions of the observed disparity, $\hat{\xi}_{dg}$, and $\hat{\xi}_{dgg'}$, all properties established for the estimators of $\hat{\xi}_{dg}$, and $\hat{\xi}_{dgg'}$ (double robustness, \sqrt{n} -consistency, asymptotic normality, and semiparametric efficiency) carry over to the final estimators of the decomposition components.

3.2. Conditional decomposition. Relative to the unconditional decomposition, inference for the conditional decomposition requires consideration of one additional generic function:

$$\xi_{dgg'g''} := E \left[E(Y^d | \mathbf{Q}, g) E(D | \mathbf{Q}, g') | g'' \right],$$

¹²The estimator of our unconditional decomposition is doubly robust with respect to the same two nuisance functions as the classic augmented-inverse-probability-of-treatment-weighting (AIPW) estimator of the ATE (Robins, Rotnitzky and Zhao, 1994). The proof for Theorem 2 is omitted, as it is similar to the proof of the double robustness of the AIPW estimator.

where (d, g, g', g'') denotes any one of eight combinations of treatment status and group memberships. The relationships between components of the conditional decomposition and the generic functions are as follows:

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

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

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

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

$$\text{Conditional Selection} = E_a(Y) - E_b(Y) - \xi_{0a} + \xi_{0b} - \xi_{1aaa} + \xi_{0aaa} + \xi_{1bbb} - \xi_{0bbb}.$$

Since estimation of ξ_{dg} was discussed in the previous subsection, we now focus on $\xi_{dgg'g''}$. The EIFs, one-step estimators, and their asymptotic distributions for the components of the conditional decomposition will then follow.¹³ Under Assumptions 1, 2, 3, and 4, we identify $\xi_{dgg'g''}$ as

$$\xi_{dgg'g''} = E \{ \omega(d, \mathbf{Q}, g) E(D | \mathbf{Q}, g') | g'' \}.$$

THEOREM 4 (EIF, conditional decomposition). *Under Assumptions 1, 2, 3, and 4, the EIF of $\xi_{dgg'g''}$ is*

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

Proof is in Appendix C where we also derive the general EIF with survey weights.

We again construct the one-step estimator by using the EIF as an estimating equation.

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

This estimator contains five nuisance functions: $p_g(\mathbf{Q})$, $\pi(d, \mathbf{X}, g)$, $\mu(d, \mathbf{X}, g)$, $E(D | \mathbf{Q}, g)$, and $\omega(d, \mathbf{Q}, g)$. As is the case for the unconditional decomposition, consistent estimation of the conditional decomposition does not require that all nuisance functions be consistently estimated. This leads to its quadruple robustness against misspecification.

THEOREM 5 (Quadruple robustness in consistency, conditional decomposition). *Under Assumptions 1, 2, 3, and 4, $\hat{\xi}_{dgg'g''}$ is consistent if one of four minimal conditions holds, as summarized in Table 1. Proof is in Appendix E.*

¹³We thereby also provide efficient and nonparametric inference for the change in disparity component in the CRED of Jackson (2021), which can be represented as $\xi_{0b} + \xi_{1bab} - \xi_{0bab} - E_b(Y)$.

TABLE 1

Quadruple robustness of $\hat{\xi}_{dgg'g''}$. For each scenario defined by g, g' and g'' , $\hat{\xi}_{dgg'g''}$ is consistent if any of four minimal sets of nuisance functions indicated by check marks is consistently estimated. The first three panels concern the conditions for the consistent estimation of the conditional prevalence ($g = g''$), conditional effect ($g' = g''$), and conditional selection ($g = g' = g''$) components, respectively. Since all conditions require that either $\pi(d, \mathbf{X}, g)$ or $\mu(d, \mathbf{X}, g)$ be consistently estimated, the conditional baseline component is also consistently estimated. The bottom panel shows the four minimal combinations of nuisance functions that must be consistently estimated so that all components of the conditional decomposition are consistently estimated simultaneously.

$\mu(d, \mathbf{X}, g)$	$\pi(d, \mathbf{X}, g)$	$\omega(d, \mathbf{Q}, g)$	$p_g(\mathbf{Q})$	$E(D \mathbf{Q}, g)$
$g = g''$ (conditional prevalence)				
✓		✓	✓	
	✓	✓	✓	
✓				✓
	✓			✓
$g' = g''$ (conditional effect)				
✓		✓		
	✓	✓		
✓			✓	✓
	✓		✓	✓
$g = g' = g''$ (conditional selection)				
✓		✓		
	✓	✓		
✓				✓
	✓			✓
All components				
✓		✓	✓	
	✓	✓	✓	
✓			✓	✓
	✓		✓	✓

As before, we recommend estimating the nuisance functions nonparametrically using ML and cross-fitting. The estimation of $\omega(d, \mathbf{Q}, g)$ deserves particular attention, because it can be doubly robust itself. Specifically, we recommend a pseudo-outcome approach (e.g., [van der Laan, 2006](#); [Semenova and Chernozhukov, 2021](#)), where the pseudo outcome for each d is defined as

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

which is motivated by the fact that $\omega(d, \mathbf{Q}, g) = E[\delta_d(Y, D, \mathbf{X}, G) | \mathbf{Q}, g]$. We first randomly draw two disjoint subsamples from the data. Then we estimate $\pi(d, \mathbf{X}, G)$ and $\mu(d, \mathbf{X}, G)$ in each subsample *without* cross-fitting and obtain the estimated pseudo outcome $\hat{\delta}_d(Y, D, \mathbf{X}, G)$. Finally, we obtain estimates of $\omega(d, \mathbf{Q}, g)$ using cross-fitting, i.e., we fit $E[\hat{\delta}_d(Y, D, \mathbf{X}, G) | \mathbf{Q}, g]$ separately in each subsample and plug in values of \mathbf{Q} from the respective other subsample. Using this procedure, we ensure that the fitting of $\omega(d, \mathbf{Q}, g)$, which relies on estimating the pseudo outcome, is done separately in each subsample. Provided that $E[\hat{\delta}_d(Y, D, \mathbf{X}, G) | \mathbf{Q}, g]$ can be consistently estimated, this approach enables consistent estimation of $\omega(d, \mathbf{Q}, g)$ if either $\mu(d, \mathbf{X}, g)$ or $\pi(d, \mathbf{X}, g)$ is consistently estimated.

To improve finite-sample performance, $\hat{\xi}_{dgg'g''}$ can be stabilized by dividing $\mathbb{1}(D = d)/\hat{\pi}(d, \mathbf{X}, g)$, $\mathbb{1}(G = g')\hat{p}_{g''}(\mathbf{Q})/\hat{p}_{g'}(\mathbf{Q})\hat{p}_{g''}$, and $\mathbb{1}(G = g)\hat{p}_{g''}(\mathbf{Q})/\hat{p}_g(\mathbf{Q})p_{g''}$ by their respective sample averages.

To establish the asymptotic distribution of the one-step estimators for the conditional decomposition, we invoke Assumptions 6, which augments Assumptions 5 with respect to the additional nuisance functions needed for the conditional decomposition.

ASSUMPTION 6a (Boundedness). With probability 1, $\hat{\pi}(d, \mathbf{X}, g) \geq \eta$, $\pi(d, \mathbf{X}, g) \geq \eta$, $\hat{p}_g(\mathbf{Q}) \geq \eta$, $p_g(\mathbf{Q}) \geq \eta$, $|Y - \hat{\mu}(d, \mathbf{X}, g)| \leq \zeta$, $|Y - \mu(d, \mathbf{X}, g)| \leq \zeta$, $|\mu(d, \mathbf{X}, g)| \leq \zeta$, $|\omega(d, \mathbf{Q}, g)| \leq \zeta$, 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 6b (Consistency). $\|\hat{\mu}(d, \mathbf{X}, g) - \mu(d, \mathbf{X}, g)\| = o_p(1)$, $\|\hat{\pi}(d, \mathbf{X}, g) - \pi(d, \mathbf{X}, g)\| = o_p(1)$, $\|\hat{\omega}(d, \mathbf{Q}, g) - \omega(d, \mathbf{Q}, g)\| = o_p(1)$, $\|\hat{\mathbb{E}}(D | \mathbf{Q}, g) - \mathbb{E}(D | \mathbf{Q}, g)\| = o_p(1)$, and $\|\hat{p}_g(\mathbf{Q}) - p_g(\mathbf{Q})\| = o_p(1)$, $\forall d, g$.

ASSUMPTION 6c (Convergence rate). First, we require $\|\hat{\pi}(d, \mathbf{X}, g) - \pi(d, \mathbf{X}, g)\| \|\hat{\mu}(d, \mathbf{X}, g) - \mu(d, \mathbf{X}, g)\| = o_p(n^{-1/2})$, $\forall d, g$. Second, depending on the specific combination of g, g' , and g'' in a $\xi_{dgg'g''}$, we require $\|\hat{\mathbb{E}}(D | \mathbf{Q}, g) - \mathbb{E}(D | \mathbf{Q}, g)\| = o_p(n^{-1/2})$, $\forall g$, when $g = g''$; $\|\hat{\omega}(d, \mathbf{Q}, g) - \omega(d, \mathbf{Q}, g)\| = o_p(n^{-1/2})$, $\forall d, g$, when $g' = g''$, and $\|\hat{\omega}(d, \mathbf{Q}, g) - \omega(d, \mathbf{Q}, g)\| \|\hat{\mathbb{E}}(D | \mathbf{Q}, g) - \mathbb{E}(D | \mathbf{Q}, g)\| = o_p(n^{-1/2})$, $\forall d, g$, when $g = g' = g''$.

THEOREM 6 (Asymptotic distribution, conditional decomposition). *Under Assumptions 1, 2, 3, 4, 6a, 6b, and 6c, the cross-fitted one-step estimator of $\xi_{dgg'g''}$ is \sqrt{n} -consistent, asymptotically normal, and semiparametrically efficient, i.e., $\sqrt{n}(\hat{\xi}_{dgg'g''} - \xi_{dgg'g''}) \rightarrow \mathcal{N}(0, \sigma_{dgg'g''}^2)$, where $\sigma_{dgg'g''}^2 := \mathbb{E}[\phi_{dgg'g''}(Y, D, \mathbf{X}, G)^2]$ is the semiparametric efficiency bound. Proof is in Appendix D.*

Since the conditional decomposition components are additive functions of the observed disparity, $\hat{\xi}_{dg}$, and $\xi_{dgg'g''}$, double robustness, \sqrt{n} -consistency, asymptotic normality, and semiparametric efficiency all carry over to the final estimators of the decomposition components. We conduct hypothesis tests and construct confidence intervals analogously to the unconditional decomposition.

Finally, for the estimation of the conditional decomposition, we note a tension between consistency on one hand, and asymptotic normality and semiparametric efficiency on the other. ML is preferred for the consistent estimation of all decomposition components because it does not impose parametric assumptions. However, ML does not typically satisfy the convergence rate conditions of Assumption 6c for establishing asymptotic normality and semiparametric efficiency for the conditional prevalence, conditional effect, and Q-distribution components. Interestingly, this tension does not exist for the baseline and conditional selection components, because their convergence rate assumptions can reasonably be achieved with ML.

4. Application.

4.1. Overview. Our application decomposes the contribution of college graduation to the perpetuation of income inequality across generations, which is also known as intergenerational income persistence, the complement to income mobility. Groups (G) are defined by parental income; the outcome (Y) is offspring's adult income; and the treatment (D) is college graduation. Although previous research has touched upon all four components of

our unconditional decomposition to some extent, this analysis is the first to present a unified decomposition.

The baseline component of our decomposition represents the part of the disparity in adult income that is unaccounted for by college graduation. Prior research enumerates multiple sources of intergenerational income persistence that may operate independently of college graduation. For example, parental income is associated with a variety of pre-college characteristics that may directly influence income attainment, such as cognitive skills and noncognitive traits in adolescence (Heckman, Stixrud and Urzua, 2006). Moreover, people from more privileged backgrounds likely benefit from their parents’ human, social, and financial capital regardless of their own formal educational attainment. Interventions on college graduation would not eliminate these channels of income persistence.

Speaking to the prevalence component in our decomposition, social scientists have long regarded education as a mediator in the intergenerational reproduction of socioeconomic inequality (Blau and Duncan, 1967, chapter 4 & 5; Featherman and Hauser, 1978, p.255-9; Ishida, Muller and Ridge, 1995). Specifically, research documents large differences in college graduation rates across parental income groups (Ziol-Guest and Lee, 2016; Bailey and Dynarski, 2011), and simulations suggest that rising educational inequality has strengthened intergenerational income persistence over time (Bloome, Dyer and Zhou, 2018).

Speaking to our effect component, there is an active literature on heterogeneity in the effect of college graduation on adult income, although results vary. Brand et al. (2021) and Cheng et al. (2021) found larger effects of college graduation on the incomes of people from disadvantaged backgrounds. By contrast, Zhou (2019), Fiel (2020), and Yu et al. (2021) found no statistically significant heterogeneity in the effect of college completion on income across parental income groups. We note that none of these works evaluated the extent to which income disparities can be attributed to groupwise differential effects of college.

Finally, some prior work has addressed selection into college as a function of college effects on income, i.e. $\text{Cov}(D, \tau)$. Here, too, results are mixed. Brand and Xie’s (2010) and Brand et al.’s (2021) analyses found negative selection, i.e., that those who are least likely to attend college would benefit most from it. By contrast, the instrumental variable analysis of Heckman et al. (2018) found positive selection into college. However, prior work estimates selection in the pooled population rather than within parental income groups, thereby missing the link between differences in group-specific selection and group-based outcome disparities that our decomposition identifies. (Appendix F clarifies and synthesizes related notions of “selection into college” in the social science literature.)

4.2. Data, variables and estimation. We analyze the National Longitudinal Survey of Youth 1979, a nationally representative U.S. cohort study of individuals born between 1957 and 1964. We restrict the sample to respondents who were between 14 to 17 years old at baseline in 1979 to ensure that income origin is measured prior to respondents’ college graduation. We also limit the analysis to respondents who graduated from high school by age 29. The sample size of our complete-case analysis is $N = 2,008$. Missingness mostly occurs in the outcome variable due to loss to follow-up (22%), with less missingness in other variables (<6%). Appendix Table A1 presents associations between outcome missingness and baseline covariates.

We contrast parental income-origin groups (G), defined as the top 40% and bottom 40% of family incomes averaged over the first three waves of the survey (1979, 1980, and 1981, when respondents were 14 to 20 years old) and divided by the square root of the family size to adjust for need (Zhou, 2019). Treatment (D) is a binary indicator of whether the respondent graduated from college by age 29. The outcome is the percentile rank of respondents’ adult income, averaged over five survey waves between age 35 and 44, divided by the square root

of family size. For the conditional decomposition, we define Q as the Armed Forces Qualification Test (AFQT) score, measured in 1980. The AFQT score is a widely used measure of academic achievement that predicts college completion.

We measure an extensive set of confounders (X) at baseline, including 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 score of control locus, the Rosenberg self-esteem score, language spoken at home, Metropolitan Statistical Area category, separation from mother, school satisfaction, region of residence, and mother's working status.

We present four estimates each for our unconditional and conditional decompositions, using different models for the nuisance functions to assess robustness: three alternative ML methods (gradient boosting machine [GBM], neural networks, and random forests) and one set of parametric models. Depending on whether the left-hand-side variable is continuous or binary, the parametric models are linear or logit. Specifically, for $\mu(d, X, g)$, we use all two-way interactions between D and (X, G) , along with their main effects. For $\pi(d, X, g)$ and $p_g(Q)$, the covariates are entered linearly. For $E(D | Q, g)$ and $E[\hat{\delta}_d(Y, D, X, G) | Q, g]$, we include all main effects and two-way interactions between G and Q .

4.3. Results. Figure 1 presents our main results for the components of the unconditional and the conditional decompositions across different models for the nuisance functions (see Appendix Tables A3 and A4 for details). To aid interpretation, Appendix Table A2 additionally reports estimated group-specific means of baseline potential outcomes, treatment proportions, ATEs, and covariances between treatment and treatment effect, as well as group differences in these quantities.

Descriptively, we find that individuals from lower income origin on average achieve 21 percentiles lower incomes in their 30s and 40s than individuals from higher income origin. This confirms the existence of intergenerational income persistence and represents the total disparity that we aim to decompose.

The baseline component constitutes nearly 90% of the total disparity in adult income across all models for the nuisance functions. This demonstrates that most of intergenerational income persistence is due to processes that do not involve college graduation or its effects. By definition, the baseline component is the same in the unconditional and the conditional decompositions.

The prevalence component in the unconditional decomposition is positive, substantively large, and statistically significant across all models for the nuisance functions, accounting for about 15% of the total disparity in adult income. Consequently, an intervention to equalize college graduation rates across income-origin groups without changing selection into graduation would reduce the total disparity in adult income by about 15%. In other words, such intervention would decrease intergenerational income persistence and increase income mobility. Underlying the prevalence component is the striking inequality in college graduation rates by parental income, as 34% of respondents in the higher-income group, but only 9% of the lower-income group, obtained a college degree by age 29 (see Appendix Table A2).

The effect components in the unconditional decomposition is statistically insignificant due to minimal between-group effect heterogeneity. In Appendix A1, we show that, the group-specific ATEs of college graduation on adult income range from 11 to 15 percentiles across models, all of which are statistically significant. However, the group difference in ATEs is always statistically insignificant.

Our estimates for the selection component in the unconditional decomposition are consistently negative. These estimates are statistically significant for the three ML models but not for the parametric models of the nuisance functions. Consequently, randomizing college graduation within each income-origin group to remove selection without changing the prevalence

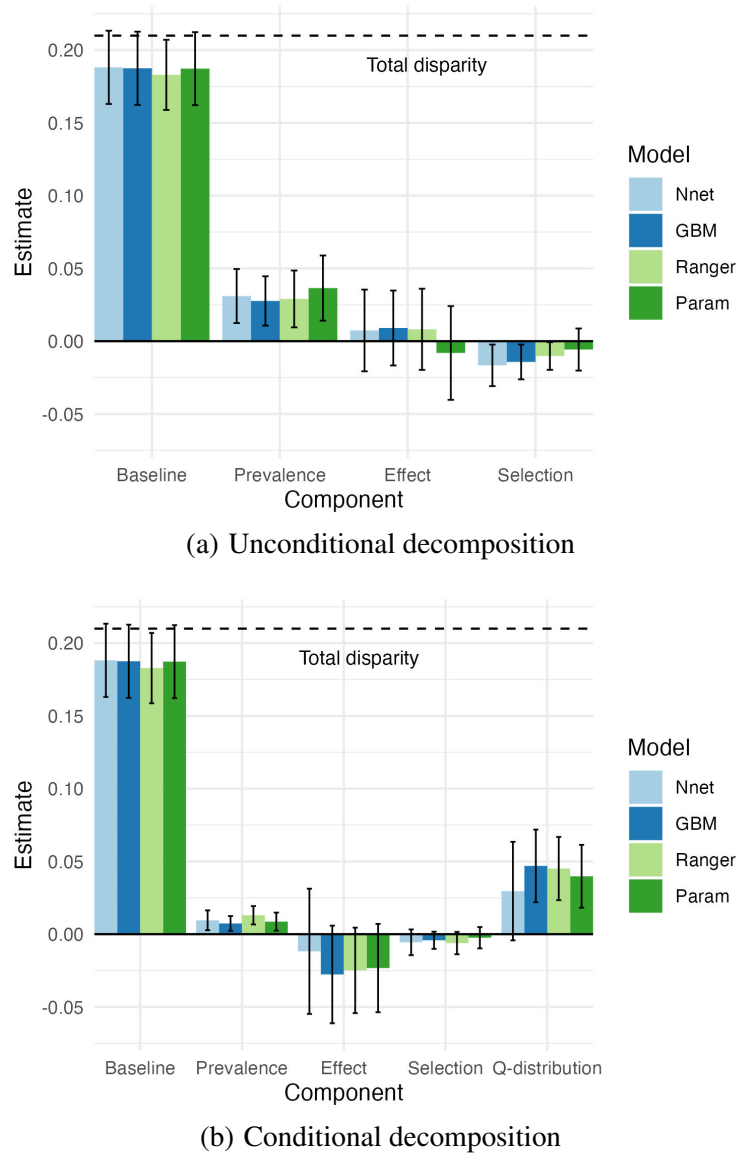


FIG 1. *Decomposition estimates. Cond=conditional, Para=parametric models, Nnet=neural networks, GBM=generalized boosting machine, Ranger=random forests. Error bars indicate 95% confidence intervals.*

of college graduation would increase the total disparity in adult income by around 7%. Thus, this intervention would increase intergenerational income persistence and decrease income mobility. Selection into college graduation decreases intergenerational income persistence because selection is positive in the disadvantaged group and negative in the advantaged group (Appendix Table A2). This lends support to stipulations from earlier research that obtaining a college degree is more of a rational decision in pursuit of economic returns among disadvantaged individuals, and more of an adherence to the social norm of college-attendance among advantaged individuals (Mare, 1980; Hout, 2012). This selection component is the central conceptual contribution of our decomposition approach and also the most novel finding of our empirical application.

In sum, our unconditional decomposition reveals that college graduation in the United States plays two contradictory roles for intergenerational income persistence. On one hand,

higher college graduation rates among high-income origin individuals increase intergenerational income persistence because college graduation increases adult income. On the other hand, part of this prevalence component is offset by our newly identified selection component.

The conditional decomposition quantifies the contributions of college graduation within levels of AFQT achievement scores. Hence, it informs settings in which policy makers cannot, or do not want to, change the factual relationship between prior achievement and college completion, perhaps due to normative constraints or meritocratic preferences.

In the conditional decomposition, the conditional prevalence component is positive and statistically significant, although it is much smaller than in the unconditional decomposition and only accounts for 3% to 6% of the total disparity. Hence, equalizing chances of graduating college within levels of prior achievement would still somewhat decrease intergenerational income persistence. All estimates for the conditional effect and conditional selection components are negative and statistically insignificant.

The Q-distribution component of the conditional composition measures the part of the unconditional contribution of college graduation and its effects to intergenerational income persistence that is associated with AFQT scores (see footnote 8 above). Here, the Q-distribution component is positive and statistically significant and shows the extent to which the unconditional contribution of college graduation exceeds its conditional contribution.

Finally, we also estimate the change in disparity components in the URED and CRED of Jackson and VanderWeele (2018) and Jackson (2021). As explained above, these decompositions estimate the impact of (marginally or conditionally) randomizing treatment and do not distinguish between group differences in prevalence and selection into college graduation. Consequently, both the URED and the CRED underestimate the extent to which differential college graduation rates alone, net of selection, contribute to intergenerational income persistence (Appendix Tables A3 and A4).

5. Discussion. We provide a new decomposition approach for quantifying the ways in which a treatment variable affects outcome disparities between ascriptive groups. Compared with previous approaches, we provide a conceptually more appropriate causal framework. Moreover, this framework reveals differential selection into treatment as a new explanation for group disparities and a novel policy lever for ameliorating disparities. We developed non-parametric estimators that are efficient, asymptotically normal, and doubly or quadruply robust.

Our approach can be extended in multiple directions. First, it is possible to develop analogous decompositions for non-binary treatments. Consider a categorical treatment variable that has J categories. For each individual, exactly one of the binary indicator variables D_1 through D_J takes the value of 1, with the others being 0. Treating D_1 as the baseline category, we have the following decomposition:

$$\begin{aligned} E_a(Y) - E_b(Y) &= E_a(Y^{D_1=1}) - E_b(Y^{D_1=1}) + \sum_{j=2}^J E_b(\tau_j) [E_a(D_j) - E_b(D_j)] \\ &\quad + \sum_{j=2}^J E_a(D_j) [E_a(\tau_j) - E_b(\tau_j)] + \sum_{j=2}^J [\text{Cov}_a(D_j, \tau_j) - \text{Cov}_b(D_j, \tau_j)], \end{aligned}$$

where $Y^{D_j=1}$ is the potential outcome of Y under the assignment of the j th category of the treatment, and $\tau_j = Y^{D_j=1} - Y^{D_1=1}$. For continuous treatments, the randomization and equalization interventions in the two-step intervention framework, 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 binary treatments, $D_1, D_2 \in \{0, 1\}$, 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, in this paper, our estimators assume conditional ignorability. However, identification of the unconditional decomposition is also possible with instrumental variables using the marginal treatment effect framework (Heckman and Vytlačil, 2005), as the unconditional decomposition can be expressed in terms of group-specific outcome means, treatment prevalence, ATTs, and ATEs (see Appendix F). Fourth, when the conditional ignorability assumption is in doubt, a sensitivity analysis would also be helpful. Fifth, as an alternative to the double ML-style (Chernozhukov et al., 2018) 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.

Acknowledgments. We thank Paul Bauer, Eric Grodsky, Aleksei Opacic, Guanghui Pan, Chan Park, Ben Rosche, Jiwei Zhao, and especially Xiang Zhou for helpful suggestions. We are also grateful for the insightful and thought-provoking comments of AOAS reviewers. Earlier versions of this paper has been presented at the Causality in the Social Sciences Workshop at GESIS in 2021, the PAA annual meeting in 2022, the ACIC in 2022, the RC28 Spring meeting in 2023, and the JSM in 2023.

Funding. The authors gratefully acknowledge core grants to the Center for Demography and Ecology (P2CHD047873) and to the Center for Demography of Health and Aging (P30AG017266) at UW-Madison, a Romnes Fellowship to Felix Elwert at UW-Madison, and a Wisconsin Partnership Program grant from the UW-Madison School of Medicine and Public Health to Christie Bartels and Felix Elwert.

SUPPLEMENTARY MATERIAL

Code for the empirical application

R Code for the empirical application is available at https://github.com/ang-yu/causal_decomposition_case_study.

Supplementary appendices

Proofs are contained in Appendices A through D. In Appendix E, we discuss the relationship between various concepts of “selection” in the social science literature. Appendix F presents supplemental tables for the empirical application.

REFERENCES

- BAILEY, M. and DYNARSKI, S. (2011). Gains and Gaps: Changing Inequality in U.S. College Entry and Completion Technical Report No. w17633, National Bureau of Economic Research, Cambridge, MA. <https://doi.org/10.3386/w17633>
- BICKEL, P. J., KLAASSEN, C. A. J., RITOV, Y. and WELLNER, J. A. (1998). *Efficient and Adaptive Estimation for Semiparametric Models*. Springer, New York, NY.
- BLAU, P. M. and DUNCAN, O. D. (1967). *The American Occupational Structure*. Free Press.
- BLINDER, A. S. (1973). Wage Discrimination: Reduced Form and Structural Estimates. *The Journal of Human Resources* **8** 436. <https://doi.org/10.2307/144855>
- BLOOME, D., DYER, S. and ZHOU, X. (2018). Educational Inequality, Educational Expansion, and Intergenerational Income Persistence in the United States. *American Sociological Review* **83** 1215–1253. <https://doi.org/10.1177/0003122418809374>

- BRADY, D., FINNIGAN, R. M. and HÜBGEN, S. (2017). Rethinking the Risks of Poverty: A Framework for Analyzing Prevalences and Penalties. *American Journal of Sociology* **123** 740–786. <https://doi.org/10.1086/693678>
- BRAND, J. E. and XIE, Y. (2010). Who Benefits Most from College?: Evidence for Negative Selection in Heterogeneous Economic Returns to Higher Education. *American Sociological Review* **75** 273–302. <https://doi.org/10.1177/0003122410363567>
- BRAND, J. E., XU, J., KOCH, B. and GERALDO, P. (2021). Uncovering Sociological Effect Heterogeneity Using Tree-Based Machine Learning. *Sociological Methodology* **51** 189–223. <https://doi.org/10.1177/0081175021993503>
- CANEN, N. and SONG, K. (2023). Synthetic Decomposition for Counterfactual Predictions. arXiv:2307.05122 [econ].
- CHA, Y., WEEDEN, K. A. and SCHNABEL, L. (2023). Is the Gender Wage Gap Really a Family Wage Gap in Disguise? *American Sociological Review* **88** 972–1001. <https://doi.org/10.1177/00031224231212464>
- CHENG, S., BRAND, J. E., ZHOU, X., XIE, Y. and HOUT, M. (2021). Heterogeneous returns to college over the life course. *Science Advances* **7** eabg7641. <https://doi.org/10.1126/sciadv.abg7641>
- CHERNOZHUKOV, V., FERNÁNDEZ-VAL, I. and LUO, Y. (2018). The Sorted Effects Method: Discovering Heterogeneous Effects Beyond Their Averages. *Econometrica* **86** 1911–1938. <https://doi.org/10.3982/ECTA14415>
- CHERNOZHUKOV, V., CHETVERIKOV, D., DEMIRER, M., DUFLO, E., HANSEN, C., NEWEY, W. and ROBINS, J. (2018). Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal* **21** C1–C68. <https://doi.org/10.1111/ectj.12097>
- DIDELEZ, V., DAWID, A. P. and GENELETTI, S. (2006). Direct and Indirect Effects of Sequential Treatments. In *Proceedings of the Twenty-Second Conference on Uncertainty in Artificial Intelligence* (R. DECHTER and T. S. RICHARDSON, eds.) 138–146. AUAI Press, Arlington, Virginia.
- DIDERICHSEN, F., HALLQVIST, J. and WHITEHEAD, M. (2019). Differential vulnerability and susceptibility: how to make use of recent development in our understanding of mediation and interaction to tackle health inequalities. *International Journal of Epidemiology* **48** 268–274. <https://doi.org/10.1093/ije/dyy167>
- DÍAZ, I., HEJAZI, N. S., RUDOLPH, K. E. and VAN DER LAAN, M. J. (2021). Nonparametric efficient causal mediation with intermediate confounders. *Biometrika* **108** 627–641. <https://doi.org/10.1093/biomet/asaa085>
- FEATHERMAN, D. L. and HAUSER, R. M. (1978). *Opportunity and Change. Studies in population*. Academic Press.
- FIEL, J. E. (2020). Great Equalizer or Great Selector? Reconsidering Education as a Moderator of Intergenerational Transmissions. *Sociology of Education* 003804072092788. <https://doi.org/10.1177/0038040720927886>
- FORTIN, N., LEMIEUX, T. and FIRPO, S. (2011). Decomposition Methods in Economics. In *Handbook of Labor Economics*, **4** 1–102. Elsevier. [https://doi.org/10.1016/S0169-7218\(11\)00407-2](https://doi.org/10.1016/S0169-7218(11)00407-2)
- GENELETTI, S. (2007). Identifying Direct and Indirect Effects in a Non-Counterfactual Framework. *Journal of the Royal Statistical Society Series B: Statistical Methodology* **69** 199–215. <https://doi.org/10.1111/j.1467-9868.2007.00584.x>
- HECKMAN, J. J., HUMPHRIES, J. E. and VERAMENDI, G. (2018). Returns to Education: The Causal Effects of Education on Earnings, Health, and Smoking. *Journal of Political Economy* **126** 50.
- HECKMAN, J. J., STIXRUD, J. and URZUA, S. (2006). The Effects of Cognitive and Noncognitive Abilities on Labor Market Outcomes and Social Behavior. *Journal of Labor Economics* **24** 411–482. <https://doi.org/10.1086/504455>
- HECKMAN, J. J. and VYTLACIL, E. (2005). Structural Equations, Treatment Effects, and Econometric Policy Evaluation I. *Econometrica* **73** 669–738. <https://doi.org/10.1111/j.1468-0262.2005.00594.x>
- HINES, O., DUKES, O., DIAZ-ORDAZ, K. and VANSTEELANDT, S. (2022). Demystifying Statistical Learning Based on Efficient Influence Functions. *The American Statistician* 1–13. <https://doi.org/10.1080/00031305.2021.2021984>
- HOLLAND, P. W. (1986). Statistics and Causal Inference. *Journal of the American Statistical Association* **81** 945–960.
- HOUT, M. (2012). Social and Economic Returns to College Education in the United States. *Annual Review of Sociology* **38** 379–400. <https://doi.org/10.1146/annurev.soc.012809.102503>
- HOWE, C. J., NAPRAVNIK, S., COLE, S. R., KAUFMAN, J. S., ADIMORA, A. A., ELSTON, B., ERON, J. J. and MUGAVERO, M. J. (2014). African American Race and HIV Virological Suppression: Beyond Disparities in Clinic Attendance. *American Journal of Epidemiology* **179** 1484–1492. <https://doi.org/10.1093/aje/kwu069>
- HUBER, M. (2015). Causal Pitfalls in the Decomposition of Wage Gaps. *Journal of Business & Economic Statistics* **33** 179–191. <https://doi.org/10.1080/07350015.2014.937437>
- ISHIDA, H., MÜLLER, W. and RIDGE, J. M. (1995). Class Origin, Class Destination, and Education: A Cross-National Study of Ten Industrial Nations. *American Journal of Sociology* **101** 145–193. <https://doi.org/10.1086/230701>

- JACKSON, J. W. (2021). Meaningful Causal Decompositions in Health Equity Research: Definition, Identification, and Estimation Through a Weighting Framework. *Epidemiology* **32** 282–290. <https://doi.org/10.1097/EDE.0000000000001319>
- JACKSON, J. W. and VANDERWEELE, T. J. (2018). Decomposition Analysis to Identify Intervention Targets for Reducing Disparities. *Epidemiology* **29** 825–835. <https://doi.org/10.1097/EDE.0000000000000901>
- KENNEDY, E. H. (2022). Semiparametric doubly robust targeted double machine learning: a review. Number: arXiv:2203.06469 arXiv:2203.06469 [stat].
- KITAGAWA, E. M. (1955). Components of a Difference Between Two Rates. *Journal of the American Statistical Association* **50** 1168–1194.
- KLINE, P. (2011). Oaxaca-Blinder as a Reweighting Estimator. *American Economic Review* **101** 532–537. <https://doi.org/10.1257/aer.101.3.532>
- LUNDBERG, I. (2022). The Gap-Closing Estimand: A Causal Approach to Study Interventions That Close Disparities Across Social Categories. *Sociological Methods & Research* 004912412110557. <https://doi.org/10.1177/00491241211055769>
- LUNDBERG, I., JOHNSON, R. and STEWART, B. M. (2021). What Is Your Estimand? Defining the Target Quantity Connects Statistical Evidence to Theory. *American Sociological Review* **86** 532–565.
- MALINSKY, D. (2018). Intervening on structure. *Synthese* **195** 2295–2312. <https://doi.org/10.1007/s11229-017-1341-z>
- MARE, R. D. (1980). Social Background and School Continuation Decisions. *Journal of the American Statistical Association* **75** 295–305.
- NAIMI, A. I., SCHNITZER, M. E., MOODIE, E. E. M. and BODNAR, L. M. (2016). Mediation Analysis for Health Disparities Research. *American Journal of Epidemiology* **184** 315–324. <https://doi.org/10.1093/aje/kwv329>
- OAXACA, R. (1973). Male-Female Wage Differentials in Urban Labor Markets. *International Economic Review* **14** 693. <https://doi.org/10.2307/2525981>
- PARK, C. and KANG, H. (2023). A groupwise approach for inferring heterogeneous treatment effects in causal inference. *Journal of the Royal Statistical Society Series A: Statistics in Society* qnad125. <https://doi.org/10.1093/jrssa/qnad125>
- PARK, S., KANG, S., LEE, C. and MA, S. (2023). Sensitivity analysis for causal decomposition analysis: Assessing robustness toward omitted variable bias. *Journal of Causal Inference* **11** 20220031. <https://doi.org/10.1515/jci-2022-0031>
- PEARL, J. (2001). Direct and Indirect Effects. In *Proceedings of the Seventeenth Conference on Uncertainty in Artificial Intelligence* 411–20. Morgan Kaufmann, San Francisco, CA.
- ROBINS, J. M., ROTNITZKY, A. and ZHAO, L. P. (1994). Estimation of Regression Coefficients When Some Regressors Are Not Always Observed. *Journal of the American statistical Association* **89** 846–866.
- RUBIN, D. B. (1974). Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of Educational Psychology* **66** 688–701. <https://doi.org/10.1037/h0037350>
- RUBIN, D. B. (1980). Randomization Analysis of Experimental Data: The Fisher Randomization Test Comment. *Journal of the American Statistical Association* **75** 591–593.
- SEMEANOVA, V. and CHERNOZHUKOV, V. (2021). Debiased machine learning of conditional average treatment effects and other causal functions. *The Econometrics Journal* **24** 264–289. <https://doi.org/10.1093/ectj/utaa027>
- SUDHARSANAN, N. and BIJLSMA, M. J. (2021). Educational note: causal decomposition of population health differences using Monte Carlo integration and the g-formula. *International Journal of Epidemiology* dyab090. <https://doi.org/10.1093/ije/dyab090>
- VAN DER LAAN, M. J. (2006). Statistical Inference for Variable Importance. *The International Journal of Biostatistics* **2**. <https://doi.org/10.2202/1557-4679.1008>
- VAN DER LAAN, M. J. and ROSE, S. (2011). *Targeted learning: Causal Inference for Observational and Experimental Data*. Springer Series in Statistics. Springer, New York.
- VAN DER VAART, A. W. (2000). *Asymptotic Statistics*. Cambridge Series on Statistical and Probabilistic Mathematics. Cambridge University Press.
- VANDERWEELE, T. J. (2014). A Unification of Mediation and Interaction: A 4-Way Decomposition. *Epidemiology* **25** 749–761. <https://doi.org/10.1097/EDE.0000000000000121>
- VANDERWEELE, T. (2015). *Explanation in Causal Inference: Methods for Mediation and Interaction*. Oxford University Press.
- VANDERWEELE, T. J. and ROBINSON, W. R. (2014). On the Causal Interpretation of Race in Regressions Adjusting for Confounding and Mediating Variables. *Epidemiology* **25** 473–484. <https://doi.org/10.1097/EDE.0000000000000105>
- VANDERWEELE, T. J. and TCHETGEN TCHETGEN, E. J. (2017). Mediation analysis with time varying exposures and mediators. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **79** 917–938. <https://doi.org/10.1111/rssb.12194>

- VANDERWEELE, T. J., VANSTEELENDT, S. and ROBINS, J. M. (2014). Effect Decomposition in the Presence of an Exposure-Induced Mediator-Outcome Confounder. *Epidemiology* **25** 300–306. <https://doi.org/10.1097/EDE.0000000000000034>
- WARD, J. B., GARTNER, D. R., KEYES, K. M., FLISS, M. D., MCCLURE, E. S. and ROBINSON, W. R. (2019). How do we assess a racial disparity in health? Distribution, interaction, and interpretation in epidemiological studies. *Annals of Epidemiology* **29** 1–7. <https://doi.org/10.1016/j.annepidem.2018.09.007>
- YAMAGUCHI, K. (2015). Decomposition of Gender or Racial Inequality with Endogenous Intervening Covariates: An Extension of the DiNardo-Fortin-Lemieux Method. *Sociological Methodology* **45** 388–428. <https://doi.org/10.1177/0081175015583985>
- YU, A. (2023). cdgd: Causal Decomposition of Group Disparities R package version 0.3.1.
- YU, A., PARK, C., KANG, H. and FLETCHER, J. (2021). Leveraging Machine Learning to Estimate Effect Modification preprint, SocArXiv. <https://doi.org/10.31235/osf.io/8ewup>
- ZHOU, X. (2019). Equalization or Selection? Reassessing the “Meritocratic Power” of a College Degree in Intergenerational Income Mobility. *American Sociological Review* **84** 459–485. <https://doi.org/10.1177/0003122419844992>
- ZIOL-GUEST, K. M. and LEE, K. T. H. (2016). Parent Income–Based Gaps in Schooling: Cross-Cohort Trends in the NLSYs and the PSID. *AERA Open* **2** 233285841664583. <https://doi.org/10.1177/2332858416645834>

SUPPLEMENTARY APPENDICES TO NONPARAMETRIC CAUSAL DECOMPOSITION OF GROUP DISPARITIES

BY ANG YU^{1,a}  AND FELIX ELWERT^{1,b} 

¹Department of Sociology, University of Wisconsin-Madison, ^aayu33@wisc.edu, ^belwert@wisc.edu

Appendix A: Proofs for Section 2.

A.1. Equation (2). Note that $R(D | g')$ denotes a randomly drawn value of treatment D from group g' .

$$\begin{aligned} & E_g(Y_{R(D|g')}) \\ &= E_g(Y^1 | R(D | g') = 1) \Pr_g(R(D | g') = 1) + E_g(Y^0 | R(D | g') = 0) \Pr_g(R(D | g') = 0) \\ &= E_g(Y^1) E_{g'}(D) + E_g(Y^0)(1 - E_{g'}(D)) \\ &= E_g(Y^0) + E_{g'}(D) E_g(\tau). \end{aligned}$$

A.2. Results in Subsection 2.4. For the CDE,

$$\begin{aligned} & E(Y^{a,0}) - E(Y^{b,0}) \\ &= E(Y^{a,0} | G = a) - E(Y^{b,0} | G = b) \\ &= E_a(Y^0) - E_b(Y^0), \end{aligned}$$

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

For the RIA of PIE,

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

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

For the RIA of PAI,

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

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

Now we derive the difference between between TE and the RIA of TE. We use $Y^{g,D^{g'}}$ to denote the potential outcome of Y under assignment of group g and the D value that would have been taken under the assignment of group g' . And $Y^{g,d}$ denotes the potential outcome of Y under assignment of fixed values g and d .

$$\begin{aligned}
& E(Y^a) - E(Y^b) \\
&= E(Y^{a,D^a}) - E(Y^{b,D^b}) \\
&= E[Y^{a,0}(1 - D^a)] + E[Y^{a,1}D^a] - E[Y^{b,0}(1 - D^a)] - E[Y^{b,1}D^a] \\
&= E(Y^{a,0}) + E[D^a(Y^{a,1} - Y^{a,0})] - E(Y^{b,0}) - E[D^b(Y^{b,1} - Y^{b,0})] \\
&= E(Y^{a,0}) - E(Y^{b,0}) + E(D^a)E(Y^{a,1} - Y^{a,0}) - E(D^b)E(Y^{b,1} - Y^{b,0}) \\
&\quad + \text{Cov}(D^a, Y^{a,1} - Y^{a,0}) - \text{Cov}(D^b, Y^{b,1} - Y^{b,0}). \\
& E(Y^{a,R(D^a)}) - E(Y^{b,R(D^b)}) \\
&= E(Y^{a,0} | R(D^a) = 0) \Pr(R(D^a) = 0) + E(Y^{a,1} | R(D^a) = 1) \Pr(R(D^a) = 1) \\
&\quad - E(Y^{b,0} | R(D^b) = 0) \Pr(R(D^b) = 0) + E(Y^{b,1} | R(D^b) = 1) \Pr(R(D^b) = 1) \\
&= E(Y^{a,0}) \Pr(D^a = 0) + E(Y^{a,1}) \Pr(D^a = 1) \\
&\quad - E(Y^{b,0}) \Pr(D^b = 0) + E(Y^{b,1}) \Pr(D^b = 1) \\
&= E(Y^{a,0}) - E(Y^{b,0}) + E(D^a)E(Y^{a,1} - Y^{a,0}) - E(D^b)E(Y^{b,1} - Y^{b,0}).
\end{aligned}$$

Hence,

$$E(Y^a) - E(Y^b) - [E(Y^{a,R(D^a)}) - E(Y^{b,R(D^b)})] = \text{Cov}(D^a, Y^{a,1} - Y^{a,0}) - \text{Cov}(D^b, Y^{b,1} - Y^{b,0}).$$

A.3. Equation (4). We first note that the outcome disparity can also be decomposed as such:

$$\begin{aligned}
& E_a(Y) - E_b(Y) \\
&= E_a(Y^0) - E_b(Y^0) + E(D)[E_a(\tau) - E_b(\tau)] + E(\tau)[E_a(D) - E_b(D)] \\
&\quad + \text{Cov}_a(D, \tau) - \text{Cov}_b(D, \tau) - [p_a - p_b][E_a(D) - E_b(D)][E_a(\tau) - E_b(\tau)],
\end{aligned}$$

where $E(\tau)$ and $E(D)$ are the overall ATE and treatment prevalence, p_g is the proportion of the population in group g . 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 change in disparity in this variant is equation (3).

A.4. Equation (5).

$$\begin{aligned}
& E_a(Y) - E_b(Y) \\
&= E_a(Y^0) - E_b(Y^0) + E_a(D\tau) - E_b(D\tau) - E_a[E_a(D | \mathbf{Q}) E_a(\tau | \mathbf{Q})] + E_b[E_b(D | \mathbf{Q}) E_b(\tau | \mathbf{Q})] \\
&\quad + E_a[E_a(D | \mathbf{Q}) E_a(\tau | \mathbf{Q})] - E_b[E_b(D | \mathbf{Q}) E_b(\tau | \mathbf{Q})] \\
&= E_a(Y^0) - E_b(Y^0) + E_a[\text{Cov}_a(D, \tau | \mathbf{Q})] - E_b[\text{Cov}_b(D, \tau | \mathbf{Q})] \\
&\quad + \int E_a(D | \mathbf{q}) E_a(\tau | \mathbf{q}) f_a(\mathbf{q}) d\mathbf{q} - \int E_b(D | \mathbf{q}) E_b(\tau | \mathbf{q}) f_b(\mathbf{q}) d\mathbf{q} \\
&= E_a(Y^0) - E_b(Y^0) \\
&\quad + \int [E_a(D | \mathbf{q}) - E_b(D | \mathbf{q})] E_b(\tau | \mathbf{q}) f_b(\mathbf{q}) d\mathbf{q} + \int E_a(D | \mathbf{q}) E_b(\tau | \mathbf{q}) [f_a(\mathbf{q}) - f_b(\mathbf{q})] d\mathbf{q} \\
&\quad + \int [E_a(\tau | \mathbf{q}) - E_b(\tau | \mathbf{q})] E_a(D | \mathbf{q}) f_a(\mathbf{q}) d\mathbf{q} + E_a[\text{Cov}_a(D, \tau | \mathbf{Q})] - E_b[\text{Cov}_b(D, \tau | \mathbf{Q})].
\end{aligned}$$

Note that the last equality uses Assumption 2 (common support).

A.5. Equation (6).

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

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

A.6. Equation (7).

$$\begin{aligned}
& E_a(Y) - E_b(Y) - [E_a(Y^{R(D|a, \mathbf{Q})}) - E_b(Y^{R(D|b, \mathbf{Q})})] \\
&= E_a(Y^0) - E_b(Y^0) + E_a(D\tau) - E_b(D\tau) - [E_a(Y^{R(D|a, \mathbf{Q})}) - E_b(Y^{R(D|b, \mathbf{Q})})] \\
&= E_a[E_a(D\tau | \mathbf{Q})] - E_a[E_a(\tau | \mathbf{Q}) E_a(\tau | \mathbf{Q})] - \{E_b[E_b(D\tau | \mathbf{Q})] - E_b[E_b(\tau | \mathbf{Q}) E_b(\tau | \mathbf{Q})]\} \\
&= E_a[\text{Cov}_a(D, \tau | \mathbf{Q})] - E_b[\text{Cov}_b(D, \tau | \mathbf{Q})].
\end{aligned}$$

Other results in equation (7) follow directly from equation (6).

A.7. Equation (8).

$$\begin{aligned}
& E_b(Y^{R(D|a, \mathbf{Q})}) - E_b(Y) \\
&= E_b(Y^0) + \int E_b(\tau | \mathbf{q}) E_a(D | \mathbf{q}) f_b(\mathbf{q}) d\mathbf{q} - E_b(Y^0) - E_b(D\tau)
\end{aligned}$$

$$\begin{aligned}
&= \int [\mathbb{E}_a(D | \mathbf{q}) - \mathbb{E}_b(D | \mathbf{q})] \mathbb{E}_b(\tau | \mathbf{q}) f_b(\mathbf{q}) d\mathbf{q} \\
&\quad + \mathbb{E}_b[\mathbb{E}_b(D | \mathbf{Q}) \mathbb{E}_b(\tau | \mathbf{Q})] - \mathbb{E}_b[\mathbb{E}_b(D\tau | \mathbf{Q})] \\
&= \int [\mathbb{E}_a(D | \mathbf{q}) - \mathbb{E}_b(D | \mathbf{q})] \mathbb{E}_b(\tau | \mathbf{q}) f_b(\mathbf{q}) d\mathbf{q} - \mathbb{E}_b[\text{Cov}_b(D, \tau | \mathbf{Q})].
\end{aligned}$$

A.8. Equation (9). In the main text, we state that the change in disparity in Lundberg's (2022) CRED is $\mathbb{E}_a(Y) - \mathbb{E}_b(Y) - [\mathbb{E}_a(Y^{R(D|\mathbf{Q})}) - \mathbb{E}_b(Y^{R(D|\mathbf{Q})})]$. Below, we start from its original form that appears in Lundberg (2022) (see his equation (2)), which is equivalent to the form expressed in randomized intervention notation in our main text.

$$\begin{aligned}
&\mathbb{E}_a(Y) - \mathbb{E}_b(Y) - \\
&\quad \{\mathbb{E}_a[\Pr(D = 0 | \mathbf{Q})Y^0 + \Pr(D = 1 | \mathbf{Q})Y^1] - \mathbb{E}_b[\Pr(D = 0 | \mathbf{Q})Y^0 + \Pr(D = 1 | \mathbf{Q})Y^1]\} \\
&= \mathbb{E}_a(Y) - \mathbb{E}_b(Y) - \{\mathbb{E}_a(Y^0) + \mathbb{E}_a[\mathbb{E}(D | \mathbf{Q})\tau] - \mathbb{E}_b(Y^0) - \mathbb{E}_b[\mathbb{E}(D | \mathbf{Q})\tau]\} \\
&= \mathbb{E}_a(D\tau) - \mathbb{E}_a[\mathbb{E}(D | \mathbf{Q}) \mathbb{E}_a(\tau | \mathbf{Q})] - \{\mathbb{E}_b(D\tau) - \mathbb{E}_b[\mathbb{E}(D | \mathbf{Q}) \mathbb{E}_b(\tau | \mathbf{Q})]\} \\
&= \int [\mathbb{E}_a(D\tau | \mathbf{q}) - \mathbb{E}(D | \mathbf{q}) \mathbb{E}_a(\tau | \mathbf{q})] f_a(\mathbf{q}) d\mathbf{q} - \int [\mathbb{E}_b(D\tau | \mathbf{q}) - \mathbb{E}(D | \mathbf{q}) \mathbb{E}_b(\tau | \mathbf{q})] f_b(\mathbf{q}) d\mathbf{q} \\
&= \mathbb{E}_a[\text{Cov}_a(D, \tau | \mathbf{Q})] - \mathbb{E}_b[\text{Cov}_b(D, \tau | \mathbf{Q})] \\
&\quad + \int [\mathbb{E}_a(D | \mathbf{q}) - \mathbb{E}(D | \mathbf{q})] \mathbb{E}_a(\tau | \mathbf{q}) f_a(\mathbf{q}) d\mathbf{q} - \int [\mathbb{E}_b(D | \mathbf{q}) - \mathbb{E}(D | \mathbf{q})] \mathbb{E}_b(\tau | \mathbf{q}) f_b(\mathbf{q}) d\mathbf{q} \\
&= \mathbb{E}_a[\text{Cov}_a(D, \tau | \mathbf{Q})] - \mathbb{E}_b[\text{Cov}_b(D, \tau | \mathbf{Q})] \\
&\quad + \int \Pr(G = b | \mathbf{q}) [\mathbb{E}_a(D | \mathbf{q}) - \mathbb{E}_b(D | \mathbf{q})] \mathbb{E}_a(\tau | \mathbf{q}) f_a(\mathbf{q}) d\mathbf{q} \\
&\quad - \int \Pr(G = a | \mathbf{q}) [\mathbb{E}_b(D | \mathbf{q}) - \mathbb{E}_a(D | \mathbf{q})] \mathbb{E}_b(\tau | \mathbf{q}) f_b(\mathbf{q}) d\mathbf{q} \\
&= \mathbb{E}_a[\text{Cov}_a(D, \tau | \mathbf{Q})] - \mathbb{E}_b[\text{Cov}_b(D, \tau | \mathbf{Q})] \\
&\quad + \int [\mathbb{E}_a(D | \mathbf{q}) - \mathbb{E}_b(D | \mathbf{q})] [\mathbb{E}_a(\tau | \mathbf{q}) f_a(\mathbf{q}) \Pr(G = b | \mathbf{q}) + \mathbb{E}_b(\tau | \mathbf{q}) f_b(\mathbf{q}) \Pr(G = a | \mathbf{q})] d\mathbf{q}.
\end{aligned}$$

Appendix B. With the randomized intervention notation, we present a graphical representation of our unconditional decomposition in Figure A1. The graph visualizes the distinct ways the four components contribute to the observed group disparity in outcome, i.e., one can vary the four components on the graph and obtain different outcome disparities. From this graph, it is clear that the selection component represents the contribution of differential effectiveness of treatment assignment across groups, where effectiveness is defined relative to the random assignment of treatment. Imai and Li (2023) have a graph which is a single-group version of Figure 1. But they do not note the covariance representation of the difference between a treatment assignment rule and the corresponding random assignment with the same treatment prevalence.

Appendix C: 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).

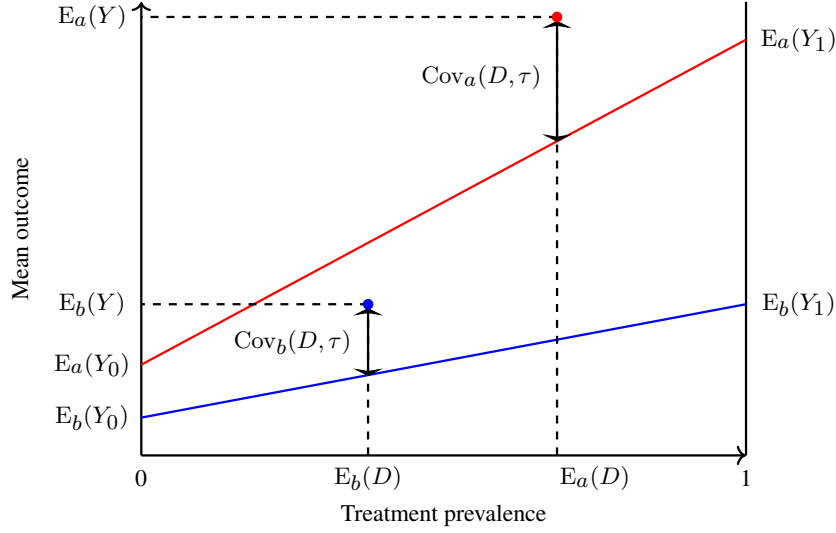


Figure A1. Illustration of the unconditional decomposition. The heights of the red dot and the blue dot respectively indicate the observed mean outcomes in group a and b . On the x axis, the position of each dot corresponds to the observed treatment prevalence for the corresponding group. The red line and the blue line respectively represent the mean potential outcomes in group a and b under hypothetically random assignment of the treatment, which vary by the prevalence of the assigned treatment. The lines are straight due to the hypothetical assignment being random. Naturally, when nobody receives the treatment, the height of each line is $E_g(Y_0)$; and when everyone receives the treatment, the height of the line is $E_g(Y_1)$. The slope of each line is hence $E_g(\tau)$. For a group g , the vertical distance between the dot and the line is $E_g(Y) - E_g(Y^{R(D|g)}) = \text{Cov}_g(D, \tau)$.

Let $\mathbb{1}_{\tilde{o}}(o)$ be the point mass density at a single empirical observation, \tilde{o} . Let subscript \mathcal{P}_t indicate a regular parametric submodel indexed by t . The subscript is omitted for the true model. 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, ξ , 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, Imbens and Ridder \(2003\)](#), 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:

$$h_g := E(w(\mathbf{X}, G) | g)$$

$$h_g(\mathbf{Q}) := E(w(\mathbf{X}, G) | \mathbf{Q}, g).$$

B.1. EIFs for the unconditional decomposition. First of all, note that we only need to derive EIFs for two generic functions, $\xi_{dg} := E\left(Y^{d \frac{w(\mathbf{X}, g)}{h_g}} | g\right)$ for an arbitrary group g ; and $\xi_{dgg'} := E\left(Y^{d \frac{w(\mathbf{X}, g)}{h_g}} | g\right) E\left(D^{\frac{w(\mathbf{X}, g')}{h_{g'}}} | g'\right)$ 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, these estimands can be written as the following functionals:

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

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

We start with ξ_{dg} .

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

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

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

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

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

Since

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

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

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

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

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

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

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

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

which, without survey weights, becomes

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

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

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

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

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

$$\begin{aligned}
& w(\mathbf{X}, g'') | g'' \Big] \frac{1}{\mathbb{E}_{\mathcal{P}_t}(w(\mathbf{X}, g) | 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) | \mathbf{Q}, g)} | \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') | \mathbf{Q}, g')} | \mathbf{Q}, g' \right) w(\mathbf{X}, 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) | \mathbf{Q}, g)} | \mathbf{Q}, g \right) \mathbb{E} \left(D \frac{w(\mathbf{X}, g')}{\mathbb{E}(w(\mathbf{X}, g') | \mathbf{Q}, g')} | \mathbf{Q}, g' \right) w(\mathbf{X}, g'') | g'' \right] \\
&\quad \frac{\partial}{\partial g''} \frac{1}{\mathbb{E}_{\mathcal{P}_t}(w(\mathbf{X}, g) | 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})} | \mathbf{Q}, g \right) \mathbb{E} \left(D \frac{w(\mathbf{X}, g')}{h_{g'}(\mathbf{Q})} | \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) | \mathbf{Q}, g)} | \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') | \mathbf{Q}, g')} | \mathbf{Q}, g' \right) w(\mathbf{X}, g'') | g'' \right] \frac{1}{h_{g''}} \\
&\quad - \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})} | \mathbf{Q}, g \right) \mathbb{E} \left(D \frac{w(\mathbf{X}, g')}{h_{g'}(\mathbf{Q})} | \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) | \mathbf{Q}, g)} | \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') | \mathbf{Q}, g')} | \mathbf{Q}, g' \right) w(\mathbf{X}, g'') | g'' \right] \frac{1}{h_{g''}}.
\end{aligned}$$

And

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

Hence,

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

$$+ \frac{w(\mathbf{X}, g'')}{h_{g''}} \frac{\mathbb{1}(G = g') p_{g''}(\mathbf{Q})}{p_{g'}(\mathbf{Q}) p_{g''}} \frac{w(\mathbf{X}, g')}{h_{g'}(\mathbf{Q})} \left[D - \mathbb{E} \left(D \frac{w(\mathbf{X}, g')}{h_{g'}(\mathbf{Q})} \mid \mathbf{Q}, g' \right) \right] \mathbb{E} \left(Y^d \frac{w(\mathbf{X}, g)}{h_g(\mathbf{Q})} \mid \mathbf{Q}, g \right),$$

which, in the absence of survey weights, simplifies to

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

Finally, recall that $\mathbb{E} \left(Y^d \mid \mathbf{Q}, g \right)$ is identified as $\omega(d, \mathbf{Q}, g) := \mathbb{E}[\mu(d, \mathbf{X}, g) \mid \mathbf{Q}, g]$.

Appendix D. Asymptotic distribution. 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 ξ . Formally, for all d, g and ξ , 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 $\hat{\mu}(d, \mathbf{X}, g)$, $\hat{\pi}(d, \mathbf{X}, g)$, $\hat{p}_g(\mathbf{Q})$, $\hat{\omega}(d, \mathbf{Q}, g)$, 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 ξ_{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,

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

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

where the second inequality uses the Cauchy–Schwarz inequality.

For the empirical process term,

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

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 $\xi_{dgg'}$,

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

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,

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

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

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

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

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

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, $(\mathbb{P}_n - \mathbb{P}) \left[\hat{\phi}_{dgg'g''}(Y, \mathbf{X}) - \phi_{dgg'g''}(Y, \mathbf{X}) \right]$, 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,

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

$$\begin{aligned}
& + \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 [\hat{\omega}(d, \mathbf{Q}, g) - \omega(d, \mathbf{Q}, 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})} D \omega(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 \omega(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 \omega(d, \mathbf{Q}, g) \\
& + \frac{\mathbb{1}(G=g')\hat{p}_{g''}(\mathbf{Q})}{\hat{p}_{g'}(\mathbf{Q})\hat{p}_{g''}} \omega(d, \mathbf{Q}, g) [\mathbb{E}(D | \mathbf{Q}, g') - \hat{\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') [\omega(d, \mathbf{Q}, g) - \hat{\omega}(d, \mathbf{Q}, g)] \\
& + \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') \omega(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') \omega(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') \omega(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 ξ terms satisfy $g = g''$.

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

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

Then it follows from similar arguments as above that $R_{2,dgg'g''} = o_p(n^{-1/2})$ under stated conditions. In particular, for the last line, note that if $g = g' = g''$,

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

so $\|\omega(d, \mathbf{Q}, g) - \hat{\omega}(d, \mathbf{Q}, g)\| \|\hat{\mathbb{E}}(D | \mathbf{Q}, g) - \mathbb{E}(D | \mathbf{Q}, g)\| = o_p(n^{-1/2})$, $\forall d, g$, is sufficient for the last line.

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

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

so the following conditions are sufficient for the last line: for some $\zeta < \infty$, $\left| \frac{\hat{p}_g(\mathbf{Q})}{\hat{p}_{g'}(\mathbf{Q})} \right| \leq \zeta$ with probability 1, $\|\hat{\mathbb{E}}(D | \mathbf{Q}, g) - \mathbb{E}(D | \mathbf{Q}, g)\| = o_p(n^{-1/2})$, and $\|\omega(d, \mathbf{Q}, g) - \hat{\omega}(d, \mathbf{Q}, g)\| = o_p(1)$, $\forall d, g, g'$.

C.2.2 When $g' = g''$. Note that all ξ terms satisfy $g' = g''$ for the conditional effect component.

$$\begin{aligned}
& R_{2,dgg'g''} \\
& = \hat{\xi}_{dgg'g'} - \xi_{dgg'g'} \\
& + \mathbb{P} \left\{ -\frac{\mathbb{1}(G=g')}{\hat{p}_{g'}} \hat{\xi}_{dgg'g'} \right\}
\end{aligned}$$

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

Under stated conditions, $R_{2,dgg'g''} = o_p(n^{-1/2})$. In particular, for the last line, the following conditions are sufficient: for some $\zeta < \infty$, $\left| \frac{\hat{p}_{g'}(\mathbf{Q})}{\hat{p}_g(\mathbf{Q})} \right| \leq \zeta$ with probability 1, $\left\| \hat{\mathbb{E}}(D | \mathbf{Q}, g) - \mathbb{E}(D | \mathbf{Q}, g) \right\| = o_p(1)$, and $\|\omega(d, \mathbf{Q}, g) - \hat{\omega}(d, \mathbf{Q}, g)\| = o_p(n^{-1/2})$, $\forall d, g, g'$.

C.2.3 Summary of convergence rate conditions. When $g = g' = g''$, which holds for conditional selection, $\|\hat{\omega}(d, \mathbf{Q}, g) - \omega(d, \mathbf{Q}, g)\| \left\| \hat{\mathbb{E}}(D | \mathbf{Q}, g) - \mathbb{E}(D | \mathbf{Q}, g) \right\| = o_p(n^{-1/2})$ is sufficient. Hence, we obtain a form of rate double robustness with respect to $\omega(d, \mathbf{Q}, g)$ and $\mathbb{E}(D | \mathbf{Q}, g)$. Second, when $g = g'' \neq g'$, which holds for conditional prevalence, the following set of conditions is sufficient¹: $\left\| \hat{\mathbb{E}}(D | \mathbf{Q}, g) - \mathbb{E}(D | \mathbf{Q}, g) \right\| = o_p(n^{-1/2})$. Third, when $g' = g'' \neq g$, a sufficient set is: $\|\hat{\omega}(d, \mathbf{Q}, g) - \omega(d, \mathbf{Q}, g)\| = o_p(n^{-1/2})$. Therefore, the assumption is weaker for conditional selection than other components, which is also the case for Theorem 5.

Appendix E. Multiple robustness. Below, we prove the multiple robustness of $\hat{\xi}_{dgg'g''}$, which is stated in Theorem 5. 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 $\mathbb{E}(D | \mathbf{Q}, g)$, the one-step estimator $\hat{\xi}_{dgg'g''}$ converges in probability to

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

¹Recall that by Assumption 6a, for some $\zeta < \infty$, $\hat{p}_g(\mathbf{Q})/\hat{p}_{g'}(\mathbf{Q}) \leq \zeta$ with probability 1, $\forall g, g'$. And by Assumption 6b, $\|\hat{\omega}(d, \mathbf{Q}, g) - \omega(d, \mathbf{Q}, g)\| = o_p(1)$, $\left\| \hat{\mathbb{E}}(D | \mathbf{Q}, g) - \mathbb{E}(D | \mathbf{Q}, g) \right\| = o_p(1)$, $\forall d, g$.

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

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

$$\begin{aligned}
& \mathbb{E} \left\{ \frac{p_g(\mathbf{Q})}{p_g} \mathbb{E} \left\{ \frac{\pi(d, \mathbf{X}, g)}{\tilde{\pi}(d, \mathbf{X}, g)} [\mu(d, \mathbf{X}, g) - \tilde{\mu}(d, \mathbf{X}, g)] + \tilde{\mu}(d, \mathbf{X}, g) \mid \mathbf{Q}, g \right\} \tilde{\mathbb{E}}(D \mid \mathbf{Q}, g') \right\} \\
& + \mathbb{E} \left\{ \frac{p_{g'}(\mathbf{Q}) \tilde{p}_g(\mathbf{Q})}{\tilde{p}_{g'}(\mathbf{Q}) p_g} \left[\mathbb{E}(D \mid \mathbf{Q}, g') - \tilde{\mathbb{E}}(D \mid \mathbf{Q}, g') \right] \tilde{\omega}(d, \mathbf{Q}, g) \right\}.
\end{aligned}$$

Next, when $g' = g''$, the probability limit becomes

$$\begin{aligned}
& \mathbb{E} \left\{ \frac{p_g(\mathbf{Q}) \tilde{p}_{g'}(\mathbf{Q})}{\tilde{p}_g(\mathbf{Q}) p_{g'}} \mathbb{E} \left\{ \frac{\pi(d, \mathbf{X}, g)}{\tilde{\pi}(d, \mathbf{X}, g)} [\mu(d, \mathbf{X}, g) - \tilde{\mu}(d, \mathbf{X}, g)] + \tilde{\mu}(d, \mathbf{X}, g) \mid \mathbf{Q}, g \right\} \tilde{\mathbb{E}}(D \mid \mathbf{Q}, g') \right\} \\
& - \mathbb{E} \left\{ \frac{p_g(\mathbf{Q}) \tilde{p}_{g'}(\mathbf{Q})}{\tilde{p}_g(\mathbf{Q}) p_{g'}} \tilde{\omega}(d, \mathbf{Q}, g) \tilde{\mathbb{E}}(D \mid \mathbf{Q}, g') \right\} + \mathbb{E} \left\{ \frac{p_{g'}(\mathbf{Q})}{p_{g'}} \mathbb{E}(D \mid \mathbf{Q}, g') \tilde{\omega}(d, \mathbf{Q}, g) \right\}.
\end{aligned}$$

Finally, when $g = g' = g''$, it becomes

$$\begin{aligned}
& \mathbb{E} \left\{ \frac{p_g(\mathbf{Q})}{p_g} \mathbb{E} \left\{ \frac{\pi(d, \mathbf{X}, g)}{\tilde{\pi}(d, \mathbf{X}, g)} [\mu(d, \mathbf{X}, g) - \tilde{\mu}(d, \mathbf{X}, g)] + \tilde{\mu}(d, \mathbf{X}, g) \mid \mathbf{Q}, g \right\} \tilde{\mathbb{E}}(D \mid \mathbf{Q}, g) \right\} \\
& + \mathbb{E} \left\{ \frac{p_g(\mathbf{Q})}{p_g} \left[\mathbb{E}(D \mid \mathbf{Q}, g) - \tilde{\mathbb{E}}(D \mid \mathbf{Q}, g) \right] \tilde{\omega}(d, \mathbf{Q}, g) \right\}.
\end{aligned}$$

Then it is relatively easy to verify the results in Theorem 5 for the case of $g = g''$. Note that the condition for conditional selection to be consistently estimated is the weakest compared with other components. In particular, for conditional selection, it is not necessary to consistently estimate $p_g(\mathbf{Q})$.

Appendix F. 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

$$\text{Cov}_g(D, \tau) = [\mathbb{E}_g(\tau \mid D = 1) - \mathbb{E}_g(\tau)] \mathbb{E}_g(D).$$

First, our definition is closely related to the “sorting on gains” concept in economics (Heckman and Vytlacil, 2005; Heckman, Humphries and Veramendi, 2018), which is defined as the difference between ATT and ATE, i.e.,

$$\mathbb{E}_g(\tau \mid D = 1) - \mathbb{E}_g(\tau).$$

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

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 selectivity 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 selectivity term.

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

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 $E_g(Y^0 | D = 1) - E_g(Y^0 | D = 0)$, which is selection on baseline outcome. Therefore, the selectivity term in the great equalizer literature captures selection on baseline outcome, on 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, Brand and Jann, 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 3 (conditional ignorability),

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

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 G. Supplemental tables.

Table A1. Relationship between Baseline Covariates and Outcome Missingness

Covariate		Missing percentage or correlation coefficient	P value
Parental income		-0.079 [#]	< 0.001
Race	White	25.27	
	Black	14.45	< 0.001
	Hispanic	12.38	
Gender	Male	21.09	
	Female	19.63	0.385
Mother's year of schooling		0.003 [#]	0.870
Parental presence	Yes	19.31	
	No	22.21	0.089
Number of siblings		0.002 [#]	0.900
Urban residence	Yes	19.65	
	No	22.86	0.108
Expecting bachelor degree or higher	Yes	19.28	
	No	21.20	0.250
AFQT score		0.021 [#]	0.294
Age		-0.015 [#]	0.453
Friends expecting bachelor degree or higher	Yes	19.15	
	No	21.31	0.190
Rotter score		-0.002 [#]	0.914
Rosenberg score		-0.013 [#]	0.517
School satisfaction		0.003	0.888
Speak foreign language at home	Yes	18.07	
	No	20.91	0.174
	Not in SMSA	22.73	
Metropolitan Statistical Area category	In SMSA, not central city	18.00	
	In SMSA, in central city	22.32	0.057
	In SMSA, central city unknown	18.59	
Separate from mother	Yes	25.83	
	No	20.08	0.158
Mother working	Yes	18.27	
	No	23.00	0.004
SMSA	Northeast	23.21	
	North central	22.42	
	South	18.75	0.031
	West	16.95	

Note: N=2580. The sample is individuals with no missing values in any baseline covariates. In the column for missing percentage or correlation coefficient, values with # are correlation coefficients between covariates and the outcome missingness indicator. P values for missing percentages are based on Chi-squared tests.

Table A2. Group-specific Estimates for the Unconditional Decomposition

	Top 40% income	Bottom 40% income	Top-Bottom
Treatment proportion	0.337 (0.310, 0.364)	0.086 (0.067, 0.105)	0.251 (0.218, 0.284)
Gradient Boosted Machine			
Baseline outcome	0.599 (0.581, 0.616)	0.411 (0.393, 0.429)	0.188 (0.162, 0.213)
ATE	0.137 (0.100, 0.174)	0.110 (0.044, 0.177)	0.027 (-0.049, 0.103)
$Cov(D, \tau)$	-0.011 (-0.021, 0.000)	0.004 (-0.002, 0.009)	-0.014 (-0.026, -0.002)
Neural Networks			
Baseline outcome	0.599 (0.582, 0.617)	0.411 (0.393, 0.430)	0.188 (0.163, 0.213)
ATE	0.145 (0.103, 0.188)	0.123 (0.052, 0.195)	0.022 (-0.061, 0.105)
$Cov(D, \tau)$	-0.014 (-0.027, -0.002)	0.002 (-0.004, 0.009)	-0.017 (-0.031, -0.002)
Random Forests			
Baseline outcome	0.593 (0.578, 0.609)	0.410 (0.392, 0.429)	0.183 (0.159, 0.207)
ATE	0.140 (0.107, 0.172)	0.116 (0.040, 0.192)	0.024 (-0.058, 0.107)
$Cov(D, \tau)$	-0.006 (-0.014, 0.001)	0.004 (-0.002, 0.010)	-0.010 (-0.020, -0.001)
Parametric Regressions			
Baseline outcome	0.598 (0.580, 0.615)	0.410 (0.393, 0.428)	0.187 (0.162, 0.212)
ATE	0.121 (0.080, 0.162)	0.145 (0.059, 0.232)	-0.024 (-0.120, 0.072)
$Cov(D, \tau)$	-0.005 (-0.017, 0.008)	0.001 (-0.006, 0.008)	-0.006 (-0.020, 0.009)

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

Table A3. Unconditional Decomposition Estimates

	GBM	Neural Networks	Random Forests	Parametric
Total	0.210 (0.188, 0.232)	0.210 (0.188, 0.232)	0.210 (0.188, 0.232)	0.210 (0.188, 0.232)
Baseline	0.188 (0.162, 0.213)	0.188 (0.163, 0.213)	0.183 (0.159, 0.207)	0.187 (0.162, 0.212)
Prevalence	0.028 (0.011, 0.045)	0.031 (0.012, 0.050)	0.029 (0.009, 0.049)	0.036 (0.014, 0.059)
Effect	0.009 (-0.017, 0.035)	0.007 (-0.021, 0.035)	0.008 (-0.020, 0.036)	-0.008 (-0.040, 0.024)
Selection	-0.014 (-0.026, -0.002)	-0.017 (-0.031, -0.002)	-0.010 (-0.020, -0.001)	-0.006 (-0.020, 0.009)
Change in disparity	0.024 (0.005, 0.043)	0.029 (0.007, 0.050)	0.025 (0.002, 0.048)	0.035 (0.008, 0.063)

Note: 95% confidence intervals are in the parentheses. Weight stabilization is used. For ML models, cross-fitting is used. The change in disparity uses the definition in Jackson and Vanderweele's (2018) URED.

Table A4: Conditional Decomposition Estimates

	GBM	Neural Networks	Random Forests	Parametric
Total	0.210 (0.188,0.232)	0.210 (0.188,0.232)	0.210 (0.188,0.232)	0.210 (0.188,0.232)
Baseline	0.188 (0.162,0.213)	0.188 (0.163,0.213)	0.183 (0.159,0.207)	0.187 (0.162,0.212)
Conditional prevalence	0.007 (0.002,0.012)	0.010 (0.003,0.016)	0.013 (0.007,0.019)	0.009 (0.002,0.015)
Conditional effect	-0.028 (-0.061,0.006)	-0.012 (-0.055,0.031)	-0.025 (-0.054,0.004)	-0.023 (-0.054,0.007)
Conditional selection	-0.004 (-0.010,0.002)	-0.006 (-0.014,0.003)	-0.006 (-0.014,0.002)	-0.002 (-0.010,0.005)
Q -distribution	0.047 (0.022,0.072)	0.030 (-0.004,0.063)	0.045 (0.023,0.067)	0.040 (0.018,0.061)
Change in disparity	0.006 (-0.001,0.013)	0.006 (-0.003,0.016)	0.007 (-0.002,0.017)	0.006 (-0.003,0.016)

Note: 95% confidence intervals are in the parentheses. Weight stabilization is used. For ML models, cross-fitting is used. The change in disparity uses the definition in Jackson’s (2021) CRED.

In addition, in Table A5, we present a set of estimates for the conditional decomposition as a robustness check. This is motivated by Assumption 6c, which requires parametric convergence rate for either $\hat{E}(D | Q, g)$ or $\hat{\omega}(d, Q, g)$ for conditional prevalence and conditional effect. Thus, to make asymptotic inference more exact, we implement the following procedure. We estimate $\mu(d, X, g)$, $\pi(d, X, g)$, and $p_g(Q)$ using ML with cross-fitting. Then we estimate $E(D | Q, g)$ and $\omega(d, Q, g)$ parametrically using linear or logistic regressions without cross-fitting. For the parametric models, we include the group indicator, the AFQT score, the squared AFQT score, and the interactions between the group indicator and the AFQT variables. For $\omega(d, Q, g)$, we apply the pseudo-outcome approach detailed in the main text, but with a different cross-fitting procedure. Here, we use cross-fitted estimates of $\mu(d, X, g)$ and $\pi(d, X, g)$ to construct the pseudo-outcomes, then we estimate $\omega(d, Q, g)$ without cross-fitting. The findings are substantively very consistent across Tables A4 and A5.

Table A5: Conditional Decomposition Estimates with Mixed Nonparametric and Parametric Models

	GBM	Neural Networks	Random Forests
Total	0.210 (0.188,0.232)	0.210 (0.188,0.232)	0.210 (0.188,0.232)
Baseline	0.182 (0.157,0.208)	0.188 (0.163,0.213)	0.183 (0.159,0.207)
Conditional prevalence	0.008 (0.002,0.013)	0.008 (0.002,0.014)	0.008 (0.002,0.014)
Conditional effect	-0.019 (-0.048,0.010)	-0.018 (-0.049,0.013)	-0.024 (-0.054,0.006)
Conditional selection	-0.001 (-0.008,0.005)	-0.006 (-0.014,0.003)	-0.004 (-0.010,0.002)
Q -distribution	0.040 (0.022,0.059)	0.037 (0.016,0.058)	0.047 (0.025,0.069)
Change in disparity	0.007 (-0.000,0.014)	0.006 (-0.002,0.015)	0.006 (-0.003,0.014)

Note: 95% confidence intervals are in the parentheses. Weight stabilization is used. $\mu(d, X, g)$, $\pi(d, X, g)$, and $p_g(Q)$ are estimated using ML with cross-fitting. $E(D | Q, g)$ and $\omega(d, Q, g)$ are estimated using parametric models without cross-fitting. The change in disparity uses the definition in Jackson’s (2021) CRED.

REFERENCES

BRAND, J. E. and XIE, Y. (2010). Who Benefits Most from College?: Evidence for Negative Selection in Heterogeneous Economic Returns to Higher Education. *American Sociological Review* **75** 273–302. <https://doi.org/10.1177/0003122410363567>

- BRAND, J. E., XU, J., KOCH, B. and GERALDO, P. (2021). Uncovering Sociological Effect Heterogeneity Using Tree-Based Machine Learning. *Sociological Methodology* **51** 189–223. <https://doi.org/10.1177/0081175021993503>
- CHERNOZHUKOV, V., CHETVERIKOV, D., DEMIRER, M., DUFLO, E., HANSEN, C. and NEWEY, W. (2017). Double/Debiased/Neyman Machine Learning of Treatment Effects. *American Economic Review* **107** 261–265. <https://doi.org/10.1257/aer.p20171038>
- FIEL, J. E. (2020). Great Equalizer or Great Selector? Reconsidering Education as a Moderator of Intergenerational Transmissions. *Sociology of Education* 003804072092788. <https://doi.org/10.1177/0038040720927886>
- FISHER, A. and KENNEDY, E. H. (2021). Visually Communicating and Teaching Intuition for Influence Functions. *The American Statistician* **75** 162–172. <https://doi.org/10.1080/00031305.2020.1717620>
- HAHN, J. (1998). On the Role of the Propensity Score in Efficient Semiparametric Estimation of Average Treatment Effects. *Econometrica* **66** 315. <https://doi.org/10.2307/2998560>
- HECKMAN, J. J., HUMPHRIES, J. E. and VERAMENDI, G. (2018). Returns to Education: The Causal Effects of Education on Earnings, Health, and Smoking. *Journal of Political Economy* **126** 50.
- HECKMAN, J. J. and VYTLACIL, E. (2005). Structural Equations, Treatment Effects, and Econometric Policy Evaluation I. *Econometrica* **73** 669–738. <https://doi.org/10.1111/j.1468-0262.2005.00594.x>
- HINES, O., DUKES, O., DIAZ-ORDAZ, K. and VANSTEELANDT, S. (2022). Demystifying Statistical Learning Based on Efficient Influence Functions. *The American Statistician* 1–13. <https://doi.org/10.1080/00031305.2021.2021984>
- HIRANO, K., IMBENS, G. W. and RIDDER, G. (2003). Efficient Estimation of Average Treatment Effects Using the Estimated Propensity Score. *Econometrica* **71** 1161–1189.
- ICHIMURA, H. and NEWEY, W. K. (2022). The influence function of semiparametric estimators. *Quantitative Economics* **13** 29–61. <https://doi.org/10.3982/QE826>
- IMAI, K. and LI, M. L. (2023). Experimental Evaluation of Individualized Treatment Rules. *Journal of the American Statistical Association* **118** 242–256. <https://doi.org/10.1080/01621459.2021.1923511>
- JACKSON, J. W. (2021). Meaningful Causal Decompositions in Health Equity Research: Definition, Identification, and Estimation Through a Weighting Framework. *Epidemiology* **32** 282–290. <https://doi.org/10.1097/EDE.0000000000001319>
- JACKSON, J. W. and VANDERWEELE, T. J. (2018). Decomposition Analysis to Identify Intervention Targets for Reducing Disparities. *Epidemiology* **29** 825–835. <https://doi.org/10.1097/EDE.0000000000000901>
- KARLSON, K. B. (2019). College as equalizer? Testing the selectivity hypothesis. *Social Science Research* **80** 216–229. <https://doi.org/10.1016/j.ssresearch.2018.12.001>
- KENNEDY, E. H. (2022). Semiparametric doubly robust targeted double machine learning: a review. Number: arXiv:2203.06469 arXiv:2203.06469 [stat].
- LUNDBERG, I. (2022). The Gap-Closing Estimand: A Causal Approach to Study Interventions That Close Disparities Across Social Categories. *Sociological Methods & Research* 004912412110557. <https://doi.org/10.1177/00491241211055769>
- MORGAN, S. L. and WINSHIP, C. (2014). *Counterfactuals and Causal Inference: Methods And Principles For Social Research*, 2nd edition ed. *Analytical Methods for Social Research*. Cambridge University Press, Cambridge, UK.
- XIE, Y., BRAND, J. E. and JANN, B. (2012). Estimating Heterogeneous Treatment Effects with Observational Data. *Sociological Methodology* **42** 314–347. <https://doi.org/10.1177/0081175012452652>
- ZHOU, X. (2019). Equalization or Selection? Reassessing the “Meritocratic Power” of a College Degree in Intergenerational Income Mobility. *American Sociological Review* **84** 459–485. <https://doi.org/10.1177/0003122419844992>