

# *Causal Inference for Time Dependent Treatments*

*Ashley I. Naimi, PhD*

## Correlation and Causation

In his *The Grammar of Science*, Karl Pearson (1911) wrote “[b]eyond such discarded fundamentals as ‘matter’ and ‘force’ lies still another fetish amidst the inscrutable arcana of modern science, namely, the category of cause and effect.” He suggested that rather than pursue an understanding of cause-effect relations, scientists would be best served by measuring correlations through tables that classify individuals into specific categories. “Such a table is termed a contingency table, and the ultimate scientific statement of description of the relation between two things can always be thrown back upon such a contingency table.”

Over a century later, a majority of statistics courses treat causal inference by simply stating that “correlation is not causation.” This treatment is hardly sufficient, for at least two reasons: 1) As scientists, our primary interest is (should be) in cause-effect relations; 2) People continue to conflate correlation with causation<sup>1</sup>. For both of these reasons, we very much need to understand the conditions that would allow us to understand causality better. This is what “causal inference” is all about.

I adopt the view that **the causal and statistical aspects of a scientific study should be kept as separate as possible**. The objective is to first articulate the conditions under which causal inference is possible, and then to understand what statistical tools will enable us to answer the causal question.<sup>2</sup> Causal inference tells us what we should estimate, and whether we can. Statistics tells us how to estimate it. By implication, we should avoid treating statistical models as if they were causal. Furthermore, to the best of our ability, we should avoid imposing unnecessary parametric assumptions on the causal models that we believe are generating the data. I will try to clarify what I intend by “imposing unnecessary parametric assumptions” later.

## Introduction to Causal Inference

“Causal inference” deals primarily with the formal mechanisms by which we can combine data, assumptions, and models to interpret a correlation (or association) as a causal relation.<sup>3</sup> The framework

<sup>1</sup> Daniel Westreich and I reviewed a book in which the allure of “Big Data” was so strong, the authors quickly forgot that correlation  $\neq$  causation. See Ashley I. Naimi and Westreich (2014)

<sup>2</sup> Loosely speaking: Causal inference is the “what?” Statistics is the “how?” Epidemiology is the “why?”

<sup>3</sup> There are a number of introductory books and articles on causal inference in the empirical sciences. Here are some excellent options: Hernán and Robins (Forthcoming 2021), Pearl, Glymour, and Jewell (2016), Imbens and Rubin (2015)

by which we define what we mean by “causal relation” or “causal effect” is the **potential outcomes framework**.

A central notion in the potential outcomes framework is the counterfactual. This notion stems from the intuitive and informal practice of interpreting cause-effect relations as **circumstances (e.g., health outcomes) that would have arisen had things (e.g., exposures) been different**.

While this intuition serves an important purpose, it is not sufficient for doing rigorous science. Suppose we ask: “what is the effect of smoking on CVD risk, irrespective of smoking’s effect on body weight?” This question seems clear and intuitive. To answer this question, we would do a study in which we collect data, enter these into a computer, perform some calculations, and obtain a number (the “effect”).

But there is a problem.<sup>4</sup> The calculations performed by the computer are **rigorously defined mathematical objects**. On the other hand, **english language sentences about cause effect relations are ambiguous**. For example, the “effect of smoking” can mean many different things:

- All people smoke any tobacco ever versus no people smoke tobacco ever.
- All people smoke 3 cigarettes per day versus all people smoke 2 cigarettes per day.
- All people who have smoked any tobacco in the last 15 years cease to smoke any tobacco whatsoever.

Similarly, “irrespective of” can mean a number of things:

- The effect of smoking on CVD risk that would be observed in a hypothetical world where smoking did not affect body mass?
- The effect of smoking on CVD risk if everyone were set to “normal” body mass?
- The effect of smoking on CVD risk if everyone were held at the body mass they had in the month prior to study entry?

<sup>4</sup> This problem was articulated by J. Robins (1987), and I am using the example from his paper.

But the numerical strings of data and the computer algorithms applied to these data are well defined mathematical objects, which do not admit such ambiguity. Depending on several choices, including the data, how variables are coded, and the modeling strategy, the computer is being told which question to answer. There is a lot of potential uncertainty in the space between the English language sentences we use to ask causal questions, and the computer algorithms we use to answer those questions. Causal inference is about clarifying this uncertainty.

## Complex Longitudinal Data

In this course, you have already encountered causal inference via potential outcomes when exposure under study is measured once (i.e., time fixed). In this lecture, we will focus on complex longitudinal data, and the complications that may arise when dealing with such data. For clarity, let's define complex longitudinal data. We will be dealing with data from a cohort study, individuals sampled from a well-defined target population, and clear study start and stop times (i.e., closed cohort). Data from such a cohort are **longitudinal** when they are measured repeatedly over time.<sup>5</sup>

Different scenarios can lead to longitudinal data:

1. exposure and covariates do not vary over time, but the study outcome can occur more than once
2. exposure and covariates vary over time, but the study outcome can only occur once
3. exposure and covariates vary over time, and the study outcome can occur more than once

Scenario 2 is the classical situation that statisticians refer to as "longitudinal" data or correlated data. Here, we will deal with data that from scenarios 2 and 3. Repeated exposure, covariate, and/or outcome measurement is what leads to "longitudinal" data. But why complex?

Repeated measurement over time creates the opportunity for us to capture complex causal relations between past and future covariates. Suppose we measure an exposure twice over follow-up, a covariate

<sup>5</sup> Another such form is when data are measured repeatedly across space. We will not be dealing with these data here.

once, and the outcome at the end of follow-up (Figure 1). If we can assume that past exposure/covariate values do not affect future exposure/covariate values (usually a very risky assumption), we might not consider these data “complex,” because we can use many standard methods we already know to analyze these data.



Figure 1: Longitudinal data that might not be considered ‘complex’ because there is no feedback between exposure and covariates.

On the other hand, if past exposure/covariates affect future exposure/covariates in such a way that prior exposures or covariates confound future exposures (Figure 2), more advanced analytic techniques are needed.

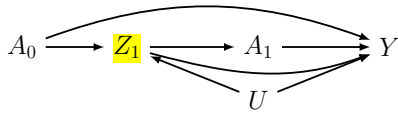


Figure 2: Causal diagram representing the relation between anti-retroviral treatment at time 0 ( $A_0$ ), HIV viral load just prior to the second round of treatment ( $Z_1$ ), anti-retroviral treatment status at time 1 ( $A_1$ ), the CD4 count measured at the end of follow-up ( $Y$ ), and an unmeasured common cause ( $U$ ) of HIV viral load and CD4.

Here, we will learn why this distinction is important, and how to use g methods to account for this type of complex time-varying confounding.

### *G Methods for Time Dependent Covariates*

Robins’ g methods enable the identification and estimation of the effects of generalized treatment, exposure, or intervention plans. G methods are a family of methods that include the g formula, marginal structural models, and structural nested models.<sup>6</sup> They provide **consistent** estimates of contrasts (e.g. differences, ratios) of average potential outcomes under a less restrictive set of identification conditions than standard regression methods (e.g. linear, logistic, Cox regression) (J. M. Robins and Hernán 2009). Specifically, standard regression **requires no feedback between time-varying**

<sup>6</sup> There are three g methods: the parametric g formula and inverse probability weighting. These two are used to estimate the parameters of a marginal structural model. Then there is g estimation (different from the g formula). This is used to estimate the parameters of a structural nested model.

treatments and time-varying confounders, while g methods do not.

Robins and Hern{a}n J. M. Robins and Hern{a}n (2009) have provided a technically comprehensive worked example of each of the three g methods. Here, we present a corresponding worked example that illustrates the need for and use of g methods, while minimizing technical details.<sup>7</sup>

### Example

Our research question concerns the effect of treatment for HIV on CD4 count. Table 1 presents data from a hypothetical observational cohort study ( $A = 1$  for treated,  $A = 0$  otherwise). Treatment is measured at baseline ( $A_0$ ) and once during follow up ( $A_1$ ). The sole covariate is elevated HIV viral load ( $Z = 1$  for those with  $> 200$  copies/ml,  $Z = 0$  otherwise), which is constant by design at baseline ( $Z_0 = 1$ ) and measured once during follow up just prior to the second treatment ( $Z_1$ ). The outcome is CD4 count measured at the end of follow up in units of cells/mm<sup>3</sup>. The CD4 outcome in Table 1 is summarized (averaged) over the participants at each level of the treatments and covariate.

$A_0$	$Z_1$	$A_1$	$Y$	$N$
0	0	0	87.29	209,271
0	0	1	112.11	93,779
0	1	0	119.65	60,654
0	1	1	144.84	136,293
1	0	0	105.28	134,781
1	0	1	130.18	60,789
1	1	0	137.72	93,903
1	1	1	162.83	210,527

The number of participants is provided in the rightmost column of Table 1. In this hypothetical study of one million participants we ignore random error and focus on identifying the parameters defining our causal effect of interest, which we describe next.

Based on Figure 2, the average outcome in our simple data generating structure may be composed of several parts: the effects of

<sup>7</sup> There are a handful of worked examples and tutorials on the use of g methods to estimate effects in complex longitudinal data. These include J. M. Robins and Hern{a}n (2009), Daniel et al. (2013), Keil et al. (2014), the paper on which these notes are based Ashley I. Naimi, Cole, and Kennedy (2017). Additionally, Hern{a}n and Robins (Forthcoming 2021) is an excellent, comprehensive, and very accessible introduction to causal inference generally, and g methods specifically.

Table 1: Prospective study data illustrating the number of subjects ( $N$ ) within each possible combination of treatment at time 0 ( $A_0$ ), HIV viral load just prior to the second round of treatment ( $Z_1$ ), and treatment status for the 2nd round of treatment ( $A_1$ ). The outcome column ( $Y$ ) corresponds to the mean of  $Y$  within levels of  $A_0$ ,  $Z_1$ ,  $A_1$ . Note that HIV viral load at baseline is high ( $Z_0 = 1$ ) for everyone by design.

$A_0$ ,  $Z_1$ , and  $A_1$ ; the two-way interactions between  $A_0$  and  $Z_1$ ,  $A_0$  and  $A_1$ , and  $A_1$  and  $Z_1$ ; and the three-way interaction between  $A_0$ ,  $Z_1$ , and  $A_1$ . These components (some whose magnitudes may be zero) can be used to “build up” a contrast of substantive interest. Here, we focus on the average causal effect of always taking treatment ( $a_0 = 1, a_1 = 1$ ) compared to never taking treatment ( $a_0 = 0, a_1 = 0$ ),<sup>8</sup>

$$\begin{aligned}\psi &= E(Y^{a_0=1, a_1=1}) - E(Y^{a_0=0, a_1=0}) \\ &= E(Y^{a_0=1, a_1=1} - Y^{a_0=0, a_1=0}),\end{aligned}\tag{1}$$

where expectations  $E(\cdot)$  are taken with respect to the target population from which our sample is a random draw. This average causal effect consists of the joint effect of  $A_0$  and  $A_1$  on  $Y$ .<sup>9</sup> Here,  $Y^{a_0, a_1}$  represents a potential outcome value that would have been observed had the exposures been set to specific levels  $a_0$  and  $a_1$ . This potential outcome is distinct from the observed (or actual) outcome.<sup>10</sup>

This average causal effect  $\psi = E(Y^{a_0, a_1} - Y^{0, 0})$  is a *marginal effect* because it averages (or marginalizes) over all individual-level effects in the population. We can write this effect as  $E(Y^{a_0, a_1} - Y^{0, 0}) = \psi_0 a_0 + \psi_1 a_1 + \psi_2 a_0 a_1$ , which states that our average causal effect  $\psi$  may be composed of two exposure main effects (e.g.,  $\psi_0$  and  $\psi_1$ ) and their two-way interaction ( $\psi_2$ ). This marginal effect  $\psi$  is indifferent to whether the  $A_1$  component ( $\psi_1 + \psi_2$ ) is modified by  $Z_1$ : whether such effect modification is present or absent, the marginal effect represents a meaningful answer to the question: what is the effect of  $A_0$  and  $A_1$  in the entire population?

Alternatively, we may wish to estimate this effect *conditional on certain values of another covariate*. A conditional effect would arise if, for example, one was specifically interested in effect measure modification by  $Z_1$ . When properly modeled, this conditional effect represents a meaningful answer to the question: what is the effect of  $A_0$  and  $A_1$  in those who receive  $Z_1 = 1$  versus those who receive  $Z_1 = 0$ ? Modeling such effect measure modification by time-varying covariates is the fundamental issue that distinguishes marginal structural from structural nested models. We thus return to this issue later. For simplicity, we define our effect of interest as  $\psi = \psi_0 + \psi_1 + \psi_2$ ,

<sup>8</sup> Alternate notation for potential outcomes includes:  $Y_x$ ,  $Y(x)$ ,  $Y \mid \text{Set}(X = x)$ , and  $Y \mid \text{do}(X = x)$ .

<sup>9</sup>

<sup>10</sup> Note this distinction is subtle, and often overlooked. Importantly, one can only equate the potential outcome with the observed outcome under the observed exposure if **counterfactual consistency** holds.

and we explore a data example with no effect modification by time-varying confounders.

## ASSUMPTIONS

Our average causal effect is defined as a function of two averages that would be observed if everybody in the population were exposed (or unexposed) at both time points. Yet we cannot directly acquire information on these averages because in any given sample, some individuals will be unexposed (or exposed). Part of our task therefore involves justifying use of averages among subsets of the population as what would be observed in the whole population.<sup>11</sup> This is accomplished by making three main assumptions.

<sup>11</sup> Understanding what this justification entails is the fundamental charge of causal inference.

Counterfactual consistency (S. R. Cole and Frangakis 2009) allows us to equate observed outcomes among those who received a certain exposure value to the potential outcomes that would be observed under the same exposure value:

$$E(Y \mid A_0 = a_0, A_1 = a_1) = E(Y^{a_0, a_1} \mid A_0 = a_0, A_1 = a_1)$$

The status of this assumption remains unaffected by the choice of analytic method (e.g., standard regression versus g methods). Rather, this assumption's validity depends on the nature of the exposure assignment mechanism (VanderWeele and Hernán 2013). Under counterfactual consistency, we partially identify our average causal effect.

Next, we assume exchangeability (Greenland and Robins 1986). Exchangeability implies that the potential outcomes under exposures  $a_0$  and  $a_1$  (denoted  $Y^{a_0, a_1}$ ) are independent of the actual (or observed) exposures  $A_0$  and  $A_1$ . We make this exchangeability assumption within levels of past covariate values (conditional) and at each time point separately (sequential):

$$\begin{aligned} E(Y^{a_0, a_1} \mid A_1, Z_1, A_0) &= E(Y^{a_0, a_1} \mid Z_1, A_0), \text{ and} \\ E(Y^{a_0, a_1} \mid A_0) &= E(Y^{a_0, a_1}). \end{aligned} \tag{2}$$

This sequential conditional exchangeability assumption would



hold if there were no uncontrolled confounding and no selection bias. The top part of equation 2 says that, within levels of prior viral load ( $Z_1$ ) and a given treatment level  $A_0$ ,  $Y^{a_0, a_1}$  does not depend on the assigned values of  $A_1$ . The bottom part of equation 2 says that  $Y^{a_0, a_1}$  does not depend on the assigned values of  $A_0$ . Note the correspondence between these two equations and the causal diagram: because in Figure 1,  $Z_1$  is a common cause of  $A_1$  and  $Y$ , the assumption in equation 2 must be made conditional on  $Z_1$ . Failing to condition for  $Z_1$  will result in uncontrolled confounding of the effect of  $A_1$ , and thus a dependence between the actual  $A_1$  value and the potential outcome. However, adjusting for  $Z_1$  using standard methods (restriction, stratification, matching, or conditioning in a linear regression model) would block part of the effect from  $A_0$  through  $Z_1$ , and potentially lead to a collider bias of the effect of  $A_0$  through  $U$  (Stephen R. Cole et al. 2010) This is the central challenge that g methods were developed to address.

The third assumption, known as positivity (Westreich and Cole 2010) requires  $0 < P(A_1 = 1 \mid Z_1 = z_1, A_0 = a_0) < 1$  and  $0 < P(A_0 = 1) < 1$ . Furthermore, this assumption must hold for all values of  $a_0$  and  $z_1$  where  $P(A_0 = a_0, Z_1 = z_1) > 0$ . This latter condition is required so that effects are not defined in strata of  $a_0$  and  $z_1$  that do not exist. Positivity is met when there are exposed and unexposed individuals within all confounder and prior exposure levels, which can be evaluated empirically.<sup>12</sup>

Under these three assumptions, our hypothetical observational study can be likened to a sequentially randomized trial in which the exposure was randomized at baseline, and randomized again at time 1 with a probability that depends on  $Z_1$ . Under these assumptions, g methods can be used to estimate counterfactual quantities with observational data.

<sup>12</sup> There are actually two types of positivity violations: stochastic and structural. In the former, one need only collect more data to alleviate concerns over stochastic positivity violations. In the latter, certain confounder values preclude the possibility of individuals being exposed or unexposed. One example of the latter is the healthy worker survivor effect.

## RESULTS

### Standard Methods

Table 2 presents results from fitting a number of standard linear regression models to the data in Table 1.

Model Parameters	Estimate ( $\hat{\beta}_1$ )
$\beta_0 + \beta_1(A_0 + A_1)/2$	60.9
$\beta_0 + \beta_1(A_0 + A_1)/2 + \beta_2 Z_1$	42.6
$\beta_0 + \beta_1 A_0$	27.1
$\beta_0 + \beta_1 A_0 + \beta_2 Z_1$	18.0
$\beta_0 + \beta_1 A_1$	38.9
$\beta_0 + \beta_1 A_1 + \beta_2 Z_1$	25.0

Table 2: Linear regression models and corresponding estimates comparing several contrasts quantifying exposed versus unexposed scenarios fit to data in Table 1.

In the first model,  $\hat{\beta} = 60.9$  cells/mm<sup>3</sup> is the crude difference in mean CD4 count for the always treated compared to the never treated. In model two,  $\hat{\beta} = 42.6$  cells/mm<sup>3</sup> is the  $Z_1$ -adjusted difference in mean CD4 count for the same contrast. Other model results are provided in Table 2, and more could be entertained.

Table 3 presents the results from fitting all three g methods to the data in Table 1.

G Method	$\hat{\psi}^a$
G Formula	50.0
IP-weighted marginal structural model	50.0
G Estimated Structural Nested Model	50.0

a  $\psi = E(Y^{1,1} - Y^{0,0})$

Table 3: G-methods and corresponding estimates comparing contrasts quantifying always exposed versus never exposed scenarios fit to data in Table 1.

The marginal structural model resulted in  $\hat{\psi} = 50.0$  cells/mm<sup>3</sup>. The g formula resulted in  $\hat{\psi} = 50.0$  cells/mm<sup>3</sup>. Finally, the structural nested model resulted in  $\hat{\psi} = 50.0$  cells/mm<sup>3</sup>. Next we discuss how we obtained these results.

### *g Methods*

The **g formula** can be used to estimate the average CD4 level that would be observed in the population under a given treatment plan. To implement the approach, we start with a mathematical representation of the data generating mechanism for all variables in Table 1.

We refer to this as the **joint density of the observed data**. We factor the joint density in a way that respects the temporal ordering of the data by conditioning each variable on its history. For example, if  $f(\cdot)$  represents the probability density function, then by the definition of conditional probabilities (Wasserman 2006, p 36) we can factor this joint density as

$$f(y, a_1, z_1, a_0) = f(y \mid a_1, z_1, a_0)P(A_1 = a_1 \mid Z_1 = z_1, A_0 = a_0) \\ P(Z_1 = z_1 \mid A_0 = a_0)P(A_0 = a_0).$$

Our interest lies in the marginal mean of  $Y$  that would be observed if  $A_0$  and  $A_1$  were set to some values  $a_0$  and  $a_1$ , respectively. To obtain this expectation, we perform two mathematical operations on the factored joint density. The first is the well-known expectation operator (Wasserman 2006, p 47), which allows us to write the conditional mean of  $Y$  in terms of its conditional density. The second is the law of total probability (Wasserman 2006, p 12), which allows us to marginalize over the distribution of  $A_1$ ,  $Z_1$  and  $A_0$ , yielding the marginal mean of  $Y$ :

$$E(Y) = \sum_{a_1, z_1, a_0} E(Y \mid A_1 = a_1, Z_1 = z_1, A_0 = a_0)P(A_1 = a_1 \mid Z_1 = z_1, A_0 = a_0) \\ P(Z_1 = z_1 \mid A_0 = a_0)P(A_0 = a_0).$$

We can now modify this equation to yield the average of potential outcomes that would be observed after intervening on the exposure [enabling us to drop out the terms for  $P(A_1 = a_1 \mid Z_1 = z_1, A_0 = a_0)$  and  $P(A_0 = a_0)$ ], yielding

$$E(Y^{a_0, a_1}) = \sum_{z_1} E(Y \mid A_1 = a_1, Z_1 = z_1, A_0 = a_0)P(Z_1 = z_1 \mid A_0 = a_0).$$

**This equation is the g formula.** Its proof, given in the Supplementary Material of Naimi et al (2017), follows from the three identifying assumptions. **In our simple scenario, the expectation  $E(Y^{0,0})$  can be calculated by summing the mean CD4 count in the never treated with  $Z_1 = 1$  (weighted by the proportion of people with  $Z_1 = 1$  in the  $A_0 = 0$  stratum) and the mean CD4 count in the never treated with  $Z_1 = 0$  (weighted by the proportion of people with  $Z_1 = 0$  in the**

$A_0 = 0$  stratum). Weighting the observed outcome's conditional expectation by the conditional probability that  $Z_1 = z_1$  enables us to account for the fact that  $Z_1$  is affected by  $A_0$ , but also confounds the effect of  $A_1$  on  $Y$ . Computing this expectation's value yields a result of  $\hat{E}(Y^{0,0}) = 100.0$ , where we use  $\hat{E}$  to denote a sample, rather than a population average, and with the understanding that  $\hat{E}(Y^{0,0})$  is equal to the g formula with  $A_0 = A_1 = 0$  (since the potential outcomes  $Y^{0,0}$  are not directly observed). We repeat the process to obtain the corresponding value for treated at time 0 only:  $\hat{E}(Y^{1,0}) = 125.0$ ; treated at time 1 only:  $\hat{E}(Y^{0,1}) = 125.0$ ; and always treated:  $\hat{E}(Y^{1,1}) = 150.0$ . Thus,  $\hat{\psi}_{GF} = 150.0 - 100.0 = 50.0$ , which is the average causal effect of treatment on CD4 cell count.

This approach to computing the value of the g formula is referred to as nonparametric maximum likelihood estimation. Several authors Edwards et al. (2014) demonstrate how simulation from parametric regression models can yield a g formula estimator, which is often required in typical population-health studies with many covariates.

Modeling each component of the joint density of the observed data (including the probability that  $Z_1 = z_1$ ) can lead to bias if any of these models are mis-specified.<sup>13</sup> To compute the expectations of interest, we can instead specify a single model that targets our average causal effect, and avoid unnecessary modeling. Marginal structural models with IP weighting map a *marginal summary* (e.g., average) of potential outcomes to the treatment and parameter of interest  $\psi$ . Unlike the g formula, they do not require a model for  $P(Z_1 = z_1 | A_0 = a_0)$ . Additionally, as we show in the Supplementary Material of Naimi et al (2017), while they cannot model it directly, they are indifferent to whether time-varying effect modification is present or absent. Because our interest lies in the marginal contrast of outcomes under always versus never treated conditions, our marginal structural model for the effect of  $A$  can be written as  $E(Y^{a_0, a_1}) = \beta_0 + \psi_0 a_0 + \psi_1 a_1 + \psi_2 a_0 a_1$ , where  $\beta_0 = E(Y^{0,0})$  is a (nuisance) intercept parameter, and  $\psi = E(Y^{1,1} - Y^{0,0}) = (\psi_0 + \psi_1 + \psi_2)$  is the effect of interest.

<sup>13</sup> One of the major limitations of the parametric g formula.

Inverse probability weighting can be used estimate marginal struc-

tural model parameters (proofs are provided in the Supplementary Material). To estimate  $\psi$  using inverse probability weighted regression, we first obtain the predicted probabilities of the observed treatments. In our example data, there are two possible  $A_1$  values (exposed, unexposed) for each of the four levels in  $Z_1$  and  $A_0$ . Additionally, there are two possible  $A_0$  values (exposed, unexposed) overall. This leads to four possible exposure regimes: never treat, treat early only, treat late only, and always treat. For each  $Z_1$  value, we require the predicted probability of the exposure that was actually received. These probabilities are computed by calculating the appropriate proportions of subjects in Table 1. Because there are no variables that affect  $A_0$ , this probability is 0.5 for all individuals in the sample. Furthermore, in our example  $A_1$  is not affected by  $A_0$  (Figure 1). Thus, the  $Z_1$  specific probabilities of  $A_1$  are constant across levels of  $A_0$ . In settings where  $A_0$  affects  $A_1$ , the  $Z_1$  specific probabilities of  $A_1$  would vary across levels of  $A_0$ .

In the stratum defined by  $Z_1 = 1$ , the predicted probabilities of  $A_1 = 0$  and  $A_1 = 1$  are 0.308 and 0.692, respectively. For example,  $\$(210,527+136,293) / (210,527+136,293+93,903+60,654) = 0.692 \$$ . Thus, the probabilities for each treatment combination are:  $0.5 \times 0.308 = 0.155$  (never treated),  $0.5 \times 0.308 = 0.155$  (treated early only),  $0.5 \times 0.692 = 0.346$  (treated late only), and  $0.5 \times 0.692 = 0.346$  (always treated). Dividing the marginal probability of each exposure category (not stratified by  $Z_1$ ) by these stratum specific probabilities gives stabilized weights of 1.617, 1.617, 0.725, and 0.725, respectively. For example, the never treated weight is  $(0.5 \times 0.501) / (0.5 \times 0.308) = 1.617$ . The same approach is taken to obtain predicted probabilities and stabilized weights in the stratum defined by  $Z_1 = 0$ . The weights and weighted data are provided in Table 4.

Fitting this model in the weighted data given in Table 4 provides the inverse-probability weighted estimates  $[\hat{\psi}_{0_{IP}} = 25.0, \hat{\psi}_{1_{IP}} = 25.0, \hat{\psi}_{2_{IP}} = 0.0]$ , thus yielding  $\hat{\psi}_{IP} = 50.0$ .

Weighting the observed data by the inverse of the probability of the observed exposure yields a “pseudo-population” (Table 4) in

$A_0$	$Z_1$	$A_1$	$Y$	$sw$	Pseudo $N$
0	0	0	87.23	0.72	151222.84
0	0	1	112.23	1.62	151680.46
0	1	0	119.79	1.62	98110.06
0	1	1	144.78	0.72	98789.4
1	0	0	105.25	0.72	97395.08
1	0	1	130.25	1.62	98321.62
1	1	0	137.8	1.62	151884.02
1	1	1	162.8	0.72	152596.51

Table 4: Pseudo-population obtained after applying inverse probability weights to data in Table 1.

which treatment at the second time point ( $A_1$ ) is no longer related to (and is thus no longer confounded by) viral load just prior to the second time point ( $Z_1$ ). Thus, weighting a conditional regression model for the outcome by the inverse probability of treatment enables us to account for the fact that  $Z_1$  both confounds  $A_1$  and is affected by  $A_0$ .

Structural nested models map a *conditional contrast* of potential outcomes to the treatment, within nested sub-groups of individuals defined by levels of  $A_1$ ,  $Z_1$ , and  $A_0$ . Our structural nested model can be written as

$$\begin{aligned}
 E(Y^{a_0, a_1} - Y^{a_0, 0} \mid A_0 = a_0, Z_1 = z_1, A_1 = a_1) &= a_1(\psi_1 + \psi_2 a_0 + \psi_3 z_1 + \psi_4 a_0 z_1) \\
 E(Y^{a_0, 0} - Y^{0, 0} \mid A_0 = a_0) &= \psi_0 a_0
 \end{aligned} \tag{3}$$

Note this model introduces two additional parameters:  $\psi_3$  for the two-way interaction between  $a_1$  and  $z_1$ , and  $\psi_4$  for the three-way interaction between  $a_1$ ,  $z_1$ , and  $a_0$ . Indeed, the ability to explicitly quantify interactions between time-varying exposures and time-varying covariates (which cannot be modeled via standard marginal structural models) is a major strength of structural nested models when effect modification is of interest.<sup>14</sup> To simplify our exposition, we set  $(\psi_3, \psi_4) = (0, 0)$  in our data example, allowing us to drop the  $\psi_3 z_1$  and  $\psi_4 a_0 z_1$  terms from the model. In effect, this renders our structural nested mean model equivalent to a semi-parametric marginal structural model. In the Supplementary Material, we explain how marginal structural and structural nested models each relate to time-varying interactions in more detail.

We can now use g-estimation to estimate  $(\psi_0, \psi_1, \psi_2)$  in the above

<sup>14</sup>

structural nested model. G-estimation is based on solving equations that directly result from the sequential conditional exchangeability assumptions in (2) and (??), combined with assumptions implied by the structural nested model. If, at each time point, the exposure is conditionally independent of the potential outcomes (sequential exchangeability) then the conditional covariance between the exposure and potential outcomes is zero.<sup>15</sup> Formally, these conditional independence relations can be written as:

$$\begin{aligned} 0 &= \text{Cov}(Y^{a_0,0}, A_1 \mid Z_1, A_0) \\ &= \text{Cov}(Y^{0,0}, A_0) \end{aligned} \tag{4}$$

where  $\text{Cov}(\cdot)$  is the well-known covariance formula.<sup>16(p52)</sup> These equalities are of little direct use for estimation, though, as they contain unobserved potential outcomes and are not yet functions of the parameters of interest. However, by counterfactual consistency and the structural nested model, we can replace these unknowns with quantities estimable from the data.

Specifically, as we prove in the Supplementary Material, the structural nested model, together with exchangeability and counterfactual consistency imply that we can replace the potential outcomes  $Y^{a_0,0}$  and  $Y^{0,0}$  in the above covariance formulas with their values implied by the structural nested model, yielding:

$$\begin{aligned} 0 &= \text{Cov}\{Y - A_1(\psi_1 + \psi_2 A_0), A_1 \mid Z_1, A_0\} \\ &= \text{Cov}\{Y - A_1(\psi_1 + \psi_2 A_0) - \psi_0 A_0, A_0\}. \end{aligned} \tag{5}$$

We provide an intuitive explanation for this substitution in the Supplementary Material. %is that it would certainly hold under a stronger version of our structural nested model assumptions, in which  $Y^{a_0,a_1} - Y^{a_0,0} = a_1(\psi_1 + \psi_2 a_0)$  and  $Y^{a_0,0} - Y^{0,0} = \psi_0 a_0$  exactly, so that  $Y^{A_0,0} = Y - A_1(\psi_1 + \psi_2 A_0)$  and  $Y^{0,0} = Y - A_1(\psi_1 + \psi_2 A_0) - \psi_0 A_0$ . We also show how these covariance relations yield three equations that can be used to solve each of the unknowns in the above structural nested model  $(\psi_0, \psi_1, \psi_2)$ .

Two of the three equations yield the following g estimators:

$$\begin{aligned}\hat{\psi}_{1GE} &= \frac{\hat{E}[(1 - A_0)Y\{A_1 - \hat{E}(A_1 | Z_1, A_0)\}]}{\hat{E}[(1 - A_0)A_1\{A_1 - \hat{E}(A_1 | Z_1, A_0)\}]} \\ \hat{\psi}_{1GE} + \hat{\psi}_{2GE} &= \frac{\hat{E}[A_0Y\{A_1 - \hat{E}(A_1 | Z_1, A_0)\}]}{\hat{E}[A_0A_1\{A_1 - \hat{E}(A_1 | Z_1, A_0)\}]} \end{aligned} \quad (6)$$

Note that to solve these equations we need to model  $E(A_1 | Z_1, A_0)$ , which in practice we might assume can be correctly specified as the predicted values from a logistic model for  $A_1$ . In our simple setting, the correctness of this model is guaranteed by saturating it (i.e., conditioning the model on  $Z_1, A_0$  and their interaction).

As we show in the Supplementary Material, implementing these equations in software can be easily done using either an instrumental variables (i.e., two-stage least squares) estimator, or ordinary least squares. %estimator that regresses  $Y$  on  $A_1$  using as an “instrument” the residual  $\{A_1 - \hat{E}(A_1 | Z_1, A_0)\}$ , where the first estimator is a two-stage least squares regression among the initially untreated with  $A_0 = 0$  and the second among the initially treated with  $A_0 = 1$ .

Once the above parameters are estimated, the next step is to subtract the effect of  $A_1$  and  $A_1A_0$  from  $Y$  to obtain  $\tilde{Y} = Y - \hat{\psi}_{1GE}A_1 - \hat{\psi}_{2GE}A_1A_0$ . We can then solve for the last parameter using a sample version of the third g estimation equality, yielding our final estimator and completing the procedure:

$$\hat{\psi}_{0GE} = \frac{\hat{E}[\tilde{Y}\{A_0 - \hat{E}(A_0)\}]}{\hat{E}[A_0\{A_0 - \hat{E}(A_0)\}]}.$$

Again the above estimator can be implemented using an instrumental variable or ordinary least squares estimator. Implementing this procedure in our example data, we obtain  $[\psi_{0GE} = 25.0, \psi_{1GE} = 25.0, \psi_{2GE} = 0.0]$ , thus yielding  $\psi_{GE} = 50.0$ .

The potential outcome under no treatment can be thought of as a given subject’s baseline prognosis: in our setting, individuals with poor baseline prognosis will have low CD4 levels, no matter what their treatment status may be. In the absence of confounding or selection bias, one expects this baseline prognosis to be independent of treatment status. G estimation exploits this independence by as-



suming no uncontrolled confounding (conditional on measured confounders), and assigning values to  $\hat{\psi}_{GE}$  that render the potential outcomes independent of the exposure. However, assigning the correct values to  $\hat{\psi}_{GE}$  depends on there being no confounding or selection bias.

## DISCUSSION

Having constructed these data using the causal diagram shown in Figure 1, we know the true effect of combined treatment is indeed 50 cells/mm<sup>3</sup> (25 cells/mm<sup>3</sup> for each exposure main effect) as well approximated by all three g methods, but not by any of the standard regression models we fit, with one exception. The final standard result presented in Table 2 correctly estimates the effect of the second treatment (an effect of 25 cells/mm<sup>3</sup>), as would be expected from the causal diagram.

For the past several years, we have used the foregoing simple example to initiate epidemiologists to g methods with some success. Once having studied this simple example in detail, we recommend working through more comprehensive examples by Robins and Hern{a}n<sup>17</sup> and Hern{a}n and Robins.<sup>18</sup> A recent tutorial<sup>19</sup> may then be of further use. G methods are becoming more common in epidemiologic research.<sup>20</sup> We hope this commentary facilitates the process of better understanding these useful methods.

17  
18  
19

20

## References

- Cole, S. R., and C. E. Frangakis. 2009. "The Consistency Statement in Causal Inference: A Definition or an Assumption?" *Epidemiol* 20 (1): 3–5.
- Cole, Stephen R., David B. Richardson, Haitao Chu, and Ashley I. Naimi. 2013. "Analysis of Occupational Asbestos Exposure and Lung Cancer Mortality Using the g Formula." *Am J Epidemiol* 177 (9): 989–96.
- Cole, Stephen R, Robert W Platt, Enrique F Schisterman, Haitao Chu, Daniel Westreich, David Richardson, and Charles Poole.

2010. "Illustrating Bias Due to Conditioning on a Collider." *Int J Epidemiol* 39 (2): 417–20.
- Daniel, R. M., S. N. Cousens, B. L. De Stavola, M. G. Kenward, and J. A. C. Sterne. 2013. "Methods for Dealing with Time-Dependent Confounding." *Stat Med* 32 (9): 1584–1618.
- Edwards, Jessie K, LJ McGrath, Buckley JP, MK Schubauer-Berigan, SR Cole, and Richardson DB. 2014. "Occupational Radon Exposure and Lung Cancer Mortality: Estimating Intervention Effects Using the Parametric g-Formula." *Epidemiol* 25 (6): 829–34.
- Greenland, Sander, and JM Robins. 1986. "Identifiability, Exchangeability, and Epidemiological Confounding." *Int J Epidemiol* 15 (3): 413–19.
- Hernán, M. A., and JM Robins. Forthcoming 2021. *Causal Inference: What If*. Boca Raton, FL: Chapman/Hall.
- Imbens, Guido W, and Donald B Rubin. 2015. *Causal Inference for Statistics, Social, and Biomedical Sciences: An Introduction*. New York, NY: Cambridge University Press.
- Keil, Alex, Jessie K Edwards, David B. Richardson, Ashley I. Naimi, and Stephen R. Cole. 2014. "The Parametric g-Formula for Time-to-Event Data: Towards Intuition with a Worked Example." *Epidemiol* 25 (6): 889–97.
- Naimi, Ashley I., and Daniel J. Westreich. 2014. "Big Data: A Revolution That Will Transform How We Live, Work, and Think." *American Journal of Epidemiology* 179 (9): 1143–44.
- Naimi, Ashley I, Stephen R Cole, and Edward H Kennedy. 2017. "An Introduction to G Methods." *Int J Epidemiol* 46 (2): 756–62.
- Pearl, Judea, Madelyn R Glymour, and Nicholas Jewell. 2016. *Causal Inference in Statistics: A Primer*. United Kingdom: Wiley.
- Pearson, Karl. 1911. *The Grammar of Science*. 3rd ed. London, J.M. Dent & sons ltd.
- Robins, James M, and Miguel Á Hernán. 2009. "Estimation of the Causal Effects of Time-Varying Exposures." In *Advances in Longitudinal Data Analysis*, edited by G Fitzmaurice, M. Davidian, G Verbeke, and G Molenberghs, 553–99. Boca Raton, FL: Chapman & Hall.

- Robins, JM. 1987. "Addendum to "a New Approach to Causal Inference in Mortality Studies with a Sustained Exposure Period—Application to Control of the Healthy Worker Survivor Effect"." *Comp Math Appl* 14 (9-12): 923–45.
- Taubman, S. L., J. M. Robins, M. A. Mittleman, and M. A. Hernán. 2009. "Intervening on Risk Factors for Coronary Heart Disease: An Application of the Parametric g-Formula." *Int J Epidemiol* 38 (6): 1599–1611.
- VanderWeele, Tyler J, and Miguel Ángel Hernán. 2013. "Causal Inference Under Multiple Versions of Treatment." *Journal of Causal Inference* 1 (1): 1–20.
- Wasserman, Larry. 2006. *All of Nonparametric Statistics*. New York; London: Springer.
- Westreich, Daniel, and Stephen R. Cole. 2010. "Invited Commentary: Positivity in Practice." *Am J Epidemiol* 171 (6): 674–77.
- Westreich, Daniel, Stephen R. Cole, Jessica G. Young, Frank Palella, Phyllis C. Tien, Lawrence Kingsley, Stephen J. Gange, and Miguel A. Hernán. 2012. "The Parametric g-Formula to Estimate the Effect of Highly Active Antiretroviral Therapy on Incident AIDS or Death." *Stat Med* 31 (18): 2000–2009.