

Causal Inference

Miguel A. Hernán, James M. Robins

May 19, 2017

Part III

Causal inference from complex longitudinal data

Chapter 19

TIME-VARYING TREATMENTS

So far this book has dealt with fixed treatments which do not vary over time. However, many causal questions involve treatments that vary over time. For example, we may be interested in estimating the causal effects of medical treatments, lifestyle habits, employment status, marital status, occupational exposures, etc. Because these treatments may take different values for a single individual over time, we refer to them as time-varying treatments.

Restricting our attention to time-fixed treatments during Parts I and II of this book helped us introduce basic concepts and methods. It is now time to consider more realistic causal questions that involve the contrast of hypothetical interventions that are played out over time. Part III extends the material in Parts I and II to time-varying treatments. This chapter describes some key terminology and concepts for causal inference with time-varying treatments. Though we have done our best to simplify those concepts (if you don't believe us, check out the causal inference literature), this is still one of the most technical chapters in the book. Unfortunately, further simplification would result in too much loss of rigor. But if you made it this far, you are qualified to understand this chapter.

19.1 The causal effect of time-varying treatments

Consider a time-fixed treatment variable A (1: treated, 0: untreated) at time zero of follow-up and an outcome variable Y measured 60 months later. We have previously defined the average causal effect of A on the outcome Y as the contrast between the mean counterfactual outcome $Y^{a=1}$ under treatment and the mean counterfactual outcome $Y^{a=0}$ under no treatment, that is, $E[Y^{a=1}] - E[Y^{a=0}]$. Because treatment status is determined at a single time (time zero) for everybody, the average causal effect does not need to make reference to the time at which treatment occurs. In contrast, causal contrasts that involve time-varying treatments need to incorporate time explicitly.

For simplicity, we will provisionally assume that no individuals were lost to follow-up or died during this period, and we will also assume that all variables were perfectly measured.

To see this, consider a time-varying dichotomous treatment A_k that may change at every month k of follow-up, where $k = 0, 1, 2, \dots, K$ with $K = 59$. For example, in a 5-year follow-up study of individuals infected with the human immunodeficiency virus (HIV), A_k takes value 1 if the individual receives antiretroviral therapy in month k , and 0 otherwise. No individuals received treatment before the start of the study at time 0, i.e., $A_{-1} = 0$ for all individuals.

For compatibility with many published papers, we use zero-based indexing for time. That is, the first time of possible treatment is $k = 0$ rather than $k = 1$.

We use an overbar to denote treatment history, that is, $\bar{A}_k = (A_0, A_1, \dots, A_k)$ is the history of treatment from time 0 to time k . When we refer to the entire treatment history through K , we often represent \bar{A}_K as \bar{A} without a time subscript. In our HIV study, an individual who receives treatment continuously throughout the follow-up has treatment history $\bar{A} = (A_0 = 1, A_1 = 1, \dots, A_{59} = 1) = (1, 1, \dots, 1)$, or $\bar{A} = \bar{1}$. Analogously, an individual who never receives treatment during the follow-up has treatment history $\bar{A} = (0, 0, \dots, 0) = \bar{0}$. Most individuals are treated during part of the follow-up only, and therefore have intermediate treatment histories with some 1s and some 0s—which we cannot represent as compactly as $\bar{1}$ and $\bar{0}$.

To keep things simple, our example considers an outcome measured at a fixed time. However, the concepts discussed in this chapter also apply to time-varying outcomes and failure time outcomes.

Remember: we use lower-case to denote possible realizations of a random variable, e.g., a_k is a realization of treatment A_k .

Suppose Y measures health status—with higher values of Y indicating better health—at the end of follow-up at time $K + 1 = 60$. We would like to estimate the average causal effect of the time-varying treatment \bar{A} on the outcome Y . But we can no longer define the average causal effect of a time-varying treatment as a contrast at a single time k , because the contrast $E[Y^{a_k=1}] - E[Y^{a_k=0}]$ quantifies the effect of treatment A_k at a single time k , not the effect of the time-varying treatment A_k at all times k between 0 and 59.

Indeed we will have to define the average causal effect as a contrast between the counterfactual mean outcomes under two treatment strategies that involve treatment at all times between the start ($k = 0$) and the end ($k = K$) of the follow-up. As a consequence, the average causal effect of a time-varying treatment is not uniquely defined. In the next section, we describe many possible definitions of average causal effect for a time-varying treatment.

19.2 Treatment strategies

A general counterfactual theory to compare treatment strategies was first articulated by Robins (1986, 1987, 1997a).

A treatment strategy—also referred to as a plan, policy, protocol, or regime—is a rule to assign treatment at each time k of follow-up. For example, two treatment strategies are “always treat” and “never treat” during the follow-up. The strategy “always treat” is represented by $\bar{a} = (1, 1, \dots, 1) = \bar{1}$, and the strategy “never treat” is represented by $\bar{a} = (0, 0, \dots, 0) = \bar{0}$. We can now define an average causal effect of \bar{A} on the outcome Y as the contrast between the mean counterfactual outcome $Y^{\bar{a}=\bar{1}}$ under the strategy “always treat” and the mean counterfactual outcome $Y^{\bar{a}=\bar{0}}$ under the strategy “never treat”, that is, $E[Y^{\bar{a}=\bar{1}}] - E[Y^{\bar{a}=\bar{0}}]$.

But there are many other possible causal effects for the time-varying treatment \bar{A} , each of them defined by a contrast of outcomes under two particular treatment strategies. For example, we might be interested in the average causal effect defined by the contrast $E[Y^{\bar{a}}] - E[Y^{\bar{a}'}]$ that compares the strategy “treat at every other month” $\bar{a} = (1, 0, 1, 0, \dots)$ with the strategy “treat at all months except the first one” $\bar{a}' = (0, 1, 1, 1, \dots)$. The number of possible contrasts is very large: we can define at least 2^K treatment strategies because there are 2^K possible combinations of values (a_0, a_1, \dots, a_K) for a dichotomous a_k . In fact, as we next explain, these 2^K such strategies do not exhaust all possible treatment strategies.

To define even more treatment strategies in our HIV example, consider the time-varying covariate L_k which denotes CD4 cell count (in cells/ μ L) measured at month k in all individuals. The variable L_k takes value 1 when the CD4 cell count is low, which indicates a bad prognosis, and 0 otherwise. At time zero, all individuals have a high CD4 cell count, $L_0 = 0$. We could then consider the strategy “do no treat while $L_k = 0$, start treatment when $L_k = 1$ and treat continuously after that time”. This treatment strategy is different from the ones considered in the previous paragraph because we cannot represent it by a rule $\bar{a} = (a_0, a_1, \dots, a_K)$ under which all individuals get the same treatment a_0 at time $k = 0$, a_1 at time $k = 1$, etc. Now, at each time, some individuals will be treated and others will be untreated, depending of the value of their evolving L_k . This is an example of a *dynamic treatment strategy*, a rule in which the treatment a_k at time k depends on the evolution of an individual’s time-varying covariate(s) \bar{L}_k . Strategies \bar{a} for which treatment does not depend

Fine Point 19.1

Deterministic and random treatment strategies. A dynamic treatment strategy is a rule $g = [g_0(\bar{a}_{-1}, \bar{l}_0), \dots, g_K(\bar{a}_{K-1}, \bar{l}_K)]$, where $g_k(\bar{a}_{k-1}, \bar{l}_k)$ specifies the treatment assigned at k to an individual with past history $(\bar{a}_{k-1}, \bar{l}_k)$. An example in our HIV study: $g_k(\bar{a}_{k-1}, \bar{l}_k)$ is 1 if an individual's CD4 cell count (a function of \bar{l}_k) was low at or before k ; otherwise $g_k(\bar{a}_{k-1}, \bar{l}_k)$ is 0. A static treatment strategy is a rule $g = [g_0(\bar{a}_{-1}), \dots, g_K(\bar{a}_{K-1})]$, where $g_k(\bar{a}_{k-1})$ does not depend on \bar{l}_k . We will often abbreviate $g_k(\bar{a}_{k-1}, \bar{l}_k)$ as $g(\bar{a}_{k-1}, \bar{l}_k)$.

Most static and dynamic strategies we are interested in comparing are *deterministic treatment strategies*, which assign a particular value of treatment (0 or 1) to each individual at each time. More generally, we could consider *random treatment strategies* that do not assign a particular value of treatment, but rather a probability of receiving a treatment value. Random treatment strategies can be static (e.g., “independently at each month, treat a subject with probability 0.3 and do not treat with probability 0.7”) or dynamic (e.g., “independently at each month treat subjects whose CD4 cell count is low with probability 0.3, but do not treat subjects with high CD4 cell count”).

We refer to the strategy g for which the mean counterfactual outcome $E[Y^g]$ is maximized (when higher values of outcome are better) as the optimal treatment strategy. For a drug treatment, the optimal strategy will be almost always be dynamic because treatment needs to be discontinued when toxicity develops. Also, no random strategy can ever be preferred to the optimal deterministic strategy. However, random strategies (i.e., ordinary randomized trials and sequentially randomized trials) remain scientifically necessary because, before the trial, it is unknown which deterministic regime is optimal. See Young et al. (2014) for a taxonomy of treatment strategies. In the text, except if noted otherwise, the letter g will refer only to deterministic treatment strategies.

on covariates are non-dynamic or *static treatment strategies*. See Fine Point 19.1 for a formal definition.

Causal inference with time-varying treatments involves the contrast of counterfactual outcomes under two or more treatment strategies. The average causal effect of a time-varying treatment is only well-defined if the treatment strategies of interest are specified. In our HIV example, we can define an average causal effect based on the difference $E[Y^{\bar{a}}] - E[Y^{\bar{a}'}]$ that contrasts strategy \bar{a} (say, “always treat”) versus strategy \bar{a}' (say, “never treat”), or on the difference $E[Y^{\bar{a}}] - E[Y^g]$ that contrasts strategy \bar{a} (“always treat”) versus strategy g (say, “treat only after CD4 cell count is low”). Note we will often use g to represent any—static or dynamic—strategy. When we use it to represent a static strategy, we sometimes write $Y^{g=\bar{a}}$ rather than just Y^g or $Y^{\bar{a}}$.

That is, there is not a single definition of causal effect for time-varying treatments. Even when only two treatment options—treat or do not treat—exist at each time k , we can still define as many causal effects as pairs of treatment strategies exist. In the next section, we describe a study design under which all these causal effects can be validly estimated: the sequentially randomized experiment.

19.3 Sequentially randomized experiments

Recall that, by definition, a causal graph must always include all common causes of any two variables on the graph.

The causal diagrams in Figures 19.1, 19.2, and 19.3 summarize three situations that can occur in studies with time-varying treatments. In all three diagrams, A_k represents the time-varying treatment, L_k the set of measured variables, Y the outcome, and U_k the set of unmeasured variables at k that are common causes of at least two other variables on the causal graph. Because the covariates U_k are not measured, their values are unknown and therefore un-

Technical Point 19.1

On the definition of dynamic strategies. Each dynamic strategy $g = [g_0(\bar{a}_{-1}, \bar{l}_0), \dots, g_K(\bar{a}_{K-1}, \bar{l}_K)]$ that depends on past treatment and covariate history is associated with a dynamic strategy $g' = [g'_0(\bar{l}_0), \dots, g'_K(\bar{l}_K)]$ that depends only on past covariate history. By consistency (see Technical Point 19.2), an individual will have the same treatment, covariate, and outcome history when following strategy g from time zero as when following strategy g' from time zero. In particular, $Y_g = Y_{g'}$ and $\bar{L}_g(K) = \bar{L}_{g'}(K)$. Specifically, g' is defined in terms of g recursively by $g'_0(l_0) = g_0(\bar{a}_{-1} = 0, l_0)$ (by convention, \bar{a}_{-1} can only take the value zero) and $g'_k(\bar{l}_k) = g_k[g'_k(\bar{l}_{k-1}), \bar{l}_k]$. For any strategy g for which treatment at each k already does not depend on past treatment history, g and g' are the identical set of functions. The above definition of g' in terms of g guarantees that an individual has followed strategy g through time t in the observed data, i.e., $A_k = g_k(\bar{A}_{k-1}, \bar{L}_k)$ for $k \leq t$, if and only if the individual has followed strategy g' through t , i.e., $A_k = g'_k(\bar{L}_k)$ for $k \leq t$.

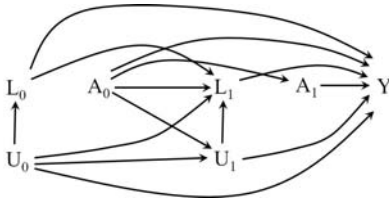


Figure 19.1

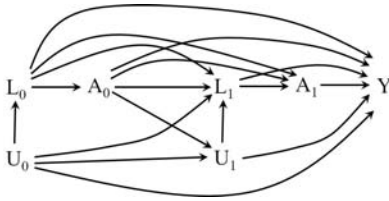


Figure 19.2

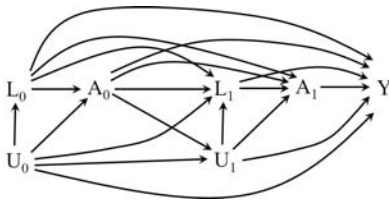


Figure 19.3

available for the analysis. In our HIV study, the time-varying covariate CD4 cell count L_k is a consequence of the true, but unmeasured, chronic damage to the immune system U_k . The greater an individual's immune damage U_k , the lower her CD4 cell count L_k and her health status Y . For simplicity, the causal diagrams include only the first two times of follow-up $k = 0$ and $k = 1$, and we will assume that all participants adhered to the assigned treatment (see Fine Point 19.2).

The causal diagram in Figure 19.1 lacks arrows from either the measured covariates \bar{L}_k or the unmeasured covariates \bar{U}_k into treatment A_k . The causal diagram in Figure 19.2 has arrows from the measured covariates \bar{L}_k , but not from the unmeasured covariates \bar{U}_k , into treatment A_k . The causal diagram in Figure 19.3 has arrows from both the measured covariates \bar{L}_k and the unmeasured covariates \bar{U}_k into treatment A_k .

Figure 19.1 could represent a randomized experiment in which treatment A_k at each time k is randomly assigned with a probability that depends only on prior treatment history. Our HIV study would be represented by Figure 19.1 if, for example, an individual's treatment value at each month k were randomly assigned with probability 0.5 for individuals who did not receive treatment during the previous month ($A_{k-1} = 0$), and with probability 1 for individuals who did receive treatment during the previous month ($A_{k-1} = 1$). When interested in the contrast of static treatment strategies, Figure 19.1 is the proper generalization of no confounding by measured or unmeasured variables for time-varying treatments. Under this causal diagram, the counterfactual outcome mean $E[Y^{\bar{a}}]$ if everybody had followed the static treatment strategy \bar{a} is simply the mean outcome $E[Y|\bar{A} = \bar{a}]$ among those who followed the strategy \bar{a} . (Interestingly, the same is not true for dynamic strategies. The counterfactual mean $E[Y^g]$ under a dynamic strategy g that depends on the variables L is only the mean outcome among those who followed the strategy g if the probability of receiving treatment $A_k = 1$ is exactly 0.5 at all times k at which treatment A_k depends on \bar{L}_k . Otherwise, identifying $E[Y^g]$ requires the application of g-methods to data on \bar{L} , \bar{A} , and Y under either Figure 19.1 or Figure 19.2.)

Figure 19.2 could represent a randomized experiment in which treatment A_k at each time k is randomly assigned by the investigators with a probability that depends on prior treatment *and* measured covariate history. Our study would be represented by Figure 19.2 if, for example, an individual's treatment value at each month k were randomly assigned with probability 0.4 for untreated

Fine Point 19.2

Per-protocol effects to compare treatment strategies. Many randomized trials assign individuals to a treatment at baseline with the intention that they will keep taking it during the follow-up, unless the treatment becomes toxic. That is, the protocol of the trial implicitly or explicitly aims at the comparison of dynamic treatment strategies. For example, the goal of a trial of statin therapy among healthy individuals may be the comparison of the dynamic strategies “initiate statin therapy at baseline and keep taking it during the study unless rhabdomyolysis occurs” versus “do not take statin therapy during the study unless LDL-cholesterol is high or coronary heart disease is diagnosed.”

Estimating the per-protocol effect—the effect that would have been observed if everybody had adhered to their assigned strategy—in this randomized trial raises the same issues as any comparison of treatment strategies in an observational study. Specifically, valid estimation of the per-protocol effect generally demands that trial investigators collect data on, and adjust for, time-varying (post-randomization) variables. Baseline randomization makes us expect baseline exchangeability for the assigned treatment, not sequential exchangeability for the received treatment.

individuals with high CD4 cell count ($A_{k-1} = 0, L_k = 1$), 0.8 for untreated individuals with low CD4 cell count ($A_{k-1} = 0, L_k = 0$), and 0.5 for previously treated individuals, regardless of their CD4 cell count ($A_{k-1} = 1$). In Figure 19.2, there is confounding by measured, but not unmeasured, variables for the time-varying treatment.

An experiment in which treatment is randomly assigned at each time k to each individual is referred to as a *sequentially randomized experiment*. Therefore Figures 19.1 and 19.2 could represent sequentially randomized experiments. On the other hand, Figure 19.3 cannot represent a randomized experiment: the value of treatment A_k at each time k depends partly on unmeasured causes of L_k and Y ; but unmeasured variables obviously cannot be used by investigators to assign treatment. That is, a sequentially randomized experiment can be represented by a causal diagram with many time points $k = 0, 1 \dots K$ and with no direct arrows from the unmeasured prognostic factors U into treatment A_k at any time k .

In observational studies, decisions about treatment often depend on outcome predictors such as prognostic factors. Therefore, observational studies will be typically represented by either Figure 19.2 or Figure 19.3 rather than Figure 19.1. For example, suppose our HIV follow-up study were an observational study (not an experiment) in which the lower the CD4 cell count L_k , the more likely a patient is to be treated. Then our study would be represented by Figure 19.2 if, at each month k , treatment decisions in the real world were made based on the values of prior treatment and CD4 cell count history (\bar{A}_{k-1}, \bar{L}_k), but not on the values of any unmeasured variables \bar{U}_k . Thus, an observational study represented by Figure 19.2 would differ from a sequentially randomized experiment only in that the assignment probabilities are unknown (but could be estimated from the data). Unfortunately, it is impossible to show empirically whether an observational study is represented by the causal diagram in either Figure 19.2 or Figure 19.3. Observational studies represented by Figure 19.3 have unmeasured confounding, as we describe later.

Sequentially randomized experiments are not frequently used in practice. However, the concept of sequentially randomized experiment is helpful to understand some key conditions for valid estimation of causal effects of time-varying treatments. The next section presents these conditions formally.

19.4 Sequential exchangeability

As described in Parts I and II, valid causal inferences about time-fixed treatments typically require conditional exchangeability $Y^a \perp\!\!\!\perp A|L$. When exchangeability $Y^a \perp\!\!\!\perp A|L$ holds, we can obtain unbiased estimates of the causal effect of treatment A on the outcome Y if we appropriately adjust for the variables in L via standardization, IP weighting, g-estimation, or other methods. We expect conditional exchangeability to hold in conditionally randomized experiments—a trial in which individuals are assigned treatment with a probability that depends on the values of the covariates L . Conditional exchangeability holds in observational studies if the probability of receiving treatment depends on the measured covariates L and, conditional on L , does not further depend on any unmeasured, common causes of treatment and outcome.

Similarly, causal inference with time-varying treatments requires adjusting for the time-varying covariates \bar{L}_k to achieve conditional exchangeability at each time point, that is, sequential conditional exchangeability. For example, in a study with two time points, sequential conditional exchangeability is the combination of conditional exchangeability at both the first time and the second time of the study. That is, $Y^g \perp\!\!\!\perp A_0|L_0$ and $Y^g \perp\!\!\!\perp A_1|A_0 = g(L_0), L_0, L_1$. (For brevity, in this book we drop the word “conditional” and simply say sequential exchangeability.) We will refer this set of conditional independences as global sequential exchangeability for Y^g under any strategy g that involves interventions on both components of the time-varying treatment (A_0, A_1) . We use the word “global” because the sequential exchangeability holds for *all* static and dynamic strategies g .

A sequentially randomized experiment—an experiment in which treatment A_k at each time k is randomly assigned with a probability that depends only on the values of their prior covariate history \bar{L}_k and treatment history \bar{A}_{k-1} —implies global sequential exchangeability for Y^g . That is, for any strategy g , the treated and the untreated at each time k are exchangeable for Y^g conditional on prior covariate history \bar{L}_k and any observed treatment history $\bar{A}_{k-1} = g(\bar{A}_{k-2}, \bar{L}_{k-1})$ compatible with strategy g . Formally, *global sequential exchangeability* for Y^g is defined as

$$Y^g \perp\!\!\!\perp A_k|\bar{A}_{k-1} = g(\bar{A}_{k-2}, \bar{L}_{k-1}), \bar{L}_k \text{ for all strategies } g \text{ and } k = 0, 1, \dots, K$$

This form of sequential exchangeability (there are others, as we will see) always holds in any causal graph which, like Figure 19.2, has no arrows from the unmeasured variables U into the treatment variables A . Therefore global sequential exchangeability for Y^g holds in sequentially randomized experiments and observational studies in which the probability of receiving treatment at each time depends on their treatment and measured covariate history $(\bar{A}_{k-1}, \bar{L}_k)$ and, conditional on this history, does not depend on any unmeasured causes of the outcome.

That is, in observational studies represented by Figure 19.2 the mean of the counterfactual outcome $E[Y^g]$ under all strategies g is identified, whereas in observational studies represented by Figure 19.3 no mean counterfactual outcome $E[Y^g]$ is identified. In observational studies represented by other causal diagrams, the mean counterfactual outcome $E[Y^g]$ under some but not all strategies g is identified.

For example, consider an observational study represented by the causal diagram in Figure 19.4, which includes an unmeasured variable W_0 . In our HIV example, W_0 could be an indicator for a scheduled clinic visit at time 0 that was not recorded in our database. In that case W_0 would be a common

Individuals with $A_0 = g(L_0)$ are those with observed treatment A_0 equal to (i.e., compatible with) the treatment $g(L_0)$ they would have received under strategy g . For those with observed treatment history $[A_0 = g(L_0), A_1 = g(A_0, L_0, L_1)]$ compatible with strategy g through the end of follow-up, the counterfactual outcome Y^g is equal to the observed outcome Y and therefore also to the counterfactual outcome under the strategy $a_0 = A_0, a_1 = A_1$.

In a randomized experiment represented by Figure 19.1, we expect that sequential unconditional exchangeability for Y holds, that is, $Y^{\bar{a}} \perp\!\!\!\perp A_k|\bar{A}_{k-1} = \bar{a}_{k-1}$ for all static strategies \bar{a}

Technical Point 19.2

Positivity and consistency for time-varying treatments. The positivity condition needs to be generalized from the fixed version “if $f_L(l) \neq 0$, $f_{A|L}(a|l) > 0$ for all a and l ” to the sequential version

$$\text{If } f_{\bar{A}_{k-1}, \bar{L}_k}(\bar{a}_{k-1}, \bar{l}_k) \neq 0, \text{ then } f_{A_k | \bar{A}_{k-1}, \bar{L}_k}(a_k | \bar{a}_{k-1}, \bar{l}_k) > 0 \text{ for all } (\bar{a}_k, \bar{l}_k)$$

In a sequentially randomized experiment, positivity will hold if the randomization probabilities at each time k are never either 0 nor 1, no matter the past treatment and covariate history. If we are interested in a particular strategy g , the above positivity condition needs only hold for treatment histories compatible with g , i.e., for each k , $a_k = g(\bar{a}_{k-1}, \bar{l}_k)$.

The consistency condition also needs to be generalized from the fixed version “If $A = a$ for a given individual, then $Y^a = Y$ for that individual” to the sequential version

$$Y^{\bar{a}} = Y^{\bar{a}^*} \text{ if } \bar{a}^* = \bar{a}; Y^{\bar{a}} = Y \text{ if } \bar{A} = \bar{a}; \bar{L}_k^{\bar{a}} = \bar{L}_k^{\bar{a}^*} \text{ if } \bar{a}_{k-1}^* = \bar{a}_{k-1}, \bar{L}_k^{\bar{a}} = \bar{L}_k \text{ if } \bar{A}_{k-1} = \bar{a}_{k-1}$$

where $\bar{L}_k^{\bar{a}}$ is the counterfactual L -history through time k under strategy \bar{a} . Technically, the identification of effects of time-varying treatments on Y requires weaker consistency conditions: “If $\bar{A} = \bar{a}$ for a given individual, then $Y^{\bar{a}} = Y$ for that individual” is sufficient for static strategies, and “For any strategy g , if $A_k = g_k(\bar{A}_{k-1}, \bar{L}_k)$ at each time k for a given individual, then $Y^g = Y$ ” is sufficient for dynamic strategies. However, the stronger sequential consistency is a natural condition that we will always accept.

Note that, if we expect that the interventions “treat in month k ” corresponding to $A_k = 1$ and “do not treat in month k ” corresponding to $A_k = 0$ are sufficiently well defined at all times k , then all static and dynamic strategies involving A_k will be similarly well defined.

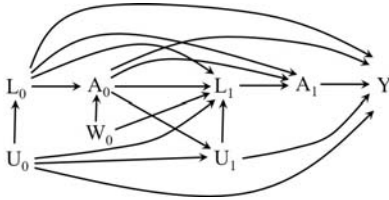


Figure 19.4

cause of treatment A_0 and of (scheduling and thus) obtaining a somewhat noisy measurement L_1 of CD4 cell count at time 1, with U_1 representing the underlying but unknown true value of CD4 cell count. Even though W_0 is unmeasured, the mean counterfactual outcome is still identified under any static strategy $g = \bar{a}$; however, the mean counterfactual outcome $E[Y^g]$ is not identified under any dynamic strategy g with treatment assignment depending on L_1 . To illustrate why identification is possible under some but not all strategies, we will use SWIGs in the next section.

In addition to some form of sequential exchangeability, causal inference involving time-varying treatments also requires a sequential version of the conditions of positivity and consistency. In a sequentially randomized experiment, both sequential positivity and consistency are expected to hold (see Technical Point 19.2). Below we will assume that sequential positivity and consistency hold. Under the three identifiability conditions, we can identify the mean counterfactual outcome $E[Y^g]$ under a strategy of interest g as long as we use methods that appropriately adjust for treatment and covariate history $(\bar{A}_{k-1}, \bar{L}_k)$, such as the g-formula (standardization), IP weighting, and g-estimation.

19.5 Identifiability under some but not all treatment strategies

Pearl and Robins (1995) proposed a generalized backdoor criterion for static strategies. Robins (1997) extended the procedure to dynamic strategies.

In Chapter 7, we presented a graphical rule—the backdoor criterion—to assess whether exchangeability holds for a time-fixed treatment under a particular causal diagram. The backdoor criterion can be generalized for time-varying treatments. For example, for static strategies, a sufficient condition for identification of the causal effect of treatment strategies is that, at each time k ,

all backdoor paths into A_k that do not go through any future treatment are blocked.

However, the *generalized backdoor criterion* does not directly show the connection between blocking backdoor paths and sequential exchangeability, because the procedure is based on causal directed acyclic graphs that do not include counterfactual outcomes. An alternative graphical check for identifiability of causal effects is based on SWIGs, also discussed in Chapter 7. SWIGs are especially helpful for time-varying treatments.

Consider the causal diagrams in Figures 19.5 and 19.6, which are simplified versions of those in Figures 19.2 and 19.4. We have omitted the nodes U_0 and L_0 and the arrow from A_0 to U_1 . In addition, the arrow from L_1 to Y is absent so L_1 is no longer a direct cause of Y . Figures 19.5 and 19.6 (like Figures 19.2 and 19.4) differ in whether A_k and subsequent covariates L_t for $t > k$ share a cause W_k .

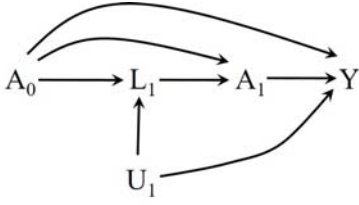


Figure 19.5

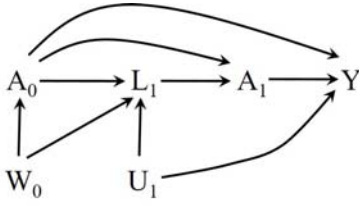


Figure 19.6

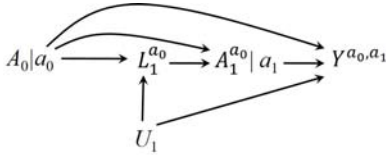


Figure 19.7

As discussed in Part I of this book, a SWIG represents a counterfactual world under a particular intervention. The SWIG in Figure 19.7 represents the world in Figure 19.5 if all individuals had received the static strategy (a_0, a_1) , where a_0 and a_1 can take values 0 or 1. For example, Figure 19.7 can be used to represent the world under the strategy “always treat” ($a_0 = 1, a_1 = 1$) or under the strategy “never treat” ($a_0 = 0, a_1 = 0$). To construct this SWIG, we first split the treatment nodes A_0 and A_1 . The right side of the split treatments represents the value of treatment under the intervention. The left side represents the value of treatment that would have been observed when intervening on all previous treatments. Therefore, the left side of A_0 is precisely A_0 because there are no previous treatments to intervene on, and the left side of A_1 is the counterfactual treatment $A_1^{a_0}$ that would be observed after setting A_0 to the value a_0 . All arrows into a given treatment in the original causal diagram now point into the left side, and all arrows out of a given treatment now originate from the right side. The outcome variable is the counterfactual outcome Y^{a_0, a_1} and the covariates L are replaced by their corresponding counterfactual variables. Note that we write the counterfactual variable corresponding to L_1 under strategy (a_0, a_1) as $L_1^{a_0}$, rather than $L_1^{a_0, a_1}$, because a future intervention on A_1 cannot affect the value of earlier L_1 .

Unlike the directed acyclic graph in Figure 19.5, the SWIG in Figure 19.7 does include the counterfactual outcome, which means that we can visually check for exchangeability using d-separation.

In Figure 19.7, d-separation shows that both $Y^{a_0, a_1} \perp\!\!\!\perp A_0$ and $Y^{a_0, a_1} \perp\!\!\!\perp A_1^{a_0} | A_0, L_1^{a_0}$ hold for any static strategy (a_0, a_1) . Note that this second conditional independence holds even though there seems to be an open path $A_1^{a_0} \leftarrow a_0 \rightarrow L_1^{a_0} \leftarrow U_1 \rightarrow Y^{a_0, a_1}$. However, this path is actually blocked for the following reason. In the counterfactual world, a_0 is a constant and in probability statements constants are always implicitly conditioned on even though, by convention, they are not shown in the conditioning event. However, when checking d-separation we need to remember that constants are conditioned on, blocking the above path.

The second conditional independence $Y^{a_0, a_1} \perp\!\!\!\perp A_1^{a_0} | A_0, L_1^{a_0}$ implies, by definition, $Y^{a_0, a_1} \perp\!\!\!\perp A_1^{a_0} | A_0 = a_0, L_1^{a_0}$ in the subset of individuals who received treatment $A_0 = a_0$. Therefore, by consistency, we conclude that $Y^{a_0, a_1} \perp\!\!\!\perp A_0$ and $Y^{a_0, a_1} \perp\!\!\!\perp A_1 | A_0 = a_0, L_1$ holds under the causal diagram in Figure 19.5, which corresponds to the SWIG in Figure 19.7 where we can actually check for exchangeability. If there were multiple time points, we would say that

$$Y^{\bar{a}} \perp\!\!\!\perp A_k | \bar{A}_{k-1} = \bar{a}_{k-1}, \bar{L}_k \text{ for } k = 0, 1 \dots K$$

$Y^{a_0, a_1} \perp\!\!\!\perp A_1^{a_0} | A_0 = a_0, L_1^{a_0}$ equals $Y^{a_0, a_1} \perp\!\!\!\perp A_1 | A_0 = a_0, L_1$ because, by consistency, $L_1^{a_0} = L_1$ and $A_1^{a_0} = A_1$ when $A_0 = a_0$.

Technical Point 19.3

The many forms of sequential exchangeability. Consider a sequentially randomized experiment of a time-varying treatment A_k with multiple time points $k = 0, 1, \dots, K$. The SWIG that represents this experiment is just a longer version of Figure 19.7. The following conditional independence can be directly read from the SWIG:

$$(Y^{\bar{a}}, \underline{L}_{k+1}^{\bar{a}}) \perp\!\!\!\perp A_k^{\bar{a}_{k-1}} | \bar{A}_{k-1}^{\bar{a}_{k-2}}, \bar{L}_k^{\bar{a}_{k-1}}$$

where $\underline{L}_{k+1}^{\bar{a}}$ is the counterfactual covariate history from time $k+1$ through the end of follow-up. The above conditional independence implies $(Y^{\bar{a}}, \underline{L}_{k+1}^{\bar{a}}) \perp\!\!\!\perp A_k^{\bar{a}_{k-1}} | \bar{A}_{k-1}^{\bar{a}_{k-2}} = \bar{a}_{k-1}, \bar{L}_k^{\bar{a}_{k-1}}$ for the particular instance $\bar{A}_{k-1}^{\bar{a}_{k-2}} = \bar{a}_{k-1}$, with \bar{a}_{k-1} being a component of strategy \bar{a} . Because of consistency, the last conditional independence statement equals

$$(Y^{\bar{a}}, \underline{L}_{k+1}^{\bar{a}}) \perp\!\!\!\perp A_k | \bar{A}_{k-1} = \bar{a}_{k-1}, \bar{L}_k$$

which we will refer to as *sequential exchangeability* (for Y^g and all other counterfactual variables) in this book.

Interestingly, the sequential exchangeability condition only refers to static strategies $g = \bar{a}$, but it is sufficient to identify the outcome distribution under any static and dynamic strategies g (Robins 1986). This identification results from the joint conditional independence between $(Y^{\bar{a}}, \underline{L}_{k+1}^{\bar{a}})$ and A_k . Note that, for dynamic strategies, sequential exchangeability does not follow the separate independences $Y^{\bar{a}} \perp\!\!\!\perp A_k | \bar{A}_{k-1} = \bar{a}_{k-1}, \bar{L}_k$ and $\underline{L}_{k+1}^{\bar{a}} \perp\!\!\!\perp A_k | \bar{A}_{k-1} = \bar{a}_{k-1}, \bar{L}_k$.

Stronger conditional independences are expected to hold in a sequentially randomized experiment, but they (i) cannot be read from SWIGs and (ii) are not necessary for identification of the causal effect in the population. For example, sequential randomization implies $(Y^{\bar{a}}, \underline{L}_{k+1}^{\bar{a}}) \perp\!\!\!\perp A_k | \bar{A}_{k-1}, \bar{L}_k$, a condition stronger than our sequential exchangeability because it holds under any treatment history \bar{A}_{k-1} , including those that cannot possibly exist in a the counterfactual world under strategy \bar{a} (i.e., it is cross-world). Of course $(Y^{\bar{a}}, \underline{L}_{k+1}^{\bar{a}}) \perp\!\!\!\perp A_k | \bar{A}_{k-1}, \bar{L}_k$ is equal to our sequential exchangeability when we instantiate it for $\bar{A}_{k-1} = \bar{a}_{k-1}$.

An even stronger condition that is expected to hold in sequentially randomized experiments is

$$(Y^{\bar{A}}, \bar{L}^{\bar{A}}) \perp\!\!\!\perp A_k | \bar{A}_{k-1}, \bar{L}_k$$

where, for a dichotomous treatment A_k , \bar{A} denotes the set of all 2^K static strategies \bar{a} , $Y^{\bar{A}}$ denotes the set of all counterfactual outcomes $Y^{\bar{a}}$, and $\bar{L}^{\bar{A}}$ denotes the set of all counterfactual covariate histories. Using a terminology analogous to that of Technical Point 2.1, we refer to this joint independence condition as *full sequential exchangeability*.

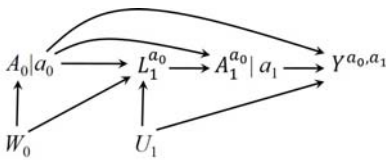


Figure 19.8

We refer to the above condition as *static sequential exchangeability* for $Y^{\bar{a}}$, which is weaker than global sequential exchangeability for Y^g , because it only requires conditional independence between counterfactual outcomes $Y^{\bar{a}}$ indexed by static strategies $g = \bar{a}$ and treatment A_k . Static sequential exchangeability is sufficient to identify the mean counterfactual outcome under any static strategy $g = \bar{a}$. See also Technical Point 19.3.

Static sequential exchangeability also holds under the causal diagram in Figure 19.6, as can be checked by applying d-separation to its corresponding SWIG in Figure 19.8. Therefore, in an observational study represented by Figure 19.6, we can identify the mean counterfactual outcome under any static strategy (a_0, a_1) .

Let us return to Figure 19.5. Let us now assume that the arrow from L_1 to A_1 were missing. In that case, the arrow from $L_1^{a_0}$ to $A_1^{a_0}$ would also be missing from the SWIG in Figure 19.7. It would then follow by d-separation that unconditional sequential exchangeability holds, and therefore that the mean counterfactual outcome under any static strategy could be identified without data on L_1 . Now let us assume that, in Figure 19.5, there was an

Fine Point 19.3

Baseline confounders. For simplicity, the causal graphs depicted in this chapter did not include a baseline confounder L_0 . If we included L_0 in Figure 19.9, then we could have considered a strategy in which the random variable representing the intervention $g(L_0)$ replaces g_0 . Then, when checking d-separation between A_1^g and Y^g on the graph, $Y^g \perp\!\!\!\perp A_1^g | A_0, g_0(L_0), L_0, L_1^g$, we need to condition on the entire past, including $g(L_0)$. If we instantiate this expression at $A_0 = g_0(L_0)$, then the intervention variable can be removed from the conditioning event because $g_0(L_0)$ is now equal to the observed A_0 and thus is redundant. That is, we have now $Y^g \perp\!\!\!\perp A_1^g | A_0 = g_0(L_0), L_0, L_1^g$ which, by consistency, is $Y^g \perp\!\!\!\perp A_1 | A_0 = g_0(L_0), L_0, L_1$. This conditional independence is sequential exchangeability for Y^g and treatment A_1 when there is also a baseline confounder L_0 .

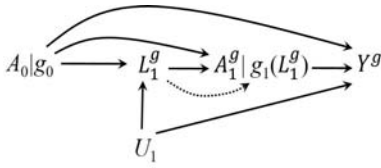


Figure 19.9

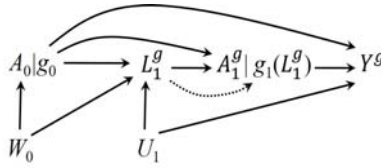


Figure 19.10

Technically, what we read from the SWIG is $Y^g \perp\!\!\!\perp A_1^g | A_0, L_1^g$ which, by consistency, implies $Y^g \perp\!\!\!\perp A_1 | A_0 = g_0, L_1$

arrow from U_1 to A_1 . Then the SWIG in Figure 19.7 would include an arrow from U_1 to $A_1^{a_0}$, and that no form of (unconditional or conditional) sequential exchangeability would hold. Therefore the counterfactual mean would not be identified under any strategy.

We now discuss the SWIGs for Figures 19.5 and 19.6 under dynamic regimes. The SWIG in Figure 19.9 represents the world of Figure 19.5 under a dynamic treatment strategy $g = [g_0, g_1(L_1)]$ in which treatment A_0 is assigned a fixed value g_0 (either 0 or 1), and treatment A_1 at time $k = 1$ is assigned a value $g_1(L_1^g)$ that depends on the value of L_1^g that was observed after having assigned treatment value g_0 at time $k = 0$. For example, g may be the strategy “do not treat at time 0, treat at time 1 only if CD4 cell count is low, i.e., if $L_1^g = 1$ ”. Under this strategy $g_0 = 0$ for everybody, and $g_1(L_1^g) = 1$ when $L_1^g = 1$ and $g_1(L_1^g) = 0$ when $L_1^g = 0$. Therefore the SWIG includes an arrow from L_1^g to $g_1(L_1^g)$. This arrow was not part of the original causal graph; it is the result of the intervention. We therefore draw this arrow differently from the others, even though we need to treat as any other arrow when evaluating d-separation. The outcome in the SWIG is the counterfactual outcome Y^g under the dynamic treatment strategy g .

By applying d-separation to the SWIG in Figure 19.9, we find that both $Y^g \perp\!\!\!\perp A_0$ and $Y^g \perp\!\!\!\perp A_1 | A_0 = g_0, L_1$ hold for any strategy g . That is, global sequential exchangeability for Y^g holds, which means that we can identify the mean counterfactual outcome under all strategies g (see also Fine Point 19.3). This result, however, does not hold for the causal diagram in Figure 19.6.

The SWIG in Figure 19.10 represents the world of Figure 19.6 under a dynamic treatment strategy $g = [g_0, g_1(L_1)]$. By applying d-separation to the SWIG in Figure 19.9, we find that $Y^g \perp\!\!\!\perp A_0$ does not hold because of the open path $A_0 \leftarrow W_0 \rightarrow L_1^g \rightarrow g_1(L_1^g) \rightarrow Y^g$. That is, global sequential exchangeability for Y^g does not hold, which means that we cannot identify the mean counterfactual outcome under all strategies g .

In summary, in observational studies (or sequentially randomized trials) represented by Figure 19.5, global sequential exchangeability for Y^g holds, and therefore the data can be used to validly estimate causal effects involving static and dynamic strategies. On the other hand, in observational studies represented by Figure 19.6, only the weaker condition for static regimes holds, and therefore the data can be used to validly estimate causal effects involving static, but not dynamic, strategies. Another way to think about this is that, under Figure 19.5, one can consistently estimate the total effect of every treatment component A_k , whereas that is not the case under Figure 19.6. Specifically, as described in Fine Point 7.2, the effect of treatment A_0 alone cannot be

consistently estimated.

One last example. Consider Figure 19.11 which is equal to Figure 19.6 except for the presence of an arrow from L_1 to Y , and its corresponding SWIG under a static strategy in Figure 19.12. We can use d-separation to show that neither global sequential exchangeability for Y^g nor static sequential exchangeability for $Y^{\bar{a}}$ hold. Therefore, in observational study represented by Figure 19.11, we cannot use the data to validly estimate causal effects involving any strategies.

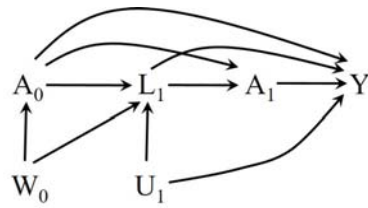


Figure 19.11

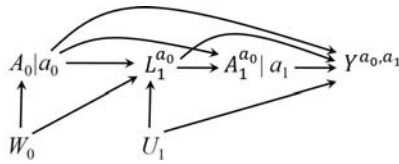


Figure 19.12

No form of sequential exchangeability is guaranteed to hold in observational studies. Achieving approximate exchangeability requires expert knowledge, which will guide investigators in the design of their studies to measure as many of the relevant variables \bar{L}_k as possible. For example, in an HIV study, experts would agree that time-varying variables like CD4 cell count, viral load, symptoms need to be appropriately measured and adjusted for.

But the question “Are the measured covariates sufficient to ensure sequential exchangeability?” can never be answered with certainty. Yet we can use our expert knowledge to organize our beliefs about exchangeability and represent them in a causal diagram. Figures 19.1 to 19.4 are examples of causal diagrams that summarize different scenarios. Note that we drew these causal diagrams in the absence of selection (e.g., censoring by loss to follow-up), which we defer to later chapters so that we can concentrate on confounding here.

Consider Figure 19.5. Like in Part I of this book, suppose that we are interested in the effect of the time-fixed treatment A_1 on the outcome Y . We say that there is confounding for the effect of A_1 on Y because A_1 and Y share the cause U , i.e., because there is an open backdoor path between A_1 and Y through U . To estimate this effect without bias, we need to adjust for confounders of the effect of the treatment A_1 on the outcome Y , as explained in Chapter 7. In other words, we need to be able to block all open backdoor paths between A_1 and Y . This backdoor path $A_1 \leftarrow L_1 \leftarrow U \rightarrow Y$ cannot be blocked by conditioning on the common cause U because U is unmeasured and therefore unavailable to the investigators. However, this backdoor path can be blocked by conditioning on L_1 , which is measured. Thus, if the investigators collected data on L_1 for all individuals, there would be no unmeasured confounding for the effect of A_1 . We then say that L_1 is a confounder for the effect of A_1 , even though the actual common cause of A_1 and Y was the unmeasured U (re-read Section 7.3 if you need to refresh your memory about confounding and causal diagrams).

As discussed in Chapter 7, the confounders do not have to be direct causes of the outcome. In Figure 19.5, the arrow from the confounder L_1 to the outcome Y does not exist. Then the source of the confounding (i.e., the causal confounder) is the unmeasured common cause U . Nonetheless, because data on L_1 suffice to block the backdoor paths from A_1 to Y and thus to control confounding, we refer to L_1 as a confounder for the effect of A_1 on Y .

Now imagine the very long causal diagram that contains all time points k , and in which L_k affects subsequent treatments $A_k, A_{k+1} \dots$. Suppose that we want to estimate the causal effects of each time-fixed treatment A_k on the outcome Y . Then L_1 and L_2 are necessary to block all backdoor paths between A_2 and Y , (L_1, L_2, L_3) are necessary to block all backdoor paths

A second backdoor path gets open after conditioning on collider L_1 :
 $A_1 \leftarrow A_0 \rightarrow L_1 \leftarrow U \rightarrow Y$
 This second backdoor path can be safely blocked by conditioning on prior treatment A_0 , which should be available to investigators.

Fine Point 19.4

A precise definition of time-varying confounding. In the absence of selection bias, we say there is confounding for causal effects involving $E[Y^{\bar{a}}]$ if $E[Y^{\bar{a}}] \neq E[Y|A = \bar{a}]$, that is, if the mean outcome had, contrary to fact, all individuals in the study followed strategy \bar{a} differs from the mean outcome among the subset of individuals who followed strategy \bar{a} in the actual study.

We say the confounding is solely time-fixed (i.e., wholly attributable to baseline covariates) if $E[Y^{\bar{a}}|L_0] = E[Y|A = \bar{a}|L_0]$, as would be the case if the only arrows pointing into A_1 in Figure 19.2 were from A_0 and L_0 . In contrast, if the identifiability conditions hold, but $E[Y^{\bar{a}}|L_0] \neq E[Y|A = \bar{a}|L_0]$, we say that time-varying confounding is present. If the identifiability conditions do not hold, as in Figure 19.3, we say that there is unmeasured confounding.

The concept of time-varying confounder was briefly introduced in Fine Point 7.1. Time-varying confounders are sometimes referred to as time-dependent confounders.

between A_3 and Y , and so on. That is, at each time k , the covariate history \bar{L}_k will be needed, together with the treatment history \bar{A}_{k-1} , to block the backdoor paths between subsequent treatment A_k and the outcome Y . Thus, if the investigators collected data on \bar{L}_k for all individuals, there would be no unmeasured confounding for the effect of \bar{A} . We then say that the time-varying covariates in \bar{L}_k are *time-varying confounders* for the effect of the time-varying treatment \bar{A} on Y at several (or, in our example, all) times k in the study. See Fine Point 19.4 for a precise definition of time-varying confounding.

Unfortunately, we cannot empirically confirm that all confounders, whether time-fixed or time-varying, are measured. That is, we cannot empirically differentiate between Figure 19.2 with no unmeasured confounding and Figure 19.3 with unmeasured confounding. Interestingly, even if all confounders were correctly measured and modeled, most adjustment methods may still result in biased estimates when comparing treatment strategies. The next chapter explains why g-methods are the appropriate approach to adjust for time-varying confounders.

Chapter 20

TREATMENT-CONFOUNDER FEEDBACK

The previous chapter identified sequential exchangeability as a key condition to identify the causal effects of time-varying treatments. Suppose that we have a study in which the strongest form of sequential exchangeability holds: the measured time-varying confounders are sufficient to validly estimate the causal effect of any treatment strategy. Then the question is what confounding adjustment method to use. The answer to this question highlights a key problem in causal inference about time-varying treatments: treatment-confounder feedback.

When treatment-confounder feedback exists, using traditional adjustment methods may introduce bias in the effect estimates. That is, even if we had all the information required to validly estimate the average causal effect of any treatment strategy, we would be generally unable to do so. This chapter describes the structure of treatment-confounder feedback and the reasons why traditional adjustment methods fail.

20.1 The elements of treatment-confounder feedback

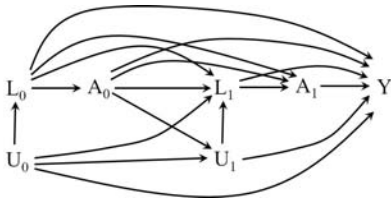


Figure 20.1

Consider again the sequentially randomized trial of HIV-positive individuals that we discussed in the previous chapter. For every person in the study, we have data on treatment A_k (1: treated, 0: untreated) and covariates L_k at each month of follow-up $k = 0, 1, 2, \dots, K$, and on an outcome Y that measures health status at month $K + 1$. The causal diagram in Figure 20.1 (which is equal to the one in Figure 19.2) represents the first two months of the study. The time-varying covariates L_k are time-varying confounders.

Something else is going on in Figure 20.1. Not only is there an arrow from CD4 cell count L_k to treatment A_k , but also there is an arrow from treatment A_{k-1} to future CD4 cell count L_k —because receiving treatment A_{k-1} increases future CD4 cell count L_k . That is, the confounder affects the treatment *and* the treatment affects the confounder. There is *treatment-confounder feedback* (see also Fine Point 20.1).

Note that time-varying confounding can occur without treatment-confounder feedback. The causal diagram in Figure 20.2 is the same as the one in Figure 20.1, except that the arrows from treatment A_{k-1} to future L_k and U_k have been deleted. In a setting represented by this diagram, the time-varying covariates L_k are time-varying confounders, but they are not affected by prior treatment. Therefore, there is time-varying confounding, but there is no treatment-confounder feedback.

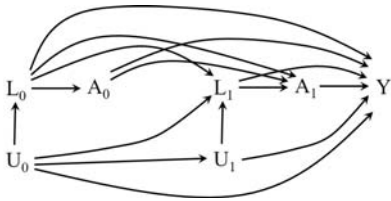


Figure 20.2

Treatment-confounder feedback creates an interesting problem for causal inference. To state the problem in its simplest form, let us simplify the causal diagram in Figure 20.1 a bit more. Figure 20.3 is the smallest subset of Figure 20.1 that illustrates treatment-confounder feedback in a sequentially randomized trial with two time points. When drawing the causal diagram in Figure 20.3, we made two simplifications:

- Because our interest is in the implications of confounding by L_1 , we did not bother to include a node L_0 for baseline CD4 cell count. Just

Fine Point 20.1

Representing feedback cycles with acyclic graphs. Interestingly, an *acyclic* graph—like the one in Figure 20.1—can be used to represent a treatment-confounder feedback loop or *cycle*. The trick to achieve this visual representation is to elaborate the treatment-confounder feedback loop in time. That is, $A_{k-1} \rightarrow L_k \rightarrow A_k \rightarrow L_{k+1}$ and so on.

The representation of feedback cycles with acyclic graphs also requires that time be considered as a discrete variable. That is, we say that treatment and covariates can change during each interval $[k, k + 1)$ for $k = 0, 1, \dots, K$, but we do not specify when exactly during the interval the change takes place. This discretization of time is not a limitation in practice: the length of the intervals can be chosen to be as short as the granularity of the data requires. For example, in a study where individuals see their doctors once per month or less frequently (as in our HIV example), time may be safely discretized into month intervals. In other cases, year intervals or day intervals may be more appropriate. Also, as we said in Chapter 17, time is typically measured in discrete intervals (years, months, days) any way, so the discretization of time is often not even a choice.

suppose that treatment A_0 is marginally randomized and treatment A_1 is conditionally randomized given L_1 .

- The unmeasured variable U_0 is not included.
- There is no arrow from A_0 to A_1 , which implies that treatment is assigned using information on L_1 only.
- There are no arrows from A_0 , L_1 , and A_1 to Y , which would be the case if treatment has no causal effect on the outcome Y of any individual, i.e., the sharp null hypothesis holds.

None of these simplifications affect the arguments below. A more complicated causal diagram would not add any conceptual insights to the discussion in this chapter; it would just be harder to read.

Now suppose that treatment has no effect on any individual's Y , which implies the causal diagram in Figure 20.3 is the correct one, but the investigators do not know it. and that we have data on treatment A_0 in month 0 and A_1 in month 1, on the confounder CD4 cell count L_1 at the start of month 1, and on the outcome Y at the end of follow-up. We wish to use these data to estimate the average causal effect of the static treatment strategy “always treat”, $(a_0 = 1, a_1 = 1)$, compared with the static treatment strategy “never treat”, $(a_0 = 0, a_1 = 0)$ on the outcome Y , that is, $E[Y^{a_0=1, a_1=1}] - E[Y^{a_0=0, a_1=0}]$. According to Figure 20.3, the true, but unknown to the investigator, average causal effect is 0 because there are no forward-directed paths from either treatment variable to the outcome. That is, one cannot start at either A_0 or A_1 and, following the direction of the arrows, arrive at Y .

Figure 20.3 can depict a sequentially randomized trial because there are no direct arrows from the unmeasured U into the treatment variables. Therefore, as we discussed in the previous chapter, we should be able to use the observed data on A_0 , L_1 , A_1 , and Y to conclude that $E[Y^{a_0=1, a_1=1}] - E[Y^{a_0=0, a_1=0}]$ is equal to 0. However, as we explain in the next section, we will not generally be able to correctly estimate the causal effect when we adjust for L_1 using traditional methods, like stratification, outcome regression, and matching. That is, in this example, an attempt to adjust for the confounder L_1 using these methods will generally result in an effect estimate that is different from 0, and thus invalid.

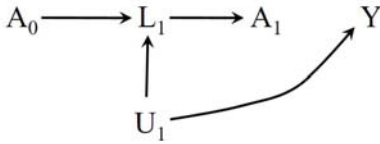


Figure 20.3

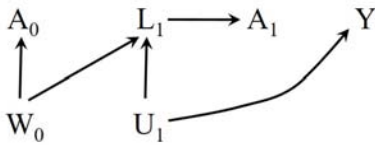


Figure 20.4

Figure 20.3 represents either a sequentially randomized trial or an observational study with no unmeasured confounding; Figure 20.4 represents an observational study.

In other words, when there are time-varying confounders and treatment-confounder feedback, traditional methods cannot be used to correctly adjust for those confounders. Even if we had sufficient longitudinal data to ensure sequential exchangeability, traditional methods would not generally provide a valid estimate of the causal effect of any treatment strategies. In contrast, g-methods appropriately adjust for the time-varying confounders even in the presence of treatment-confounder feedback.

This limitation of traditional methods applies to settings in which the time-varying confounders are affected by prior treatment as in Figure 20.3, but also to settings in which the time-varying confounders share causes W with prior treatment as in Figure 20.4, which is a subset of Figure 19.4. We refer to both Figures 20.3 and 20.4 (and Figures 19.2 and 19.4) as examples of treatment-confounder feedback. The next section explains why traditional methods cannot adequately handle treatment-confounder feedback.

20.2 The bias of traditional methods

This is an ideal trial with full adherence to the assigned treatment and no losses to follow-up.

Table 20.1

N	A_0	L_1	A_1	Mean Y
2400	0	0	0	84
1600	0	0	1	84
2400	0	1	0	52
9600	0	1	1	52
4800	1	0	0	76
3200	1	0	1	76
1600	1	1	0	44
6400	1	1	1	44

If there were additional times k at which treatment A_k were affected by L_k , then L_k would be a time-varying confounder

Specifically, there is no direct arrow from L_1 to Y ; otherwise A_0 would have an effect on Y through L_1

To illustrate the bias of traditional methods, let us consider a (hypothetical) sequentially randomized trial with 32,000 HIV-positive individuals and two time points $k = 0$ and $k = 1$. Treatment $A_0 = 1$ is randomly assigned at baseline with probability 0.5. Treatment A_1 is randomly assigned in month 1 with a probability that depends only on the value of CD4 cell count L_1 at the start of month 1—0.4 if $L_1 = 0$ (high), 0.8 if $L_1 = 1$ (low). The outcome Y , which is measured at the end of follow-up, is a function of CD4 cell count, concentration of virus in the serum, and other clinical measures, with higher values of Y signifying better health.

Table 20.1 shows the data from this trial. To save space, the table displays one row per combination of values of A_0 , L_1 , and A_1 , rather than one row per individual. For each of the eight combinations, the table provides the number of subjects N and the mean value of the outcome $E[Y|A_0, L_1, A_1]$. Thus, row 1 shows that the mean of the 2400 individuals with $(A_0 = 0, L_1 = 1, A_1 = 0)$ was $E[Y|A_0 = 0, L_1 = 0, A_1 = 0] = 84$. In this sequentially randomized trial, the identifiability conditions—sequential exchangeability, positivity, and consistency—hold. By design, there are no confounders for the effect of A_0 on Y , and L_1 is the only confounder for the effect of A_1 on Y so (conditional on L_1) sequential exchangeability holds. By inspection of Table 20.1, we can conclude that the positivity condition is satisfied, because otherwise one or more of the eight rows would have zero individuals.

The causal diagram in Figure 20.3 depicts this sequentially randomized experiment when the sharp null hypothesis holds. To check whether the data in Table 20.1 are consistent with the causal diagram in Figure 20.3, we can separately estimate the average causal effects of each of the time-fixed treatments A_0 and A_1 within levels of past covariates and treatment, which should all be null. In the calculations below, we will ignore random variability.

A quick inspection of the table shows that the average causal effect of treatment A_1 is indeed zero in all four strata defined by A_0 and L_1 . Consider the effect of A_1 in the 4000 individuals with $A_0 = 0$ and $L_1 = 0$, whose data are shown in rows 1 and 2 of Table 20.1. The mean outcome among those who did not receive treatment at time 1, $E[Y|A_0 = 0, L_1 = 0, A_1 = 0]$, is 84, and the mean outcome among those who did receive treatment at time 1,

Technical Point 20.1

G-null test. Suppose the sharp null hypothesis is true. Then any counterfactual outcome Y^g is the observed outcome Y . In this setting, global sequential exchangeability for Y^g can be written as $Y \perp\!\!\!\perp A_0 | L_0$ and $Y \perp\!\!\!\perp A_1 | A_0 = g(L_0), L_0, L_1$ in a study with two time points. The first independence implies no causal effect of A_0 in any strata defined by L_0 , and the second independence implies no causal effect of A_1 in any strata defined by L_1 and A_0 . Therefore, under sequential exchangeability, a test of these conditional independences is a test of the sharp null. This is the g-null test.

Conversely, the g-null theorem (Robins 1986) says that, if these conditional independences hold, then the distribution of Y^g and therefore the mean $E[Y^g]$ is the same for all g , and also equal to the distribution and mean of the observed Y . Note, however, that equality of distributions under all g only implies the sharp null hypothesis under a strong form of faithfulness that forbids perfect cancellations of effects. As discussed in Fine Point 6.2, we assume faithfulness throughout this book unless we say otherwise.

$E[Y|A_0 = 0, L_1 = 0, A_1 = 1]$, is also 84. Therefore the difference

$$E[Y|A_0 = 0, L_1 = 0, A_1 = 1] - E[Y|A_0 = 0, L_1 = 0, A_1 = 0]$$

is zero. Because the identifiability conditions hold, this associational difference validly estimates the average causal effect

$$E[Y^{a_1=1}|A_0 = 0, L_1 = 0] - E[Y^{a_1=0}|A_0 = 0, L_1 = 0]$$

in the stratum $(A_0 = 0, L_1 = 0)$. Similarly, it is easy to check that the average causal effect of treatment A_1 on Y is zero in the remaining three strata $(A_0 = 0, L_1 = 1)$, $(A_0 = 1, L_1 = 0)$, $(A_0 = 1, L_1 = 1)$, by comparing the mean outcome between rows 3 and 4, rows 5 and 6, and rows 7 and 8, respectively.

We can now show that the average causal effect of A_0 is also zero. To do so, we need to compute the associational difference $E[Y|A_0 = 1] - E[Y|A_0 = 0]$ which, because of randomization, is a valid estimator of the causal contrast $E[Y^{a_0=1}] - E[Y^{a_0=0}]$. The mean outcome $E[Y|A_0 = 0]$ among the 16,000 individuals treated at time 0 is the weighted average of the mean outcomes in rows 1, 2, 3 and 4, which is 60. And $E[Y|A_0 = 1]$, computed analogously, is also 60. Therefore, the average causal effect of A_0 is zero.

We have confirmed that the causal effects of A_0 and A_1 (conditional on the past) are zero when we treat A_0 and A_1 separately as time-fixed treatments. What if we now treat the joint treatment (A_0, A_1) as a time-varying treatment and compare two treatment strategies? For example, let us say that we want to compare the strategies “always treat” versus “never treat”, that is $(a_0 = 1, a_1 = 1)$ versus $(a_0 = 0, a_1 = 0)$. Because the identifiability conditions hold, the data in Table 20.1 should suffice to validly estimate this effect.

Because the effect for each of the individuals components of the strategy, a_0 and a_1 , is zero, it follows from the g-null theorem that the average causal effect $E[Y^{a_0=1, a_1=1}] - E[Y^{a_0=0, a_1=0}]$ is zero. But is this what we conclude from the data if we use conventional analytic methods? To answer this question, let us conduct two data analyses. In the first one, we do not adjust for the confounder L_1 , which should give us an incorrect effect estimate. In the second one, we do adjust for the confounder L_1 via stratification.

1. We compare the mean outcome in the 9600 individuals who were treated at both times (rows 6 and 8 of Table 20.1) with that in the 4800 individuals who were untreated at both times (rows 1 and 3). The respective averages are $E[Y|A_0 = 1, A_1 = 1] = 54.7$, and $E[Y|A_0 = 0, A_1 = 0] =$

The weighted average is

$$\frac{2400}{16000} \times 84 + \frac{1600}{16000} \times 84 + \frac{2400}{16000} \times 52 + \frac{9600}{16000} \times 52 = 60$$

$$\begin{aligned} E[Y|A_0 = 1, A_1 = 1] \\ \frac{3200}{9600} \times 76 + \frac{6400}{9600} \times 44 = 54.7 \\ E[Y|A_0 = 0, A_1 = 0] \\ \frac{2400}{4800} \times 84 + \frac{2400}{4800} \times 52 = 68.0 \end{aligned}$$

Note that, because the effect is -8 in both strata of L_1 , it is not possible that a weighted average of the stratum-specific effects will yield the correct value 0.

68. The associational difference is $54.7 - 68 = -13.3$ which, if interpreted causally, would mean that not being treated at either time is better than being treated at both times. This analysis gives the wrong answer—a non-null difference—because $E[Y|A_0 = a_0, A_1 = a_1]$ is not a valid estimator of $E[Y^{a_0, a_1}]$. Adjustment for the confounder L_1 is needed.

2. We adjust for L_1 via stratification. That is, we compare the mean outcome in individuals who were treated with that in individuals who were untreated at both times, within levels of L_1 . For example, take the stratum $L_1 = 0$. The mean outcome in the treated at both times, $E[Y|A_0 = 1, L_1 = 0, A_1 = 1]$, is 76 (row 6). The mean outcome in the untreated at both times, $E[Y|A_0 = 0, L_1 = 0, A_1 = 0]$, is 84 (row 1). The associational difference is $76 - 84 = -8$ which, if interpreted causally, would mean that, in the stratum $L_1 = 0$, not being treated at either time is better than being treated at both times. Similarly, the difference $E[Y|A_0 = 0, L_1 = 1, A_1 = 0]$ in the stratum is also -8 .

What? We said that the effect estimate should be 0, not -8 . How is it possible that the analysis adjusted for the confounder also gives a wrong answer? This estimate reflects the bias of traditional methods to adjust for confounding when there is treatment-confounder feedback. The next section explains why the bias arises.

20.3 Why traditional methods fail

Table 20.1 shows data from a sequentially randomized trial with treatment-confounder feedback, as represented by the causal diagram is represented in Figure 20.3. Even though no data on the unmeasured variable U_1 (immunosuppression level) is available, all three identifiability conditions hold: U_1 is not needed if we have data on the confounder L_1 . Therefore, as discussed in Chapter 19, we should be able to correctly estimate causal effects involving any static or dynamic treatment strategies. And yet our analyses in the previous section did not yield the correct answer, whether or not we adjusted for L_1 .

The problem was that we did not use the correct method to adjust for confounding. Stratification is a commonly used method to adjust for confounding, but it cannot handle treatment-confounder feedback. Stratification means estimating the association between treatment and outcome in subsets—strata—of the study population defined by the confounders— L_1 in our example. Because the variable L_1 can take only two values—1 if the CD4 cell count is low, and 0 otherwise—there are two such strata in our example. To estimate the causal effect in those with $L_1 = l$, we selected (i.e., conditioned or stratified on) the subset of the population with value $L_1 = l$.

But stratification can have unintended effects when the association measure is computed within levels of a variable L_1 that is caused by prior treatment A_0 . Indeed Figure 20.5 shows that conditioning on L_1 —a collider—opens the path $A_0 \rightarrow L_1 \leftarrow U_1 \rightarrow Y$. That is, stratification induces a noncausal association between the treatment A_0 at time 0 and the unmeasured variable U_1 , and therefore between A_0 and the outcome Y , within levels of L_1 . Among those with low CD4 count ($L_1 = 1$), being on treatment ($A_0 = 1$) becomes a marker for severe immunosuppression (high value of U_1); among those with a high level of CD4 ($L_1 = 0$), being off treatment ($A_0 = 0$) becomes a marker for milder

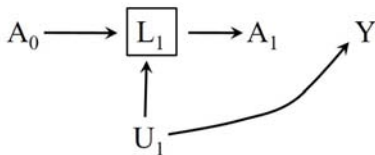


Figure 20.5

Fine Point 20.2

Confounders on the causal pathway. Conditioning on confounders L_1 which are affected by previous treatment can create selection bias even if the confounder is not on a causal pathway between treatment and outcome. In fact, no such causal pathway exists in Figures 20.5 and 20.6.

On the other hand, in Figure 20.7 the confounder L_1 for subsequent treatment A_1 lies on a causal pathway from earlier treatment A_0 to outcome Y , i.e., the path $A_0 \rightarrow L_1 \rightarrow Y$. Were the potential for selection bias not present in Figure 20.7 (e.g., were U_1 not a common cause of L_1 and Y), the associational differences within strata of L_1 could be an unbiased estimate of the direct effect of A_0 on Y not through L_1 , but still would not be an unbiased estimate of the overall effect of \bar{A} on Y , because the effect of A_0 mediated through L_1 is not included.

It is sometimes said that variables on a causal pathway between treatment and outcome cannot be considered as confounders, because adjusting for those variables will result in a biased effect estimate. However, this characterization of confounders is inaccurate for time-varying treatments. Figure 20.7 shows that a confounder for subsequent treatment A_1 can be on a causal pathway between past treatment A_0 and the outcome. As for whether adjustment for confounders on a causal pathway induces bias for the effect of a treatment strategy, that depends on the choice of adjustment method. Stratification will indeed induce bias; g-methods will not.

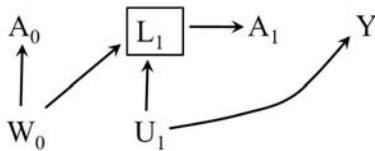


Figure 20.6

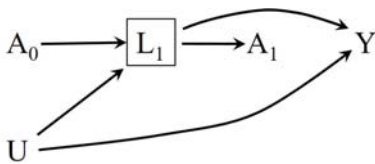


Figure 20.7

immunosuppression (low value of U_1). Thus, the side effect of stratification is to induce an association between treatment A_0 and outcome Y .

In other words, stratification eliminates confounding for A_1 at the cost of introducing selection bias for A_0 . The associational differences

$$E[Y|A_0 = 1, L_1 = l, A_1 = 1] - E[Y|A_0 = 0, L_1 = l, A_1 = 0]$$

may be different from 0 even if, as in our example, treatment has no effect on the outcome of any individuals at any time. This bias arises from choosing a subset of the study population by selecting on a variable L_1 affected by (a component A_0 of) the time-varying treatment. The net bias depends on the relative magnitude of the confounding that is eliminated and the selection bias that is created.

Technically speaking, the bias of traditional methods will occur not only when the confounders are affected by prior treatment (in randomized experiments or observational studies), but also when the confounders share an unmeasured cause W with prior treatment (in observational studies). In the observational study depicted in Figure 20.6, conditioning on the collider L_1 opens the path $A_0 \leftarrow W_0 \rightarrow L_1 \leftarrow U_1 \rightarrow Y$. For this reason, we referred to both settings in Figures 20.3 and 20.4—which cannot be distinguished using the observed data—as examples of treatment-confounder feedback.

The causal diagrams that we have considered to describe the bias of traditional methods are all very simple. They only represent settings in which treatment does not have a causal effect on the outcome. However, conditioning on a confounder in the presence of treatment-confounder feedback also induces bias when treatment has a non-null effect, as in Figure 20.1. The presence of arrows from A_0 , A_1 , or L_1 to Y does not change the fact that conditioning on L_1 creates an association between A_0 and Y that does not have a causal interpretation (see also Fine point 20.2). Also, our causal diagrams had only two time points and a limited number of nodes, but the bias of traditional methods will also arise from high-dimensional data with multiple time points and variables. In fact, the presence of time-varying confounders affected by previous treatment at multiple times increases the possibility of a large bias.

In general, valid estimation of the effect of treatment strategies is only possible when the joint effect of the treatment components A_k can be estimated

simultaneously and without bias. As we have just seen, this may be impossible to achieve using stratification, even when data on all time-varying confounders are available.

20.4 Why traditional methods cannot be fixed

We showed that stratification cannot be used as a confounding adjustment method when there is treatment-confounder feedback. But what about other traditional methods? For example, we could have used parametric outcome regression, rather than nonparametric stratification, to adjust for confounding. Would outcome regression succeed where plain stratification failed?

This question is particularly important for settings with high-dimensional data, because in high-dimensional settings we will be unable to conduct a simple stratified analysis like we did in the previous section. In Table 20.1, treatment A_k occurs at two months $k = 0, 1$, which means that there are only 2^2 static treatment strategies \bar{a} . But when the treatment A_k occurs at multiple points $k = 0, 1, \dots, K$, we will not be able to present a table with all the combinations of treatment values. If, as is not infrequent in practice, K is of the order of 100, then there are 2^{100} static treatment strategies \bar{a} , a staggering number that far exceeds the sample size of any study. The total number of treatment strategies is much greater when we consider dynamic strategies as well.

The number of data combinations is even greater because there are multiple confounders L_k measured at each time point k .

As we have been arguing since Chapter 11, we will need modeling to estimate average causal effects involving $E[Y^{\bar{a}}]$ when there are many possible treatment strategies \bar{a} . To do so, we will need to hypothesize a dose-response function for the effect of treatment history \bar{a} on the mean outcome Y . One possibility would be to assume that the effect of treatment strategies \bar{a} increases linearly as a function of the cumulative treatment under each strategy. Under this assumption, all strategies that assign treatment for exactly three months have the same effect, regardless of the period when those three months of treatment occur, e.g., the first 3 months of follow-up, the last 3 months of follow-up, etc. The price paid for modelling is yet another threat to the validity of our estimates due to possible model misspecification of the dose-response function.

Unfortunately, regression modeling does not remove the bias of traditional methods in the presence of treatment-confounder feedback, as we now show. For the data in Table 20.1, let us define cumulative treatment $cum(\bar{A}) = A_0 + A_1$, which can take 3 values: 0 (if the individual remains untreated at both times), 1 (if the subject is treated at time 1 only or at time 2 only), and 2 (if the subject is treated at both times). The treatment strategies of interest can then be expressed as “always treat” $cum(\bar{a}) = 2$, and “never treat” $cum(\bar{a}) = 0$, and the average causal effect as $E[Y^{cum(\bar{a})=2}] - E[Y^{cum(\bar{a})=0}]$. Again, any valid method should estimate that the value of this difference is 0.

Under the assumption that the mean outcome $E[Y|\bar{A}, L_1]$ depends linearly on the covariate $cum(\bar{A})$, we could fit the outcome regression model

$$E[Y|\bar{A}, L_1] = \theta_0 + \theta_1 cum(\bar{A}) + \theta_2 L_1$$

The associational difference $E[Y|cum(\bar{A}) = 2, L_1] - E[Y|cum(\bar{A}) = 0, L_1]$ is equal to $\theta_1 \times 2$. (The model correctly assumes that the difference is the same in the strata $L_1 = 1$ and $L_1 = 0$.) Therefore some might want to interpret $\theta_1 \times 2$

We invite readers to check for themselves that θ_1 is not zero by fitting this outcome regression model to the data in Table 20.1.

as the average causal effect of “always treat” versus “never treat” within levels of the covariate L_1 . But such causal interpretation is unwarranted because, as Figure 20.5 shows, conditioning on L_1 induces an association between A_0 , a component of treatment $cum(\bar{A})$, and the outcome Y . This implies that θ_1 —and therefore the associational difference of means—is non-zero even if the true causal effect is zero. A similar argument can be applied to matching. G-methods are needed to appropriately adjust for time-varying confounders in the presence of treatment-confounder feedback.

20.5 Adjusting for past treatment

One more thing before we discuss g-methods. For simplicity, we have so far described treatment-confounder feedback under simplified causal diagrams in which past treatment does not directly affect subsequent treatment. That is, the causal diagrams in Figures 20.3 and 20.4 did not include an arrow from A_0 to A_1 . We now consider the more general case in which past treatment may directly affect subsequent treatment.

As an example, suppose doctors in our HIV study use information on past treatment history \bar{A}_{k-1} when making a decision about whether to prescribe treatment A_k at time k . To represent this situation, we add an arrow from A_0 to A_1 to the causal diagrams in Figures 20.3 and 20.4, as depicted in Figures 20.8 and 20.9.

The causal diagrams in Figures 20.8 and 20.9 show that, in the presence of treatment-confounder feedback, conditioning on L_1 is insufficient to block all backdoor paths between treatment A_1 and outcome Y . Indeed conditioning on L_1 opens the path $A_1 \leftarrow A_0 \rightarrow L_1 \leftarrow U \rightarrow Y$ in Figure 20.8, and the path $A_1 \leftarrow A_0 \leftarrow W_0 \rightarrow L_1 \leftarrow U \rightarrow Y$ in Figure 20.9. Of course, regardless of whether treatment-confounder feedback exists, conditioning on past treatment history is always required when past treatment has a non-null effect of the outcome, as in the causal diagram of Figure 20.10. Under this diagram, treatment A_0 is a confounder of the effect of treatment A_1 .

Therefore, sequential exchangeability at time k generally requires conditioning on treatment history \bar{A}_{k-1} before k ; conditioning only on the covariates L is not enough. That is why, in this and in the previous chapter, all the conditional independence statements representing sequential exchangeability were conditional on treatment history.

Past treatment plays an important role in the estimation of effects of time-fixed treatments too. Suppose we are interested in estimating the effect of the time-fixed treatment A_1 —as opposed to the effect of a treatment strategy involving both A_0 and A_1 —on Y . (Sometimes the effect of A_1 is referred to as the short-term effect of the time-varying treatment \bar{A} .) Then lack of adjustment for past treatment A_0 will generally result in selection bias if there is treatment-confounder feedback, and in confounding if past treatment A_0 directly affects the outcome Y . In other words, the difference $E[Y|A_1 = 1, L_1] - E[Y|A_1 = 0, L_1]$ would not be zero even if treatment A_1 had no effect on any individual’s outcome Y , as in Figures 20.8-20.10. In practice, when making causal inferences about time-fixed treatments, bias may arise in analyses that compare current users ($A_1 = 1$) versus nonusers ($A_1 = 0$) of treatment. To avoid the bias, one can adjust for prior treatment history or to restrict the analysis to individuals with a particular treatment history. This is the idea behind “new-user designs” for time-fixed treatments: restrict the

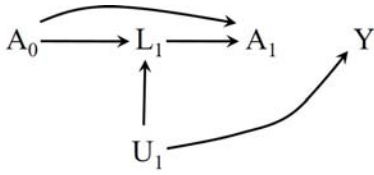


Figure 20.8

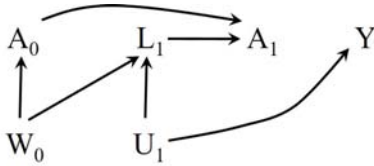


Figure 20.9

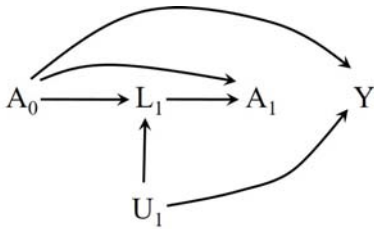


Figure 20.10

analysis to individuals who had not used treatment in the past.

The requirement to adjust for past treatment has additional bias implications when past treatment is mismeasured. As discussed in Section 9.3, a mismeasured confounder may result in effect estimates that are biased, either upwards or downwards. In our HIV example, suppose investigators did not have access to the study participants' medical records. Rather, to ascertain prior treatment, investigators had to ask participants via a questionnaire. Since not all participants provided an accurate recollection of their treatment history, treatment A_0 was measured with error. Investigators had data on the mismeasured variable A_0^* rather than on the variable A_0 . To depict this setting in Figures 20.8-20.10, we add an arrow from the true treatment A_0 to the mismeasured treatment A_0^* , which shows that conditioning on A_0^* cannot block the biasing paths between A_1 and Y that go through A_0 . Investigators will then conclude that there is an association between A_1 to Y , even after adjusting for A_0^* and L_1 , despite the lack of an effect on A_1 on Y . Therefore, contrary to a widespread belief, mismeasurement treatment may exaggerate effect estimates—even if the measurement error is independent non-differential—because of imperfect bias adjustment.

