

CAUSAL INFERENCE FOR CONTINUOUS-TIME PROCESSES WHEN COVARIATES ARE OBSERVED ONLY AT DISCRETE TIMES

BY MINGYUAN ZHANG, MARSHALL M. JOFFE¹ AND DYLAN S. SMALL²

University of Pennsylvania

Most of the work on the structural nested model and g-estimation for causal inference in longitudinal data assumes a discrete-time underlying data generating process. However, in some observational studies, it is more reasonable to assume that the data are generated from a continuous-time process and are only observable at discrete time points. When these circumstances arise, the sequential randomization assumption in the observed discrete-time data, which is essential in justifying discrete-time g-estimation, may not be reasonable. Under a deterministic model, we discuss other useful assumptions that guarantee the consistency of discrete-time g-estimation. In more general cases, when those assumptions are violated, we propose a controlling-the-future method that performs at least as well as g-estimation in most scenarios and which provides consistent estimation in some cases where g-estimation is severely inconsistent. We apply the methods discussed in this paper to simulated data, as well as to a data set collected following a massive flood in Bangladesh, estimating the effect of diarrhea on children's height. Results from different methods are compared in both simulation and the real application.

1. Introduction and motivation. In this paper, we study assumptions and methods for making causal inferences about the effect of a treatment that varies in continuous time when its time-dependent confounders are observed only at discrete times. Examples of settings in which this problem arises are given in Section 1.2. In such settings, standard discrete-time methods such as g-estimation usually do not work, except when certain conditions are assumed for the continuous-time process. In this paper, we formulate such conditions. When these conditions do not hold, we propose a controlling-the-future method which can produce consistent estimates when g-estimation is consistent and which is still consistent in some cases when g-estimation is severely inconsistent.

First, we review the approach of James Robins and collaborators to making causal inferences about the effect of a treatment that varies at discrete, observed times.

Received September 2009; revised April 2010.

¹Supported in part by NIH Grant CA095415.

²Supported in part by NSF Grant SES-0961971.

AMS 2000 subject classifications. Primary 62P10; secondary 62M99.

Key words and phrases. Causal inference, continuous-time process, deterministic model, diarrhea, g-estimation, longitudinal data, structural nested model.

1.1. *Review of Robins' causal inference approach for treatments varying at discrete, observed times.* In a cross-sectional observational study of the effect of a treatment on an outcome, a usual assumption for making causal inferences is that there are no unmeasured confounders, that is, that conditional on the measured confounders, the data is generated as if the treatment were assigned randomly. Under this assumption, a consistent estimate of the average causal effect of the treatment can be obtained from a correct model of the association between the treatment and the outcome conditional on the measured confounders [Cochran (1965)]. In a longitudinal study, the analog of the "no unmeasured confounders" assumption is that at the time of each treatment assignment, there are no unmeasured confounders; this is called the *sequential randomization* or *sequential ignorability* assumption, given as follows.

(A1) The longitudinal data of interest are generated as if the treatment is randomized in each period, conditional on the current values of measured covariates and the history of the measured covariates and the treatment.

The sequential randomization assumption implies that decision on treatment assignment is based on observable history and contemporaneous covariates, and that people have no ability to see into the future. Robins (1986) has shown that for a longitudinal study, unlike for a cross-sectional study, even if the sequential randomization assumption holds, the standard method of estimating the causal effect of the treatment by the association between the outcome and the treatment history conditional on the confounders can provide a biased and inconsistent estimate. This bias can occur when we are interested in estimating the joint effects of all treatment assignments and when the following conditions hold:

(c1) conditional on past treatment history, a time-dependent variable is a predictor of the subsequent mean of the outcome and also a predictor of subsequent treatment;

(c2) past treatment history is an independent predictor of the time-dependent variable.

Here, "independent predictor" means that prior treatment predicts current levels of the covariate, even after conditioning on other covariates. An example in which the standard methods are biased is the estimation of the causal effect of the drug AZT (zidovudine) on CD4 counts in AIDS patients. Past CD4 count is a time-dependent confounder for the effect of AZT on future CD4 count since it not only predicts future CD4 count, but also subsequent initiation of AZT therapy. Also, AZT history is an independent predictor of subsequent CD4 count [e.g., Hernán, Brumback and Robins (2002)].

To eliminate the bias of standard methods for estimating the causal effect of treatment in longitudinal studies where sequential randomization holds but there are time-dependent confounders satisfying conditions (c1) and (c2) (e.g., past CD4 counts), Robins (1986, 1992, 1994, 1998, 2000) developed a number of innovative

methods. We focus here on structural nested models (SNMs) and their associated methods of g-testing and g-estimation. The basic idea of the g-test is the following. Given a hypothesized treatment effect and a deterministic model of the treatment effect, we can calculate the *potential outcome* that a subject would have had if she never received the treatment. Such an outcome is also known as a *counterfactual outcome*, which is the outcome under a treatment history that might be contrary to the realized treatment history. If the hypothesized treatment effect is the true treatment effect, then this potential outcome will be independent of the actual treatment the subject received conditional on the confounder and treatment history, under the sequential randomization assumption (A1). g-estimation involves finding the treatment effect that makes the g-test statistic have its expected null value. For simplicity, our exposition focuses on deterministic rank-preserving structural nested distribution models; g-estimation also works for nondeterministic structural nested distribution models.

The SNM and g-estimation were developed for settings in which treatment decisions are being made at discrete times at which all the confounders are observed. In some settings, the treatment is varying in continuous time, but confounders are only observed at discrete times.

1.2. *Examples of treatments varying in continuous time where covariates are observed only at discrete times.*

EXAMPLE 1 (The effect of diarrhea on children's height). Diarrheal disease is one of the leading causes of childhood illness in developing regions [Kosek, Bern and Guerrant (2003)]. Consequently, there is considerable concern about the effects of diarrhea on a child's physical and cognitive development [Moore et al. (2001), Guerrant et al. (2002)]. A data set which provides the opportunity to study the impact of diarrhea on a child's height is a longitudinal household survey conducted in Bangladesh in 1998–1999 after Bangladesh was struck by its worst flood in over a century in the summer of 1998 [del Ninno et al. (2001), del Ninno and Lundberg (2005)]. The survey was fielded in three waves from a sample of 757 households: round 1 in November, 1998; round 2 in March–April, 1999; round 3 in November, 1999. The survey recorded all episodes of diarrhea for each child in the household in the past six months or since the last interview by asking the families at the time of each interview. In addition, the survey recorded at each of the three interview times several important time-dependent covariates for the effect of diarrhea on a child's future height: the child's current height and weight; the amount of flooding in the child's home and village; the household's economic and sanitation status. In particular, the child's current height and weight are time-dependent confounders that satisfy conditions (c1) and (c2), making standard longitudinal data analysis methods biased [see Martorell and Ho (1984) and Moore et al. (2001) for discussion of evidence for and reasons why current height and weight satisfy conditions (c1) and (c2)]. The time-dependent confounders of current height and

weight are available only at the time of the interview, and changes in their value that might affect the exposure of the child to the “treatment” of diarrhea, which varies in continuous time, are not recorded in continuous time.

EXAMPLE 2 [The effect of AZT (Zidovudine) on CD4 counts]. The Multi-center AIDS Cohort Study [MACS, Kaslow et al. (1987)] has been used to study the effect of AZT on CD4 counts [Hernán, Brumback and Robins (2002), Brumback et al. (2004)]. Participants in the study are asked to come semi-annually for visits at which they are asked to complete a detailed interview, including a complete history of their AZT use, as well as to take a physical examination. Decisions on AZT use are made by subjects and their physicians, and switches of treatment might happen at any time between two visits. These decisions are based on the values of diagnostic variables, possibly including CD4 and CD8 counts, and the presence of certain symptoms. However, these covariates are only measured by MACS at the time of visits; the values of these covariates at the exact times that treatment decisions are made between visits are not available.

1.3. *A model data generating process.* In both the examples of AZT and diarrhea, the exposure or treatment process happens continuously in time and a complete record of the process is available, but the time-dependent confounders are only observed at discrete times. There could be various interpretations of the relationship between the data at the treatment decision level and the data at the observational time level. To clarify the problem of interest in this paper, we consider a model data generating process that satisfies all of the following assumptions:

- (a1) a patient takes a certain medicine under the advice of a doctor;
- (a2) a doctor continuously monitors and records a list of health indicators of her patient and decides the initiation and cessation of the medicine solely based on current and historical records of these conditions, the historical use of the medicine and possibly random factors unrelated to the patient’s health;
- (a3) a third party organization asks a collection of patients from various doctors to visit the organization’s office semi-annually; the organization measures the same list of health indicators for the patients during their visits and asks the patients to report the detailed history of the use of the medicine between two visits;
- (a4) we are only provided with the third party’s data.

Note that in (a2), we assume the sequential randomization assumption (A1) at the treatment decision level.

The AZT example can be approximated by the above data generating process. In the AZT example, (a1) and (a2) approximately describe the joint decision-making process by the patient and the doctor in the real world. (a3) can be justified by reasonably assuming that the staff at the MACS receive similar medical training and use similar medical equipment as the patients’ doctors. In the diarrhea example,

the patient's body, rather than a doctor, determines whether the patient gets diarrhea. Assumption (a3), then, is saying that the third party organization (the survey organization) collects enough health data and that if all the histories of such health data are available, the occurrence of diarrhea is conditionally independent of the potential height.

1.4. *Difficulties posed by treatments varying in continuous time when covariates are observed only at discrete times.* Suppose our data are generated as in the previous section and we apply discrete-time g-estimation at the discrete times at which the time-dependent covariates are observed; we will denote these observation times by $0, \dots, K$. In discrete-time g-estimation, we are testing whether the observed treatment at time t ($t = 0, \dots, K$) is, conditional on the observed treatments at times $0, 1, \dots, t - 1$ and observed covariates at times $0, \dots, t$, independent of the *putative* potential outcomes at times $t + 1, \dots, K$, calculated under the hypothesized treatment effect, where the putative potential outcomes considered are what the subject's outcome would be at times $t + 1, \dots, K$ if the subject never received treatment at any time point. The difficulty with this procedure is that even if sequential randomization holds when the measured confounders are measured in continuous time [as is assumed in (a2)], it may not hold when the measured confounders are measured only at discrete times. For the discrete-time data, there can be *unmeasured confounders*. In the MACS example, the diagnostic measures at the time of AZT initiation are missing unless the start of AZT initiation occurred exactly at one of the discrete times that the covariates are observed; the diagnostic measures at the initiation time are clearly important confounders for the treatment status at the subsequent observational time. In the diarrhea example, the nutrition status of the child before the start of a diarrhea episode is missing unless the start of the diarrhea episode occurred exactly at one of the discrete times that covariates are observed; this nutrition status is also an important confounder for the diarrhea status at the subsequent observational time. Continuous-time sequential randomization does not, in general, justify sequential randomization holding for the discrete-time data, meaning that discrete-time g-estimation can produce inconsistent estimates, even when continuous-time sequential randomization holds.

In this paper, we approach this problem from two perspectives. First, we give conditions on the underlying continuous-time processes under which discrete-time sequential randomization is implied, warranting the use of discrete-time g-estimation. Second, we propose a new estimation method, called the *controlling-the-future method*, that can produce consistent estimates whenever discrete-time g-estimation is consistent and can produce consistent estimates in some cases where discrete-time g-estimation is inconsistent.

Our discussion focuses on a binary treatment and repeated continuous outcomes. We also assume that the cumulative amount of treatment between two visits is observed. This is true for Examples 1 and 2, the AZT and diarrhea studies,

respectively. If cumulative treatment is not observed, there will often be a measurement error problem in the amount of treatment, which is beyond the scope of this paper and an issue which we are currently researching.

The organization of the paper is as follows: Section 2 reviews the standard discrete-time structural nested model and g-estimation, describes a modified application when the underlying process is in continuous time and proposes conditions on the continuous-time processes when it works; Section 3 describes our controlling-the-future method; Section 4 presents a simulation study; Section 5 provides an application to the diarrhea study discussed in Example 1; Section 6 concludes the paper.

2. A modified g-estimation for discretely observed continuous-time processes. In this section, we first review the discrete-time structural nested model and the standard g-estimation, and mathematically formalize the setting we described in Section 1.3. Then, with a slight modification and different interpretation of notation, the g-estimation can be applied to the discrete-time observations from the continuous-time model. We will show that under certain conditions, this estimation method is consistent.

2.1. Review of discrete-time structural nested model and g-estimation. To reduce notation for the continuous-time setting, we use a star superscript on every variable in this section.

Assuming that all variables can only change values at time $0, 1, 2, \dots, K$, we use A_k^* to denote the binary treatment decision at time k . Under the discrete-time setup, A_k^* is assumed to be the constant level of treatment between time k and time $(k + 1)$. We use Y_k^{0*} to denote the baseline potential outcome of the study at time k if the subject does not receive any treatment throughout the study and Y_k^* to denote the actual outcome at time k . In this paper, we assume that all Y_k^{0*} 's and Y_k^* 's are continuous variables. Let L_k^* be the vector of covariates collected at time k . As a convention, Y_k^* is included in L_k^* .

We consider a simple deterministic model for the purposes of illustration,

$$(1) \quad Y_k^* = Y_k^{0*} + \Psi \sum_{i=0}^{k-1} A_i^*,$$

where Ψ is the causal parameter of interest and can be interpreted as the effect of one unit of the treatment on the outcome.

Model (1) is known as a *rank-preserving* model [Robins (1992)]. Under this model, for subjects i and j who have the same observed treatment history up to time k , if we observe $Y_{k,i} < Y_{k,j}$, then we must have $Y_{k,i}^{0*} < Y_{k,j}^{0*}$. It is also stronger than a more general rank-preserving model since Y_k^* depends deterministically only on Y_k^{0*} and the A_i^* 's.

Causal inference aims to estimate Ψ from the observables, the A_k^* 's and L_k^* 's. One way to achieve the identification of Ψ is to assume sequential randomization (A1). Given this notation and model (1), a mathematical formulation of (A1) is

$$(2) \quad P(A_k^* | \bar{L}_k^*, \bar{A}_{k-1}^*, \underline{Y}_{k+}^{0*}) = P(A_k^* | \bar{L}_k^*, \bar{A}_{k-1}^*),$$

where $\bar{L}_k^* = (L_0, L_1, \dots, L_k)$, $\bar{A}_{k-1}^* = (A_0, A_1, \dots, A_{k-1})$ and $\underline{Y}_{k+}^{0*} = (Y_{k+1}^{0*}, Y_{k+2}^{0*}, \dots, Y_K^{0*})$.

For any hypothesized value of Ψ , we define a putative potential outcome,

$$Y_k^{0*}(\Psi) = Y_k^* - \Psi \sum_{i=0}^{k-1} A_i.$$

Then, under (1) and (2), the correct Ψ should solve

$$(3) \quad E[U(\Psi)] \equiv E \left\{ \sum_{\substack{k < m \leq K \\ 1 \leq i \leq N}} [A_{i,k}^* - p_k(X_{i,k}^*)] g(Y_{i,m}^{0*}(\Psi), X_{i,k}^*) \right\} = 0,$$

where i is the index for each subject where there are N subjects, $X_{i,k}^* = (\bar{L}_{i,k}^*, \bar{A}_{i,k-1}^*)$, $p_k(X_{i,k}^*) = P(A_{i,k}^* = 1 | X_{i,k}^*)$ is the propensity score for subject i at time k and g is any function. This estimating equation can be generalized, with g being a function of any number of future $Y_{i,m}^{0*}(\Psi)$'s and $X_{i,k}^*$.

To estimate Ψ , we solve the empirical version of (3):

$$(4) \quad U(\Psi) \equiv \sum_{\substack{k < m \leq K \\ 1 \leq i \leq N}} [A_{i,k}^* - p_k(X_{i,k}^*)] g(Y_{i,m}^{0*}(\Psi), X_{i,k}^*) = 0.$$

If the true propensity score model is unknown and is parameterized as $p_k(X_k^*, \beta)$, additional estimating equations are needed to identify β . For example, the following estimating equations could be used:

$$(5) \quad U(\Psi, \beta) = \sum_{\substack{k < m \leq K \\ 1 \leq i \leq N}} [A_{i,k}^* - p_k(X_{i,k}^*)] [g(Y_{i,m}^{0*}(\Psi), X_{i,k}^*), X_{i,k}^*]^T = 0.$$

The method is known as *g-estimation*. The efficiency of the estimate depends on the functional form of g . The optimal g function that produces the most efficient estimation can be derived [Robins (1992)]. The formulas for estimating the covariance matrix of $(\hat{\Psi}, \hat{\beta})$ are given in Appendix A. A short discussion of the existence of the solution to the estimating equation and identification can be found in Appendix B.

2.2. *A continuous-time deterministic model and continuous-time sequential randomization.* We now extend the model in Section 2.1 to a continuous-time model and define a continuous-time version of the sequential randomization assumption (A1) as a counterpart of (2).

We now assume that the variables can change their values at any real time between 0 and K . The model in Section 2.1 is then extended as follows:

- $\{Y_t; 0 \leq t \leq K\}$ is the continuous-time, continuously-valued outcome process;
- $\{L_t; 0 \leq t \leq K\}$ is the continuous-time covariate process—it can be multi-dimensional and Y_t is an element of L_t ;
- $\{A_t; 0 \leq t \leq K\}$ is the continuous-time binary treatment process;
- $\{Y_t^0; 0 \leq t \leq K\}$ is the continuous-time, continuously-valued potential outcome process if the subject does not receive any treatment from time 0 to time K —it can be thought of as the *natural process* of the subject, free of treatment/intervention.

As a regularity condition, we further assume that all of the continuous-time stochastic processes are càdlàg processes (i.e., continuous from the right, having limits from the left) throughout this paper.

A natural extension of model (1) is

$$(6) \quad Y_t = Y_t^0 + \Psi \int_0^t A_s ds,$$

where Ψ is the causal parameter of interest. Ψ can be interpreted as the effect rate of the treatment on the outcome.

In this continuous-time model, a continuous-time version of the sequential randomization assumption (A1) or, equivalently, assumption (a2), can be formalized, although it does not have a simple form similar to equation (2). It was noted by Lok (2008) that a direct extension of the formula (2) involves “conditioning null events on null events.”

Lok (2008) formally defined continuous-time sequential randomization when there is only one outcome at the end of the study. We propose a similar definition for studies with repeated outcomes under the deterministic model (6).

Let $Z_t = (L_t, A_t, Y_t^0)$. Let $\sigma(Z_t)$ be the σ -field generated by Z_t , that is, the smallest σ -field that makes Z_t measurable. Let $\sigma(\bar{Z}_t)$ be the σ -field generated by $\bigcup_{u \leq t} \sigma(Z_u)$. Similarly, $\sigma(\bar{Z}_t, \underline{Y}_{t+}^0)$ is the σ -field generated by $\sigma(\bar{Z}_t) \cup \sigma(\underline{Y}_{t+}^0)$, where $\sigma(\underline{Y}_{t+}^0)$ is the σ -field generated by $\bigcup_{u > t} \sigma(Y_u^0)$. By definition, the sequence of $\sigma(\bar{Z}_t)$, $0 \leq t \leq K$, forms a filtration. The sequence of $\sigma(\bar{Z}_t, \underline{Y}_{t+}^0)$, $0 \leq t \leq K$, also forms a filtration because $\sigma(\bar{Z}_t, \underline{Y}_{t+}^0) \subset \sigma(\bar{Z}_s, \underline{Y}_{s+}^0)$ for $t < s$ [note that this is true under the deterministic model (6), but not in general].

Let N_t be a counting process determined by A_t . It counts the number of jumps in the A_t process. Let λ_t be a version of the intensity process of N_t with respect to $\sigma(\bar{Z}_t)$. $M_t = N_t - \int_0^t \lambda_s ds$ will be a martingale with respect to $\sigma(\bar{Z}_t)$.

DEFINITION 1. With N_t and M_t defined as above, the càdlàg process $Z_t \equiv (L_t, A_t, Y_t^0)$, $0 \leq t \leq K$, is said to satisfy the *continuous-time sequential randomization* assumption, or *CTSR*, if M_t is also a martingale with respect to $\sigma(\bar{Z}_t, \underline{Y}_{t+}^0)$. Or, equivalently, there exists a λ_t that is the intensity of N_t , with respect to both the filtration of $\sigma(\bar{Z}_t, \underline{Y}_{t+}^0)$ and the filtration of $\sigma(\bar{Z}_t)$.

In this definition, given A_0 , the counting process $\{N_t\}_0^T$ offers an alternative description of the treatment process $\{A_t\}_0^T$. The intensity process λ_t , which models the jumping rate of N_t , plays the same role as the propensity scores in the discrete-time model, which models the switching of the treatment process. Definition 1 formalizes assumption (A1) in the continuous-time model, by stating that λ_t does not depend on future potential outcomes.

The definition can be generalized if A_t has more than two levels, where N_t can be a multivariate counting process, each element counts a type of jump of the A_t process and λ_t is the multivariate intensity process for N_t under both the filtration of $\sigma(\bar{Z}_t)$ and the filtration of $\sigma(\bar{Z}_t, \underline{Y}_{t+}^0)$; see Lok (2008).

2.3. *A modified g-estimation.* In this paper, we assume that the continuous process defined in Section 2.2 can only be observed at integer times, namely, times $0, 1, 2, \dots, K$. We use the same starred notation as in Section 2.1, but interpret instances of this as discrete-time observations from the model in Section 2.2. Specifically:

- $\{A_k^*, k = 0, 1, 2, \dots, K\}$ denotes the set of treatment assignments observable at times $0, 1, 2, \dots, K$. We use \bar{A}_k^* to denote the observed history of observed discrete-time treatment up to time k , that is, $(A_0^*, A_1^*, \dots, A_k^*)$. Additionally, we use $\text{cum } A_k^* = \int_0^{k-} A_s ds$ to denote the cumulative amount of treatment up to time k . Note that in the continuous-time model, $\text{cum } A_k^* \neq \sum_{k'=0}^{k-1} A_{k'}^*$, as it would in discrete-time models. We let $\overline{\text{cum } A_k^*} = (\text{cum } A_1^*, \text{cum } A_2^*, \dots, \text{cum } A_k^*)$. We note that, in practice, people sometimes use $\tilde{A}_k^* = \text{cum } A_{k+1}^* - \text{cum } A_k^*$ as the treatment at time k when applying discrete-time g-estimation to discrete-time observational data. Under deterministic models, such use of g-estimation usually requires stronger conditions than the conditions discussed in this paper. Throughout this paper, we define the treatment at time k as A_k^* .

- We define L_k^* , the observed covariates at time k , to be L_{k-} , the left limit of L at time k , following the convention that in the discrete model, people usually assume that the covariates are measured just before the treatment decision at time k . Y_k^* and Y_k^{0*} are also defined as Y_{k-} and Y_{k-}^0 , respectively, following the same convention. \bar{L}_k^* denotes $(L_0^*, L_1^*, \dots, L_k^*)$, and \bar{Y}_k^* and \bar{Y}_k^{0*} are defined accordingly. $\underline{Y}_{k+}^{0*} = (Y_{k+1}^{0*}, Y_{k+2}^{0*}, \dots, Y_K^{0*})$.

With this notation and in the spirit of g-estimation, which controls all observed history in the propensity score model for the treatment, we propose the following

working estimating equation:

$$(7) \quad U(\Psi) \equiv \sum_{\substack{k < m \leq K \\ 1 \leq i \leq N}} [A_{i,k}^* - p_k(X_{i,k}^*)] g(Y_{i,m}^{0*}(\Psi), X_{i,k}^*) = 0,$$

where $X_{i,k}^*$ is the collection of $\bar{L}_{i,k}^*$, $\bar{A}_{i,k-1}^*$ and $\overline{\text{cum } A_{i,k}^*}$, $p_k(X_{i,k}^*) = P(A_{i,k}^* = 1 | X_{i,k}^*)$ and $Y_{i,m}^{0*}(\Psi) = Y_{i,m}^* - \Psi \overline{\text{cum } A_{i,k}^*}$.

In practice, $p_k(X_{i,k}^*)$ is unknown and has to be parameterized as $p_k(X_{i,k}^*; \beta)$, and we use different functions g to identify all of the parameters, as in Section 2.1. The covariance matrix of estimated parameters can be estimated as in Appendix A. A discussion of the existence of a solution and identification can be found in Appendix B.

The estimating equation has the same form as (4), except for two important differences. First, the propensity score model in this section conditions on the additional $\overline{\text{cum } A_{i,k}^*}$. In the discrete-time model of Section 2.1, $\overline{\text{cum } A_{i,k}^*}$ would be a transformed version of $\bar{A}_{i,k-1}^*$ and was redundant information. However, with continuous-time underlying processes, $\overline{\text{cum } A_{i,k}^*}$ provides new information on the treatment history. Second, the putative potential outcome $Y_{i,m}^{0*}(\Psi)$ is calculated by subtracting the $\overline{\text{cum } A_{i,k}^*}$ from $Y_{i,m}^*$, instead of $\sum_{l=0}^{k-1} A_{i,l}^*$. We will later refer to the g-estimation in this section as the *modified g-estimation* (although it is in the true spirit of g-estimation). The justification and limitation of using the modified g-estimation will be discussed in Section 2.4.

We refer to the g-estimation in Section 2.1 as *naive g-estimation* when it is applied to data from a continuous-time model. When the data come from a continuous-time model, the naive g-estimation can be severely biased, as we will show in our simulation study and the diarrhea application. One source of bias is a measurement error problem, $\sum_{l=0}^{k-1} A_{i,l}^*$ is not the correct measure of the treatment; another source of bias is that the important information $\overline{\text{cum } A_{i,k}^*}$ is not conditioned on in the propensity score. Although we would not expect researchers to use naive g-estimation when the true cumulative treatments are available, we present the simulation and real application results using this method as a reference to show how severely biased the estimates would be had we not known the true cumulative treatments and the measurement error problem had dominated.

2.4. Justification of the modified g-estimation. Given discrete-time observational data from continuous-time underlying processes, solving equation (7) provides an estimate for Ψ . For this Ψ estimate to be consistent, an analog to condition (2) is needed:

$$(8) \quad P(A_k^* | \bar{L}_k^*, \bar{A}_{k-1}^*, \overline{\text{cum } A_k^*}, \underline{Y}_{k+}^{*0}) = P(A_k^* | \bar{L}_k^*, \bar{A}_{k-1}^*, \overline{\text{cum } A_k^*}).$$

Condition (8) is a requirement on variables at observational time points. Its validity for a given study relies on how the data are collected, in addition to the

underlying continuous-time data generating process. It is not clear, without conditions on the underlying continuous-time data generating process, how one would go about collecting data in a way such that (8) would hold while the standard ignorability (2) is not true. Here, we will seek conditions at the continuous-time process level that imply condition (8) and hence justify the estimating equation (7). In particular, we consider two such conditions.

2.4.1. *Sequential randomization at any finite subset of time points.* Recall the data generating process described in Section 1.3. The third party organization periodically (e.g., semi-annually) collects the health data and treatment records of the patients. Suppose that a researcher thinks (8) holds for the time points at which the third party organization collects these data. If the time points have not been chosen in a special way to make (8) hold, then the researcher will often be willing to make the stronger assumption that (8) would hold for any finite subset of time points at which the third party organization chose to collect data. For example, for the diarrhea study, the survey was actually conducted in November, 1998, March–April, 1999 and November, 1999. If a researcher thought (8) held for these three time points, then she might be willing to assume that (8) should also hold if the survey was instead conducted in December, 1998, February, 1999, May, 1999 and October, 1999.

Before formalizing the researcher’s assumption on any finite subset of time points, we make the following observation.

PROPOSITION 2. *Under the deterministic model assumption (6), the propensity score has the following property:*

$$(9) \quad P(A_k^* = 1 | \bar{L}_k^*, \bar{A}_{k-1}^*, \overline{\text{cum } A_k^*}) = P(A_k^* = 1 | \bar{L}_k^*, \bar{A}_{k-1}^*, \bar{Y}_k^{0*}).$$

PROOF. Under the deterministic assumption (6) and the correct Ψ , $(\bar{L}_k^*, \bar{A}_{k-1}^*, \overline{\text{cum } A_k^*})$ is a one-to-one transformation of $(\bar{L}_k^*, \bar{A}_{k-1}^*, \bar{Y}_k^{0*})$. \square

Using Proposition 2, we state the sequential randomization assumption at any finite subset of time points as follows.

DEFINITION 3. A càdlàg process $Z_t \equiv (L_t, A_t, Y_t^0)$, $0 \leq t \leq K$, is said to satisfy the *finite-time sequential randomization* assumption, or *FTSR*, if, for any finite subset of time points, $0 \leq t_1 < t_2 < \dots < t_n < t_{n+1} < \dots < t_{n+l} \leq K$, we have

$$(10) \quad P(A_{t_n} | \bar{L}_{t_n-}, \bar{A}_{t_{n-1}}, \bar{Y}_{t_n-}^0, \underline{Y}_{t_n+}^0) = P(A_{t_n} | \bar{L}_{t_n-}, \bar{A}_{t_{n-1}}, \bar{Y}_{t_n-}^0),$$

where $\bar{L}_{t_n-} = (L_{t_1-}, L_{t_2-}, \dots, L_{t_n-})$, $\bar{A}_{t_{n-1}} = (A_{t_1}, A_{t_2}, \dots, A_{t_{n-1}})$, $\bar{Y}_{t_n-}^0 = (Y_{t_1-}^0, Y_{t_2-}^0, \dots, Y_{t_n-}^0)$ and $\underline{Y}_{t_n+}^0 = (Y_{t_{n+1}-}^0, Y_{t_{n+2}-}^0, \dots, Y_{t_{n+l}-}^0)$.

It should be noted that for the conditional densities in (9) and (10), and the conditional densities in the following sections, we always choose the version that is the ratio of joint density to marginal density.

The finite-time sequential randomization assumption clearly implies condition (2) and thus justifies the modified g-estimation equation (7). We have also proven a result that shows the relationship between the FTSR assumption and the CTSR assumption.

THEOREM 4. *If a continuous-time càdlàg process Z_t satisfies finite-time sequential randomization, then, under some regularity conditions, it will also satisfy continuous-time sequential randomization.*

PROOF. See Appendix C. The regularity conditions are also stated in Appendix C. \square

The result of Theorem 4 is natural. As mentioned in Section 1.4, the continuous-time sequential randomization does not imply FTSR because, in discrete-time observations, we do not have the full continuous-time history to control. To compensate for the incomplete data problem, some stronger assumption on the continuous-time processes must be made if identification is to be achieved.

2.4.2. A Markovian condition. Given the finite-time sequential randomization assumption described above, two important questions arise. First, Theorem 4 shows that the FTSR assumption is stronger than the continuous-time sequential randomization assumption. It is natural to ask how much stronger it is than the CTSR assumption. Second, the FTSR assumption, unlike the CTSR assumption (A1), is not an assumption on the data generating process itself and so it is not clear how to incorporate domain knowledge about the data generating process to justify it. Is there a condition at the data generating process level which will be more helpful in deciding whether g-estimation is valid?

We partially answer both questions in the following theorem.

THEOREM 5. *Assuming that the process (Y_t^0, L_t, A_t) satisfies the continuous-time sequential randomization assumption, and that the process (Y_{t-}^0, L_{t-}, A_t) is Markovian, for any time t and $t + s$, $s > 0$, we have*

$$(11) \quad P(A_t | L_{t-}, Y_{t-}^0, \underline{Y}_{t+}^0) = P(A_t | L_{t-}, Y_{t-}^0),$$

which implies the finite-time sequential randomization assumption. Here, $\underline{Y}_{t+}^0 = (Y_{t_1-}^0, Y_{t_2-}^0, \dots, Y_{t_n-}^0)$ and $t < t_1 < t_2 < \dots < t_n$.

PROOF. The proof can be found in the Appendix D. \square

The theory states that the Markov condition and the CTSR assumption together imply the FTSR condition. Therefore, they imply condition (2) and thus justify the modified g -estimation equation (7).

We make the following comments on the theorem.

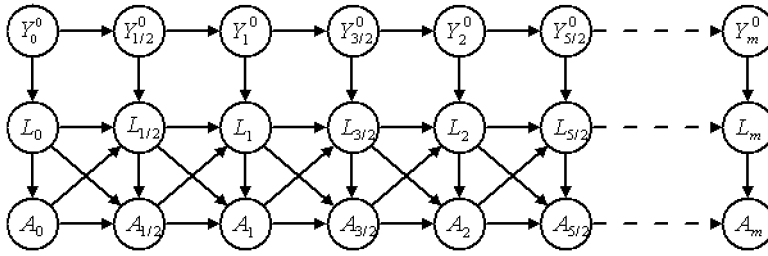
- The theorem partially answers our first question—the FTSR assumption is stronger than the CTSR assumption, but the gap between the two assumptions is less than a Markovian assumption. The result is not surprising since, with missing covariates between observational time points, we would hope that the variables at the observational time points well summarize the missing information. The Markovian assumption guarantees that variables at an observational time point summarize all information prior to that time point.

- The theorem also partially answers our second question. The CTSR assumption is usually justified by domain knowledge of how treatments are decided. Theorem 5 suggests that the researchers could further look for biological evidence that the process is Markovian to validate the use of g -estimation. The Markovian assumption can also be tested. One could first use the modified g -estimation to estimate the causal parameter, construct the Y^0 process at the observational time points and then test whether the full observational data of A, L, Y^0 come from a Markov process. A strict test of whether the discretely observed longitudinal data come from a continuous-time (usually nonstationary) Markov process could be difficult and is beyond the scope of this paper. As a starting point, we suggest Singer's trace inequalities [Singer (1981)] as a criterion to test for the Markovian property. A weaker test for the Markovian property is to test conditional independence of past observed values and future observed values conditioning on current observed values.

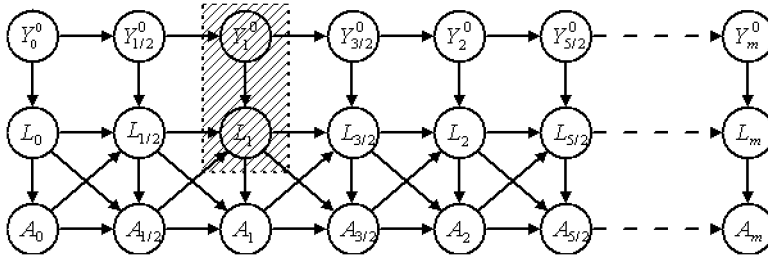
- In the theorem, equation (11) looks like an even stronger version of the continuous-time sequential randomization assumption—the treatment decision seems to be based only on current covariates and current potential outcomes. One could, of course, directly assume this stronger version of randomization and apply g -estimation. However, Theorem 5 is more useful since we are assuming a weaker untestable CTSR assumption and a Markovian assumption that is testable in principle.

- The theorem suggests that it is sufficient to control for current covariates and current potential outcomes for g -estimation to be consistent. In practice, we advise controlling for necessary past covariates and treatment history. The estimate would still be consistent if the Markovian assumption were true and it might reduce bias when the Markovian assumption was not true. As a result, we do control for previous covariates and treatments in our simulation and application to the diarrhea data.

- It is worth noting that the labeling of time is arbitrary. In practice, researchers can label whatever they have controlled for in their propensity score as the “current” covariates, which could include covariates and treatments that are



(a) DAG of a Markovian process.



(b) Verification of equation (9).

FIG. 1. Directed acyclic graph.

measured or assigned previously. In this case, the dimension of the process that needs to be tested for the Markovian property should also be expanded to include older covariates and treatments.

- Finally, we note that a discrete-time version of the theorem is implied by Corollary 4.2 of Robins (1997) if we set, in his notation, U_{ak} to be the covariates between two observational time points and U_{bk} to be the null set.

As a discretized example, we illustrate the idea of Theorem 5 by a directed acyclic graph (DAG) in part (a) of Figure 1, which assumes that all variables can only change values at time points $0, 1/2, 1, 3/2, 2, \dots, m$. Note that we do not distinguish the left limits of variables and the variables themselves in all DAGs of this paper, for reasons discussed in Appendix C. We also assume that the process can only be observed at times $0, 1, 2, \dots, m$. It is easy to verify that the DAG satisfies sequential randomization at the $0, 1/2, 1, 3/2, 2, \dots, m$ time level. The DAG is also Markovian in time. For example, if we control A_1, L_1, Y_1^0 , any variable prior to time 1 will be d-separated from any other variable after time 1.

Part (b) of Figure 1 verifies that A_1 is d-separated from $Y_m^0, m > 1$ by the shaded variables, namely, L_1 and Y_1^0 , as is implied by equation (9). By Theorem 5, the modified g-estimation works for data observed at the integer times if they are generated by the model defined by this DAG.

It is true that the Markovian condition that justifies the g-estimation equation (7) is restrictive, as will be discussed in the following section. However, our simulation

study shows that g-estimation has some level of robustness when the Markovian assumption is not seriously violated.

3. The controlling-the-future method. In this section, we consider situations in which the observational time sequential randomization fails and seek methods that are more robust to this failure than the modified g-estimation given in Section 2.3. The method we are going to introduce was proposed in Joffe and Robins (2009), which deals with a more general case of the existence of unmeasured confounders. It can be applied to deal with unmeasured confounders coming from either a subset of contemporaneous covariates or a subset of covariates that represent past time, the latter case being of interest for this paper. The method, which we will refer to as the controlling-the-future method (the reason for the name will become clear later on), gives consistent estimates when g-estimation is consistent and produces consistent estimates in some cases even when g-estimation is severely inconsistent.

In what follows, we will first describe an illustrative application of the controlling-the-future method and then discuss its relationship with our framework of g-estimation in continuous-time processes with covariates observed at discrete times.

3.1. *Modified assumption and estimation of parameters.* We assume the same continuous-time model as in Section 2.2. Following Joffe and Robins (2009), we consider a revised sequential randomization assumption on variables at the observational time points

$$(12) \quad P(A_k^* | \bar{L}_k^*, \bar{A}_{k-1}^*, \overline{\text{cum } A_k^*}, \underline{Y}_{k+}^{0*}) = P(A_k^* | \bar{L}_k^*, \bar{A}_{k-1}^*, \overline{\text{cum } A_k^*}, Y_{k+1}^{0*}).$$

This assumption relaxes (8). At each time point, conditioning on previous observed history, the treatment can depend on future potential outcomes, but only on the next period's potential outcome. In Joffe and Robins' extended formulation, this can be further relaxed to allow for dependence on more than one period of future potential outcomes, as well as other forms of dependence on the potential outcomes.

If the revised assumption (12) is true, then we obtain a similar estimating equation as (7). For each putative Ψ , we map Y_k^* to

$$Y_k^{0*}(\Psi) = Y_k^* - \Psi \text{cum } A_k^*,$$

the potential outcome if the subject never received any treatment under the hypothesized treatment effect Ψ .

Define the putative propensity score as

$$(13) \quad p_k(\Psi) \equiv P(A_k^* = 1 | \bar{L}_k^*, \bar{A}_{k-1}^*, \overline{\text{cum } A_k^*}, Y_{k+1}^{0*}(\Psi)).$$

Under assumption (12), the correct Ψ should solve

$$(14) \quad U(\Psi) = E \left\{ \sum_{\substack{1 \leq i \leq n \\ k+1 < m \leq K}} [A_{i,k}^* - p_{i,k}(\Psi)] g(Y_{i,m}^{0*}(\Psi), X_{i,k}^*, h_{i,k}(\Psi)) \right\} = 0,$$

where $X_{i,k}^* = (\bar{L}_{i,k}^*, \bar{A}_{i,k-1}^*, \overline{\text{cum } A_{i,k}^*})$, $h_{i,k}(\Psi) = Y_{i,k+1}^{0*}(\Psi)$ and g is any function and can be generalized to functions of $X_{i,k}^*$, $h_{i,k}(\Psi)$ and any number of future potential outcomes that are later than time $k + 1$, for example, $g(Y_{i,k+2}^{0*}(\Psi), Y_{i,k+3}^{0*}(\Psi), X_{i,k}^*, h_{i,k}(\Psi))$. In most real applications, the model for $p_k(\Psi) = E[A_k^* | X_k^*, h_k(\Psi)]$ is unknown and is usually estimated by a parametric model,

$$p_{i,k}(\Psi; \beta_X, \beta_h) = E[A_{i,k} | X_{i,k}^*, h_{i,k}(\Psi); \beta_X, \beta_h].$$

We can solve the following set of estimating equations to obtain the estimates of Ψ , β_X and β_h :

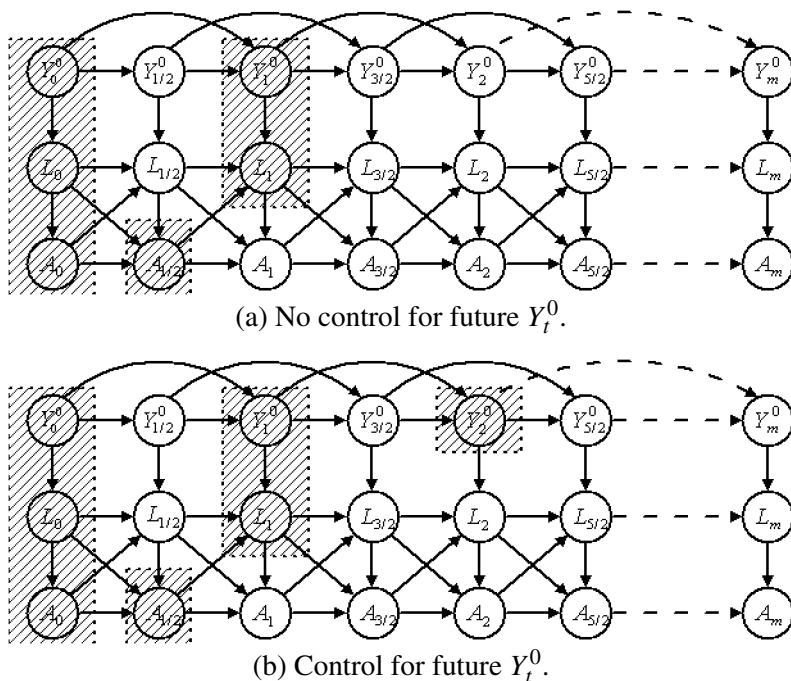
$$(15) \quad \begin{aligned} U(\Psi, \beta_X, \beta_h) &= \sum_{\substack{1 \leq i \leq n \\ k+1 < m \leq K}} (A_{i,k}^* - p_{i,k}(\Psi; \beta_X, \beta_h)) \\ &\quad \times [g(Y_{i,m}^{0*}(\Psi), X_{i,k}^*, h_{i,k}(\Psi)), X_{i,k}^*, h_{i,k}(\Psi)]^T \\ &= 0. \end{aligned}$$

The estimation of the covariance matrix of Ψ , β_X and β_h is similar to the usual standard g-estimation, which is described in Appendix A.

Two important features of estimating equation (15) distinguish it from estimating equation (7). First, in (15), there is a common parameter Ψ in both p_k 's model and $Y_m^{0*}(\Psi)$, caused by the fact that the treatment depends on a future potential outcome. Second, in (15), the sum over m and k is restricted to $m > k + 1$, while in (7), we only need $m > k$. If we use $m = k + 1$ in (15), $E\{[A_{i,k}^* - p_{i,k}(\Psi)]g(Y_{i,k+1}^{0*}(\Psi), X_{i,k}^*, h_{i,k}(\Psi))\} = 0$ usually does not lead to the identification of Ψ , unless certain functional forms of the propensity score model are assumed to be true [see Joffe and Robins (2009)].

3.2. *The controlling-the-future method and the Markovian condition.* Joffe and Robins' revised assumption (12) is an assumption on the discrete-time observational data. It relaxes the observational time sequential randomization (8) because (8) always implies (12). At the continuous-time data generating level, (12) allows less stringent underlying stochastic processes than the Markovian process in Theorem 5.

In particular, we identify two important scenarios where the relaxation happens. One scenario is to allow for more direct temporal dependence for the Y^0 process, which we will refer to as the *non-Markovian- Y^0* case. The other scenario is to allow colliders in L , which we will refer to as the *leading-indicator-in- L* case. We illustrate both cases by modifying the directed acyclic graph (DAG) example in Figure 1.

FIG. 2. Directed acyclic graph with non-Markovian Y_t^0 .

THE NON-MARKOVIAN- Y^0 CASE. Assume, for example, our data is generated from the DAG in Figure 2, where we allow the dependence of Y_2^0 on Y_1^0 , even if $Y_{3/2}^0$ is controlled. In part (a) of Figure 2, we control for observed covariates (L_0, L_1), treatment ($A_0, A_{1/2}$) and current and historical potential outcome (Y_0^0, Y_1^0) for treatment at time 1 (A_1), that is, we have controlled for all historically observed covariates, treatment and cumulative treatment as suggested in the comments accompanying Theorem 5. In this case, the modified g-estimation fails because the paths like $A_1 \leftarrow L_{1/2} \leftarrow Y_{1/2}^0 \rightarrow Y_{3/2}^0 \rightarrow Y_2^0 \rightarrow \dots \rightarrow Y_m^0$ are not blocked by the shaded variables. In part (b) of Figure 2, we control for the additional Y_2^0 . A_1 is not completely blocked from Y_m^0 , but some paths that are not blocked in part (a) are now blocked, for example, the path of $A_1 \leftarrow L_{1/2} \leftarrow Y_{1/2}^0 \rightarrow Y_{3/2}^0 \rightarrow Y_2^0 \rightarrow Y_{5/2}^0 \rightarrow \dots \rightarrow Y_m^0$. Also, no additional paths are opened by conditioning on Y_2^0 . We would usually expect that the correlation between A_1 and Y_m^0 is weakened. Under the framework of Joffe and Robins (2009), we can control for more than one period of future potential outcomes and expect to further weaken the correlation between A_1 and Y_m^0 . A modification of assumption (12) that conditions on more future potential outcomes may be approximately true.

The scenario relates to real-world problems. For instance, in the diarrhea example, Y_t^0 is the natural height growth of a child without any occurrence of diarrhea.

Height in the next month not only depends on the current month's height, but also depends on the previous month's height: the complete historical growth curve of the child provides information on genetics and nutritional status, and provides information about future natural height beyond that of current natural height alone. Therefore, the potential height process for the child is not Markovian. [For a formal argument why children's height growth is not Markovian, see Gasser et al. (1984).] By the reasoning employed above, g-estimation fails. However, if we assume that the delayed dependence of natural height wanes after a period of time (as in Figure 2), controlling for the next period potential height in the propensity score model might weaken the relationship between current diarrhea exposure and future potential height later than the next period and the assumptions of the controlling-the-future method might hold approximately.

THE LEADING-INDICATOR-IN- L CASE. In Figure 1, we do not allow any arrows from future Y^0 to previous L , which means that among all measures of the subject, there are no elements in L that contain any leading information about future Y^0 . This means that Y^0 is a measure that is ahead of all other measures, by which we mean that, for example, $L_2 \perp Y_m^0 | \bar{Y}_2^0, \bar{A}_{2-}, \bar{L}_{2-m} > 2$. This is not realistic in many real-world problems. In the example of the effect of the diarrhea on height, weight is an important covariate. While both height and weight reflect the nutritional status of a child, malnutrition usually affects weight more quickly than height, that is, the weight contains leading information for the natural height of the child. Figure 1 is thus not an appropriate model for studying the effect of diarrhea on height.

In Figure 3, we allow arrows from $Y_{1/2}^0$ to L_0 , from Y_1^0 to $L_{1/2}$ and so on, which assumes that L contains leading indicators of Y^0 , but the leading indicators are only ahead of Y^0 for less than one unit of time. Part (a) of Figure 3 shows that controlling for history of covariates, treatment and potential outcomes does not block A_1 from Y_m^0 . On the path of $A_1 \leftarrow L_{1/2} \rightarrow L_1 \leftarrow Y_{3/2}^0 \rightarrow Y_2^0 \rightarrow Y_{5/2}^0 \rightarrow \dots \rightarrow Y_m^0$, L_1 is a controlled collider. However, in part (b), if we do control for Y_2^0 additionally, the same path will be blocked. In general, if we assume that there exist leading indicators in covariates and that the leading indicators are not ahead of potential outcomes for more than one time unit, g-estimation will fail, but the controlling-the-future method will produce consistent estimates.

The fact that the controlling-the-future method can work in the leading information scenario can also be related to the discussion of Section 3.6 of Rosenbaum (1984). The main reason for g-estimation's failure in the DAG example is that $L_{1/2}$ is not observable and cannot be controlled. If $L_{1/2}$ is observed, it is easy to verify that the DAG in Figure 3 satisfies sequential randomization on the finest time grid. The idea behind the controlling-the-future method is to condition on a "surrogate" for $L_{1/2}$. The surrogate should satisfy the property that Y_m^0 is independent of the unobserved $L_{1/2}$ given the surrogate and other observed covariates [similar to formula 3.17 in Rosenbaum (1984)]. In the leading information case, when

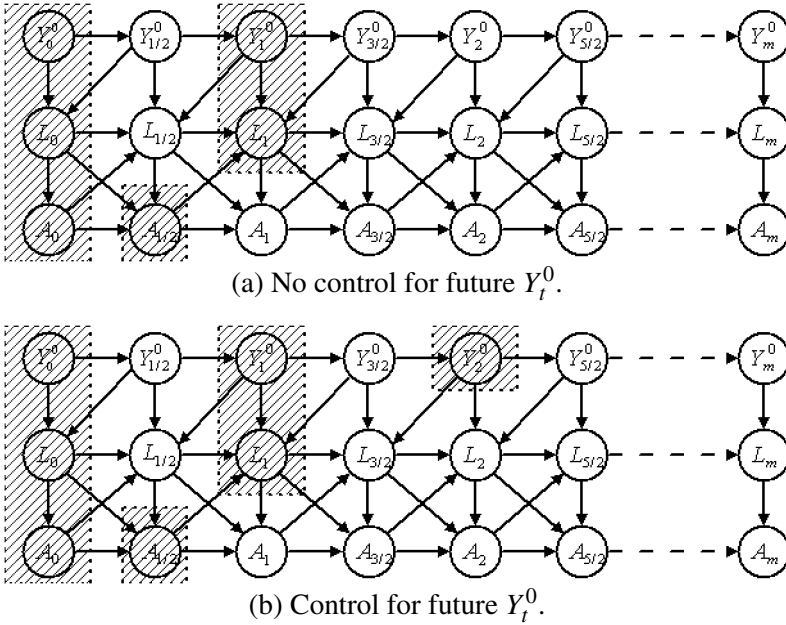


FIG. 3. Directed acyclic graph with leading indicator in L_t .

$m > k + 1$ and we have covariates \bar{L}_k that are only ahead of the potential outcome until time at most $k + 1$, the future potential outcome Y_{k+1}^0 is a surrogate. It is easy to check that in Figure 3, $L_{1/2}$ is independent of Y_m^0 , given Y_2^0, L_1, A_0 and cum A_1 (equivalently, Y_1^0).

It is worth noting that we do not need to control for anything except Y_2^0 in Figure 3 in order to get a consistent estimate. It is possible to construct more complicated DAGs in which controlling for additional past and current covariates is necessary, which involves more model specifications for the relationships among different covariates and deviates from the main point of this paper.

In Section 4, we will simulate data in cases of non-Markovian- Y_t^0 and leading-indicator-in- L_t , respectively, and show that the controlling-the-future method does produce better estimates than g-estimation. However, it is worth noting that when the modified g-estimation in Section 2.3 is consistent, the controlling-the-future estimation is usually considerably less efficient. This is because condition (12) is less stringent than (8). The semiparametric model under (8) is a submodel of the semiparametric model defined by (12). The latter will have a larger semiparametric efficiency bound than the former. Theoretically, the most efficient g-estimation will be more efficient than the most efficient controlling-the-future estimation if the g-estimation is valid. In practice, even if we are not using the most efficient estimators, controlling-the-future estimation usually estimates more parameters, for

example, coefficients for $h_{i,k}(\Psi)$ in the propensity model, and thus is less efficient. For a formal discussion, see Tsiatis (2006).

4. Simulation study. We set up a simple continuous-time model that satisfies sequential ignorability in continuous time, and simulate and record discrete-time data from variations of the simple model. We estimate causal parameters from both the modified g-estimation and the controlling-the-future estimation. We also present the estimates from naive g-estimation in Section 2.1, where we ignore the continuous-time information of the treatment processes, as a way to show the severity of the bias in the presence of the measurement error problem. The results support the discussions in Sections 2.4 and 3.

In the simulation models below, M1 satisfies the Markovian condition in Theorem 5. It also serves as a proof that there exist processes satisfying the conditions of Theorem 5.

4.1. *The simulation models.* We first consider a continuous-time Markov model which satisfies the CTSR assumption.

- Y_t^0 is the potential outcome process if the patient is not receiving any treatment. We assume that

$$Y_t^0 = g(V, t) + e_t,$$

where $g(V, t)$ is a function of baseline covariates V and time t . Let $g(V, t)$ be continuous in t and let e_t follow an Ornstein–Uhlenbeck process, that is,

$$de_t = -\theta e_t dt + \sigma dW_t,$$

where W_t is the standard Brownian motion.

- Y_t is the actual outcome process and follows the deterministic model (6):

$$Y_t = Y_t^0 + \Psi \int_0^t A_s ds.$$

- A_t is the treatment process, taking binary values. The jump of the A_t process follows the following formula:

$$P(A_s \text{ jumps once from } (t, t+h] | \bar{A}_t, \bar{Y}_t, \bar{Y}^0) = s(A_t, Y_t)h + o(h),$$

$$P(A_s \text{ jumps more than once from } (t, t+h] | \bar{A}_t, \bar{Y}_t, \bar{Y}^0) = o(h),$$

where \bar{A}_t and \bar{Y}_t are the full continuous-time history of treatment and outcome up to time t and \bar{Y}^0 is the full continuous-time path of potential outcome from time 0 to time K . By making $s(\cdot)$ independent of \bar{Y}^0 , we make our model satisfy the continuous-time sequential randomization assumption.

In this model, the only time-dependent confounder is the outcome process itself.

We also consider several variations of the above model (denoted as M1 below):

- Model (M2) extends (M1) to the non-Markovian- Y_t^0 case. Specifically, we consider the case where e_t in the model of Y_t^0 follows a non-Markovian process, namely an Ornstein–Uhlenbeck process in random environments, which is defined as the following:

(1) J_t is a continuous-time Markov process taking values in a finite set $\{1, \dots, m\}$, which is the environment process;

(2) we have $m > 1$ sets of parameters $\theta_1, \sigma_1, \dots, \theta_m, \sigma_m$;

(3) e_t follows an Ornstein–Uhlenbeck process with parameters θ_j, σ_j , when $J_t = j$; the starting point of each diffusion is chosen to be simply the endpoint of the previous one.

- Model (M3) extends (M1) to another setting of non-Markovian- Y_t^0 process, where

$$Y_t^0 = g(V, t) + 0.8e_{t-1} + 0.2e_t.$$

e_t follows the same Markovian Ornstein–Uhlenbeck process as in M1. Every other variable is the same as in M1.

- Model (M4) considers the case with more than one covariate. In M4, we keep the assumptions on Y_t^0 as in (M1) and the deterministic model of Y_t . We add one more covariate, which is generated as follows:

$$L_t^- = 0.2Y_t + 0.8Y_{t+0.5}^0 + 0.5\eta_t.$$

η_t follows an Ornstein–Uhlenbeck process independent of the Y_t^0 process. In this specification, the covariate L_t^- contains some leading information about Y^0 , but it is only ahead of Y^0 for 0.5 length of a time unit. Here, we use L_t^- instead of L_t to denote that it is the covariate excluding Y_t . The simulation model for the A_t process is given in Appendix E.

In all of these models, to simulate data, we use $g(V, t) = C$ (a constant), $\Psi = 1$, a time span from 0 to 5 and a sample size of 5000. Details of other parameter specifications can be found in Appendix E. We generate 5000 continuous paths of Y_t and A_t (and L_t^- in M4), from time 0 to time 5, and record $Y_0^*, A_0^*, Y_1^*, A_1^*, \dots, Y_4^*, A_4^*, Y_5^*$ and $\text{cum } A_1^*, \dots, \text{cum } A_5^*$ (and $L_0^{-*}, \dots, L_4^{-*}$ in M4) as the observed data.

4.2. *Estimations and results under M1.* Figure 4 shows a typical continuous-time path of Y_t^0 , Y_t and A_t . The treatment switches around time 0.7 and time 2.8.

We apply three estimating methods on data simulated from M1: the naive discrete-time g-estimation described in Section 2.1, which ignores the underlying continuous-time processes; the modified g-estimation described in Section 2.3, which controls for all the observed discrete-time history; and the controlling-the-future method in Section 3.1 of controlling for the next period’s potential outcome in addition to the discrete-time history.

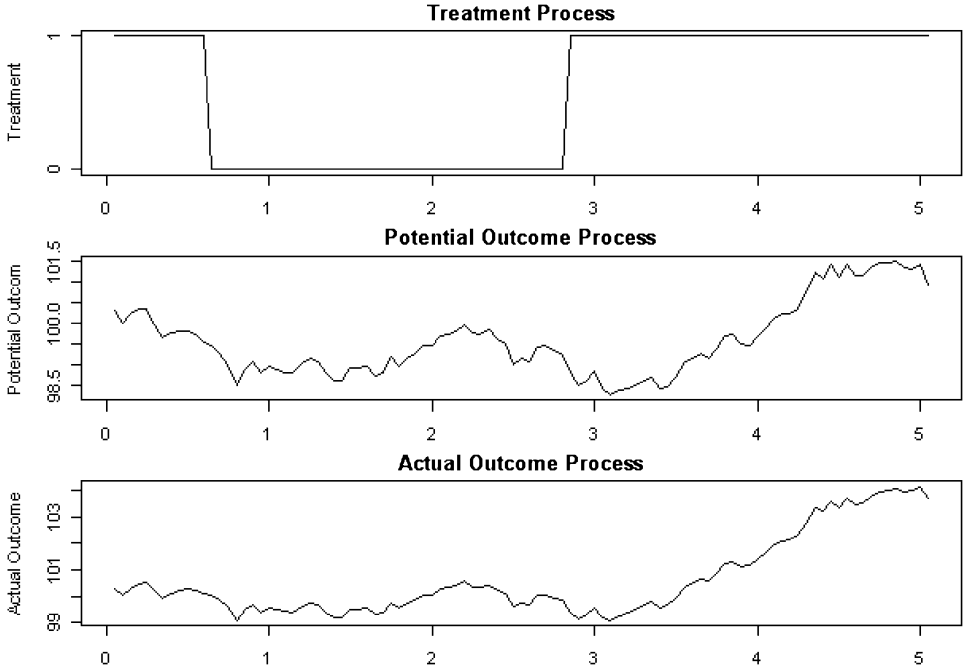


FIG. 4. Example of continuous-time paths under M1.

For estimation, even though we know the data generating process, it is too complicated to use the correct model for the propensity score, that is, the correct functional form for $p_k(\Psi) \equiv P(A_k^* | \bar{L}_k^*, \bar{Y}_k^*, \bar{A}_{k-1}^*, \overline{\text{cum } A_k^*}, Y_{k+}^{0*}(\Psi))$. Therefore, we use the following approximations (note that we control for past treatment and covariates as well—see comments for Theorem 5):

- (1) standard g-estimation ignoring continuous-time processes (naive g-estimation)

$$\text{logit}(p_k) = \beta_0 + \beta_1 A_{k-1}^* + \beta_2 Y_{k-1}^* + \beta_3 Y_k^*;$$

- (2) g-estimation controlling for all observed history (modified g-estimation)

$$\text{logit}(p_k) = \beta_0 + \beta_1 A_{k-1}^* + \beta_2 Y_{k-1}^* + \beta_3 Y_k^* + \beta_4 \text{cum } A_k^*;$$

- (3) the controlling-the-future method, controlling for next period potential outcomes (controlling-the-future estimation)

$$\text{logit}(p_k(\Psi)) = \beta_0 + \beta_1 A_{k-1}^* + \beta_2 Y_{k-1}^* + \beta_3 Y_k^* + \beta_4 \text{cum } A_k^* + \beta_5 Y_{k+1}^{0*}(\Psi).$$

We plug these models for the propensity scores into estimation equations (5), (7) and (15) respectively. [Note that in equation (5), $Y_k^{0*}(\Psi) = Y_k^* - \Psi \sum_{l=0}^{k-1} A_l^*$, while in the other two, $Y_k^{0*}(\Psi) = Y_k^* - \Psi \text{cum } A_k^*$.]

TABLE 1
Estimated causal parameters from data generated by M1–M4

	Naive g-est.	Mod. g-est.	Ctr-future est.
Simulation results from M1, <i>true parameter = 1</i>			
Mean estimate [†]	0.7728	1.0005	0.9988
S.D. of estimates [‡]	0.0183	0.0191	0.0403
S.D. of the mean estimate [*]	0.0005	0.0006	0.0013
Absolute bias ^{**}	0.2272	0.0005	0.0012
Coverage [◇]	0	0.946	0.956
Simulation results from M2, <i>true parameter = 1</i>			
Mean estimate [†]	0.7651	1.0016	1.0000
S.D. of estimates [‡]	0.0132	0.0158	0.0371
S.D. of the mean estimate [*]	0.0004	0.0005	0.0012
Absolute bias ^{**}	0.2349	0.0016	0.0000
Coverage [◇]	0	0.953	0.950
Simulation results from M3, <i>true parameter = 1</i>			
Mean estimate [†]	0.7580	0.9845	1.0026
S.D. of estimates [‡]	0.0149	0.0180	0.0487
S.D. of the mean estimate [*]	0.0005	0.0006	0.0015
Absolute bias ^{**}	0.2420	0.0155	0.0026
Coverage [◇]	0	0.855	0.956
Simulation results from M4, <i>true parameter = 1</i>			
Mean estimate [†]	0.7816	1.0853	1.0085
S.D. of estimates [‡]	0.0201	0.0289	0.0806
S.D. of the mean estimate [*]	0.0006	0.0009	0.0025
Absolute bias ^{**}	0.2184	0.0853	0.0085
Coverage [◇]	0	0.115	0.948

[†] Averaged over estimates from 1000 independent simulations of sample size 5000.

[‡] Sample standard deviation of the 1000 estimates.

^{*} Sample S.D./ $\sqrt{1000}$.

^{**} Absolute value of (1-mean estimates).

[◇] Coverage rate of 95% confidence intervals for 1000 simulations.

The first panel of Table 1 shows a summary of the estimates of causal parameters for 1000 simulations from M1. The naive g-estimation gives severely biased estimates. Controlling for all observed history and controlling for additional next period potential outcome both give us unbiased estimates. As discussed at the end of Section 3.2, the controlling-the-future method has lower efficiency.

The last row of the first panel in Table 1 shows the coverage rate of the 95% confidence interval estimated from the 1000 independent simulations. Naive g-estimation has a zero coverage rate, while the other two methods have coverage rates around 95%.

4.3. *Simulation results under M2 and M3.* The results in the second panel of Table 1 are typical for different values of parameters under M2. The naive g-estimation performs badly, while both of the other methods still work well with the data generated from M2. This shows that the modified g-estimation and the controlling-the-future method have some level of robustness to mild violations of the Markovian assumption.

The third part of Table 1 shows the results of simulation from M3, where Y^0 violates the Markov property more substantially. In this case, we can see that the mean of the modified g-estimates is biased, but the mean of the controlling-the-future estimates is almost unbiased. In the last row of the third panel, the coverage rate for the modified g-estimation drops to 0.855, while the controlling-the-future method still has a coverage rate of 0.956.

4.4. *Estimations and results under M4.* In M4, we create a covariate L_t^- that has leading information about Y_t^0 . In the data simulated from M4, the observational time sequential randomization (8) no longer holds, although the data are generated following continuous-time sequential randomization. This simulation serves as a numerical proof of the claim that continuous-time sequential randomization does not imply discrete-time sequential randomization.

To show this, we consider the following working propensity score model at time $k = 2$ and its dependence on the future potential outcome at $m = 4$:

- not controlling for the next period potential outcome (used in modified g-estimation)

$$\begin{aligned}
 & \text{logit}(P(A_k^* = 1 | \bar{A}_{k-}^*, \bar{L}_k^{-*}, \overline{\text{cum } A_k^*}, \bar{Y}_k^*, Y_m^{0*})) \\
 (16) \quad & = \beta_0 + \beta_1 \text{cum } A_k^* + \beta_2 L_{k-1}^{-*} + \beta_3 L_k^{-*} + \beta_4 A_{k-1}^* \\
 & + \beta_5 Y_k^* + \beta_6 Y_{k-1}^* + \beta_8 Y_m^{0*};
 \end{aligned}$$

- controlling for the next period potential outcome (used in controlling-the-future estimation)

$$\begin{aligned}
 & \text{logit}(P(A_k^* = 1 | \bar{A}_{k-}^*, \bar{L}_k^{-*}, \overline{\text{cum } A_k^*}, \bar{Y}_k^*, Y_{k+1}^{0*}, Y_m^{0*})) \\
 (17) \quad & = \beta_0 + \beta_1 \text{cum } A_k^* + \beta_2 L_{k-1}^{-*} + \beta_3 L_k^{-*} + \beta_4 A_{k-1}^* \\
 & + \beta_5 Y_k^* + \beta_6 Y_{k-1}^* + \beta_7 Y_{k+1}^{0*} + \beta_8 Y_m^{0*}.
 \end{aligned}$$

We can use the true values of Y_{k+1}^{0*} and Y_m^{0*} in the regression to test the discrete-time ignorability since we are simulating the data. Table 2 shows the estimates of β_7 and β_8 in both regression models. The result shows that the coefficient of Y_m^{0*} , β_8 , is significant if we do not control for the future potential outcome and is not significant if we control for the future potential outcome. This shows that observational time sequential randomization (8) does not hold, while the revised assumption (12) holds.

TABLE 2
Verification of observational time sequential randomization under M4

	Reg. model (16)*	Reg. model (17)*
β_7		0.1868
p -value		5.56e-05
β_8	0.0936	0.0134
p -value	0.0006	0.691

*Simulation sample size = 10,000.

The estimation results from M4 appear in the fourth panel of Table 1. In applying these methods, we use the following propensity score models separately:

(1) g-estimation ignoring the underlying continuous-time processes (naive g-estimation)

$$\text{logit}(p_k) = \beta_0 + \beta_1 A_{k-1}^* + \beta_2 Y_{k-1}^* + \beta_3 Y_k^* + \beta_5 L_{k-1}^{-*} + \beta_6 L_k^{-*};$$

(2) g-estimation controlling for all observed history (modified g-estimation)

$$\text{logit}(p_k) = \beta_0 + \beta_1 A_{k-1}^* + \beta_2 Y_{k-1}^* + \beta_3 Y_k^* + \beta_4 \text{cum } A_k^* + \beta_5 L_{k-1}^{-*} + \beta_6 L_k^{-*};$$

(3) the controlling-the-future method controlling for next period potential outcomes (controlling-the-future estimation)

$$\begin{aligned} \text{logit}(p_k(\Psi)) = & \beta_0 + \beta_1 A_{k-1}^* + \beta_2 Y_{k-1}^* + \beta_3 Y_k^* + \beta_4 \text{cum } A_k^* \\ & + \beta_5 L_{k-1}^{-*} + \beta_6 L_k^{-*} + \beta_7 Y_{k+1}^{0*}(\Psi). \end{aligned}$$

Both the naive g-estimation and the modified g-estimation give us estimates with severe bias and they have coverage rates of 0 and 0.115, respectively, for the 95% confidence interval constructed from them. It is worth noting that model 3 is misspecified, but, nevertheless, leads to much less biased estimates, and the controlling-the-future method has a coverage rate of 0.948.

5. Application to the diarrhea data. In this section, we apply the different approaches to the diarrhea example mentioned in Section 1 (Example 2). For illustration purposes, we ignore any informative censoring and use a set of 224 children with complete records between ages 3 and 6 from 757 households in Bangladesh around 1998. The outcomes, Y_k^* , are the heights of the children in centimeters, measured at round k of the interviews, for $k = 1, 2, 3$. The treatment A_k^* at the interview k is defined as $A_k^* = 1$ if the child was sick with diarrhea during the past two weeks of the interview and $A_k^* = 0$ otherwise. The cumulative treatment $\text{cum } A_k^*$ is the number of days that the child suffered from diarrhea from four months before the first interview (July 15th, 1998) to the k th interview. Baseline covariates V

include age in months, mother’s height and whether the household was exposed to the flood. Time-dependent covariates other than the outcome, that is, L_k^{-*} , include mid-upper arm circumference, weight for age z-score, type of toilet (open place, fixed place, unsealed toilet, water-sealed toilet or other), garbage disposal method (throwing away in own fixed place, throwing away in own nonfixed place, disposing anywhere or other method), water purifying process (filter, filter and broil, or other) and source of cooking water (from pond or river/canal, or from tube well, ring well or supply water).

We apply naive g-estimation, modified g-estimation and the controlling-the-future method to this data set. Since we only have three rounds, the actual propensity score models and the estimating equations for the three methods are as follows. Note that these estimating equations are for illustrative purpose and may not be the most efficient estimating equations for this data set.

- Naive g-estimation uses the following propensity score model:

$$\text{logit}\{P[A_k^* = 1|V, L_k^{-*}, Y_k^*]\} = \beta_0 + \beta_V V + \beta_L L_k^{-*} + \beta_Y Y_k^*,$$

where $k = 1, 2$.

The estimating equations follow the form of (5) in Section 2.1:

$$\sum_{\substack{1 \leq k < m \leq 3 \\ 1 \leq i \leq n}} [A_{k,i}^* - P(A_{k,i}^* = 1|V_i, L_{k,i}^{-*}, Y_{k,i}^*)] \begin{pmatrix} lY_{m,i}^{0*}(\Psi) \\ V_i \\ L_{k,i}^{-*} \\ Y_{k,i}^* \end{pmatrix} = 0,$$

where $Y_{m,i}^{0*}(\Psi) = Y_{m,i}^* - \Psi \sum_{l=1}^{m-1} A_l$.

- Modified g-estimation uses this propensity score model:

$$\begin{aligned} \text{logit}\{P[A_k^* = 1|V, L_k^{-*}, Y_k^*, \text{cum} A_k^*]\} \\ = \beta_0 + \beta_V V + \beta_L L_k^{-*} + \beta_Y Y_k^* + \beta_{\text{cum} A} \text{cum} A_k^*, \end{aligned}$$

where $k = 1, 2$.

The estimating equations follow the form of (7) in Section 2.3.

$$\sum_{\substack{1 \leq k < m \leq 3 \\ 1 \leq i \leq n}} [A_{k,i}^* - P(A_{k,i}^* = 1|V_i, L_{k,i}^{-*}, Y_{k,i}^*, \text{cum} A_{k,i}^*)] \begin{pmatrix} Y_{m,i}^{0*}(\Psi) \\ V_i \\ L_{k,i}^{-*} \\ Y_{k,i}^* \\ \text{cum} A_{k,i}^* \end{pmatrix} = 0,$$

where $Y_{m,i}^{0*}(\Psi) = Y_{m,i}^* - \Psi \text{cum} A_m$.

- Controlling-the-future estimation uses the following propensity score model:

$$\begin{aligned} \text{logit}\{P[A_1^* = 1|V, L_1^{-*}, Y_1^*, \text{cum} A_1^*, Y_2^{0*}(\Psi)]\} \\ = \beta_0 + \beta_V V + \beta_L L_1^{-*} + \beta_Y Y_1^* + \beta_{\text{cum} A} \text{cum} A_1^* + \beta_{Y^0} Y_2^{0*}(\Psi). \end{aligned}$$

The estimating equations follow (15) in Section 3:

$$\sum_{1 \leq i \leq n} [A_{1,i}^* - P(A_{1,i}^* | V_i, L_{1,i}^{-*}, Y_{1,i}^*, \text{cum } A_{1,i}^*, Y_{2,i}^{0*}(\Psi))] \begin{pmatrix} Y_{3,i}^{0*}(\Psi) \\ V_i \\ L_{1,i}^{-*} \\ Y_{1,i}^* \\ \text{cum } A_{1,i}^* \\ Y_{2,i}^{0*}(\Psi) \end{pmatrix} = 0,$$

where $Y_{3,i}^{0*}(\Psi) = Y_{3,i}^* - \Psi \text{cum } A_3$.

The interpretation of Ψ in the last two models is that one day of suffering from diarrhea reduces the height of the child by Ψ centimeters. For naive g-estimation, the underlying data generating model treats the exposure at the observational time as the constant exposure level for the next six months, which does not make sense in the context. It should be noted that if we apply the naive g-estimation, the estimated Ψ should not be interpreted the same way in the modified g-estimation and the controlling-the-future method. Instead, it be interpreted as the effect of having diarrhea at the time of visits. The effect of the child having diarrhea at any time between the visit and the next visit six months later, but not at the time of the visit, is not described by this Ψ .

The estimating equations are solved by a Newton–Raphson algorithm. The estimated Ψ and its standard deviation are reported in Table 3. Modified g-estimation estimates $\hat{\Psi} = -0.3481$, which means that the height of the child is reduced by 0.35 cm if the child has one day of diarrhea. Our controlling-the-future method produces an estimate of $\hat{\Psi} = -0.0840$. Although all of the estimates are not significant because of the small sample size, the sign and magnitude of the estimate from the controlling-the-future method are similar to what has been found in other research on diarrhea’s effect on height [e.g., Moore et al. (2001)].

In addition, we note that the standard deviation of the modified g-estimate is higher than that of the controlling-the-future estimate. As discussed at the end of Section 3.2, if the modified g-estimation is consistent, we would expect the controlling-the-future estimation to have larger standard deviation. The standard deviations in Table 3 provide evidence that the modified g-estimation is not consistent.

TABLE 3
Estimation of Ψ from the diarrhea data set

Method	Estimate	Std. err.
Naive g-est.	-0.3991	0.2469
Modified g-est.	-0.3481	0.2832
Controlling-the-future est.	-0.0840	0.1894

6. Conclusion. In this paper, we have studied causal inference from longitudinal data when the underlying processes are in continuous time, but the covariates are only observed at discrete times. We have investigated two aspects of the problem. One is the validity of the discrete-time g-estimation. Specifically, we investigated a modified g-estimation that is in the spirit of standard discrete-time g-estimation, but is modified to incorporate the information of the underlying continuous-time treatment process, which we have referred to as “modified g-estimation” throughout the paper. We have shown that an important condition that justifies this modified g-estimation is the finite-time sequential randomization assumption at any subset of time points, which is strictly stronger than the continuous-time sequential randomization. We have also shown that a Markovian assumption and the continuous-time sequential randomization would imply the FTSR assumption. The Markovian condition is more useful than the FTSR assumption, in the sense that it can potentially help researchers decide whether the application of the modified g-estimation is appropriate. The other aspect is the controlling-the-future method that we propose to use when the condition to warrant g-estimation does not hold. The controlling-the-future method can produce consistent estimates when g-estimation is inconsistent and is less biased in other scenarios. In particular, we identified two important cases in which controlling the future is less biased, namely, when there is delayed dependence in the baseline potential outcome process and when there are leading indicators of the potential outcome process in the covariate process.

In our simulation study, we have shown the performance of the modified g-estimation and the controlling-the-future estimation. The results confirm our discussion in earlier sections. The simulation results also indicate the danger of applying naive g-estimation, which is usually severely biased and inconsistent when its underlying assumptions are violated, as in the situations considered.

We have applied the g-estimation methods and the controlling-the-future method to estimating the effect of diarrhea on a child’s height and estimated that its effect is negative but not significant. The real application also provides some evidence that the modified g-estimation is not consistent.

All of the discussion in this paper is based on a particular form of causal model—equation (6). However, all of the arguments could apply to a class of more general rank-preserving models, with necessary adjustments in various equations. If we assume a generic rank-preserving model with $Y_t = f(Y_t^0, h(\bar{A}_{t-}); \Psi)$, where \bar{A}_{t-} is the continuous-time path of A from time 0 to $t-$, h is some functional [e.g., in our paper, $h(\bar{A}_{t-}) = \int_0^t A_s ds$] and f is some strictly monotonic function with respect to the first argument [e.g., in our paper, $f(x, y; \Psi) = x + \Psi y$], we map Y_k^* to $Y_k^{0*} = f^{-1}(Y_k^*, h(\bar{A}_{k-}); \Psi)$, where f^{-1} is the inverse of $f(x, y; \Psi)$ with respect to x for any given y . We can then substitute all cum A_k^* ’s in this paper by the $h(\bar{A}_{k-})$ ’s. All of the discussion and formulas in the paper would remain valid under the assumption that we observe all $h(\bar{A}_{k-})$ ’s, which can be easily satisfied with detailed continuous-time records of the treatment. It should be noted

that the argument does not work if a time-varying covariate modifies the effect of treatment. For example, if $Y_t = Y_t^0 + \Psi \int_0^t L_s^2 A_s ds$, where L_s is a time-varying covariate, observing the full continuous-time treatment process is not enough. Some imputation for the L_s process is necessary.

The methods considered here have several limitations. These include rank preservation, a strong assumption that the effects of treatment are deterministic. This assumption facilitates the interpretation of models. In other work on structural nested distribution and related models [e.g., Robins (2008)], rank preservation has been shown to be unnecessary in settings in which one is not modeling the joint distribution of potential outcomes under different treatments. We expect that this is also the case here, and work justifying this more formally is in progress. We also require that the cumulative amount of treatment (or the full continuous-time treatment process, if using other causal models mentioned above) between the discrete time points when the covariates are observed is known. Work is in progress on the more challenging case in which the treatment process is only observed at discrete times and the cumulative amount of treatment is measured with error. In addition, we ignore any censoring problem requiring that our data is complete, which might not be satisfied in reality. It will also be interesting to study how to accommodate censored data in our framework in future work.

APPENDIX A: ESTIMATING COVARIANCE MATRIX OF ESTIMATED PARAMETERS

The formulas in this appendix can be used to estimate the covariance matrix of the estimated parameters from naive g-estimation of Section 2.1, modified g-estimation of Section 2.3 and the controlling-the-future estimation of Section 3.1. More general results on the asymptotical covariances can be found in van der Vaart (2000).

We write $\theta = (\Psi, \beta)$. In Sections 2.1 and 2.3, β is the parameter in the propensity score model. In Section 3.1, $\beta = (\beta_X, \beta_h)$ is the parameter in the propensity score model. Let $U(\theta)$ be the vector on the left-hand side of the estimating equations [equation (5) in Section 2.1, equation (7) in Section 2.3 and equation (15) in Section 3.1, respectively]. We also define

$$U_{i,k,m}(\theta) \equiv (A_{i,k}^* - p_{i,k}(\beta)) [g(Y_{i,m}^{0*}(\Psi), X_{i,k}^*), X_{i,k}^*]^T$$

for the naive g-estimation and the modified g-estimation, and

$$U_{i,k,m}(\theta) \equiv (A_{i,k}^* - p_{i,k}(\Psi; \beta_X, \beta_h)) [g(Y_{i,m}^{0*}(\Psi), X_{i,k}^*, h_i), X_{i,k}^*, h_{i,k}]^T$$

for the controlling-the-future estimation. We then have $U(\theta) = \sum U_{i,k,m}$.

Let $B(\theta) = E[\frac{\partial U(\theta)}{\partial \theta}]$, which can be estimated as

$$\hat{B}(\theta) = - \sum_{i,k,m} \left\{ \frac{\partial U_{i,k,m}}{\partial \theta} \right\} \Big|_{\theta=\hat{\theta}}$$

where $\hat{\theta}$ is the solution from the corresponding estimating equations, $k < m$ in both g-estimations and $k < m - 1$ in controlling-the-future estimation. The covariance matrix of the estimator $\hat{\theta}$ can then be estimated as

$$\text{Cov}(\hat{\theta}) = \hat{B}^{-1}(\theta) \text{Cov}[U(\hat{\theta})] \hat{B}^{-1}(\theta)'$$

by the delta method, where $\text{Cov}[U(\theta)]$ is estimated by

$$\text{Cov}[U(\hat{\theta})] = \sum_i U_i(\hat{\theta}) U_i(\hat{\theta})'$$

with $U_i = \sum_{k,m} U_{i,k,m}(\hat{\theta})$, $k < m$ in both g-estimations and $k < m - 1$ in controlling-the-future estimation.

APPENDIX B: EXISTENCE OF SOLUTION AND IDENTIFICATION

The estimating equations in this paper, equation (5) in Section 2.1, equation (7) in Section 2.3 and equation (15) in Section 3.1, are asymptotically consistent systems of equations by definition, if the respective underlying assumptions for each estimating equation hold true. The existence of a solution is guaranteed asymptotically. In addition, we have the same number of equations as the number of parameters in each system. One would usually expect there to exist a solution for the estimating equations, even in a relatively small sample.

However, the asymptotic solution may not be unique, which leads to an identification problem. As a special case from the more general semi-parametric theory [see Tsiatis (2006)], we state the following lemma for identification, following the notation of Appendix A.

LEMMA 6. *The parameter θ is identifiable under the model*

$$E[U(\theta)] = 0$$

if both $\text{Cov}[U(\theta_0)]$ and $B(\theta_0) \equiv E[\frac{\partial U(\theta_0)}{\partial \theta}]$ are of full rank. Here, θ_0 is the value of the true parameter.

PROOF. The proof is trivial. By Appendix A, the asymptotic covariance matrix of the estimates is given by

$$B^{-1}(\theta_0) \text{Cov}[U(\theta_0)] B^{-1}(\theta_0)',$$

which will be finite and of full rank when the conditions in the lemma hold true. \square

APPENDIX C: PROOF THAT FTSR IMPLIES CTSR

We assume that Z_t is a càdlàg process, and everything we discuss is in an a.s. sense.

We first define

$$\begin{aligned}\mathcal{H}_{t-} &\equiv \sigma(\bar{Z}_{t-}), \\ \mathcal{F}_{t-,t+} &\equiv \sigma(\bar{Z}_{t-}, \underline{Y}_{t+}^0).\end{aligned}$$

Recall that N_t counts the number of jumps in A_t up to time t . We assume that a continuous version of the $\mathcal{F}_{t-,t+}$ intensity process of N_t exists, which we denote by η_t . If we define

$$r_t(\delta) = (1 - A_{t-})A_{t+\delta} + A_{t-}(1 - A_{t+\delta}).$$

Then, under certain regularity conditions [see Chapter 2 of Andersen et al. (1992)], for every t ,

$$\eta_t = \lim_{\delta \downarrow 0} \frac{E[r_t(\delta) | \mathcal{F}_{t-,t+}]}{\delta} \quad \text{a.s.}$$

For Theorem 4, we need to show that η_t is also \mathcal{H}_{t-} -measurable. This is because if this is true, then

$$\begin{aligned}& E \left[N_{t_0+s} - \int_0^{t_0+s} \eta_t dt \middle| \mathcal{H}_{t_0} \right] \\ &= E \left[\left(N_{t_0+s} - \int_0^{t_0+s} \eta_t dt \right) - \left(N_{t_0} - \int_0^{t_0} \eta_t dt \right) \middle| \mathcal{H}_{t_0} \right] \\ &\quad + E \left[\left(N_{t_0} - \int_0^{t_0} \eta_t dt \right) \middle| \mathcal{H}_{t_0} \right] \\ &= E \left\{ E \left[\left(N_{t_0+s} - \int_0^{t_0+s} \eta_t dt \right) - \left(N_{t_0} - \int_0^{t_0} \eta_t dt \right) \middle| \mathcal{F}_{t_0-,t_0+} \right] \middle| \mathcal{H}_{t_0} \right\} \\ &\quad + N_{t_0} - \int_0^{t_0} \eta_t dt \\ &= 0 + N_{t_0} - \int_0^{t_0} \eta_t dt \\ &= N_{t_0} - \int_0^{t_0} \eta_t dt.\end{aligned}$$

The second equality follows because of properties of conditional expectation and the assumption that η_t is \mathcal{H}_{t-} -measurable. The third equality holds because η_t is an $\mathcal{F}_{t-,t+}$ intensity process of N_t . The last equality shows that η_t is also a \mathcal{H}_{t-} intensity process of N_t , which agrees with the definition of CTSR.

Before proving the main result, we assume the following regularity conditions.

1. As stated before, we assume that η_t is continuous. We further assume that η_t is positive, and bounded from below and above by constants that do not depend on t . We also assume that $\frac{E[r_t(\delta)]\mathcal{F}_{t-t+\delta}}{\delta}$ is bounded by a constant for every t within a interval of $(0, \delta_0]$.

2. We assume that for any finite sequence of time points, $t_1 \leq t_2 \leq t_3 \leq \dots \leq t_n$, the density $f(Z_{t_1} = z_1, Z_{t_2} = z_2, \dots, Z_{t_n} = z_n)$ is well defined and locally uniformly bounded, that is, there exists a constant D and a rectangle $B \equiv [t_1 - \delta_1, t_1 + \delta_1] \times [t_2 - \delta_2, t_2 + \delta_2] \times \dots \times [t_n - \delta_n, t_n + \delta_n]$ such that for any $(t'_1, t'_2, \dots, t'_n)^T \in B$ and any possible value of $(z_1, z_2, \dots, z_n)^T$,

$$f(Z_{t'_1} = z_1, Z_{t'_2} = z_2, \dots, Z_{t'_n} = z_n) \leq D.$$

For any conditional expectation involving finite sequence of time points, we choose the version that is defined by the joint density.

3. Given any finite sequence of time points, $t_1 \leq t_2 \leq t_3 \leq \dots \leq t_n$ and any possible value of $(z_1, z_2, \dots, z_n)^T$, we assume that the following convergence is uniform in a closed neighborhood of $\tilde{t} \equiv (t_1, t_2, t_3, \dots, t_n)$:

$$\begin{aligned} & f(Z_{t'_1} = z_1, Z_{t'_2} = z_2, \dots, Z_{t'_n} = z_n) \\ &= \lim_{\Delta \downarrow 0} \frac{P(Z_{t'_1} \in [z_1, z_1 + \Delta_1], Z_{t'_2} \in [z_2, z_2 + \Delta_2], \dots, Z_{t'_n} \in [z_n, z_n + \Delta_n])}{\Delta_1 \times \Delta_2 \times \dots \times \Delta_n}, \end{aligned}$$

where $(t'_1, t'_2, \dots, t'_n)^T$ is in a neighborhood of \tilde{t} .

4. Given any finite sequence of time points, $t_1 \leq t_2 \leq t_3 \leq \dots \leq t_i \leq \dots \leq t_n$ and any possible value of $(z_1, z_2, \dots, z_n)^T$, we define

$$f(\delta) = \frac{P(A_{t_i+\delta} \neq A_{t_i} | Z_{t_1} = z_1, Z_{t_2} = z_2, \dots, Z_{t_n} = z_n)}{\delta}.$$

We assume that $\lim_{\delta \downarrow 0} f(\delta)$ exists and is positive and finite. We also assume that $f(\delta)$ is finite and is right-continuous in δ , and the continuity is uniform with respect to (δ, t_i) in $[0, \delta_0] \times B(t_i)$, where $B(t_i)$ is a closed neighborhood of t_i . Further, we assume that the above assumption is true if any of the Z in f is in its left-limit value rather than the concurrent value.

REMARK 7. The third regularity condition is needed when we want to prove convergence in density. For example, consider that when $\delta \downarrow 0$, we have $Z_{t_2+\delta} \rightarrow Z_{t_2}$. We can then see that

$$\begin{aligned} & \lim_{\delta \downarrow 0} f(Z_{t_1} = z_1, Z_{t_2+\delta} = z_2, Z_{t_3} = z_3) \\ &= \lim_{\delta \downarrow 0} \lim_{\substack{\Delta_1 \downarrow 0 \\ \Delta_2 \downarrow 0 \\ \Delta_3 \downarrow 0}} \frac{P(Z_{t_1} \in [z_1, z_1 + \Delta_1], Z_{t_2+\delta} \in [z_2, z_2 + \Delta_2], Z_{t_3} \in [z_3, z_3 + \Delta_3])}{\Delta_1 \Delta_2 \Delta_3} \end{aligned}$$

$$\begin{aligned}
&= \lim_{\substack{\Delta_1 \downarrow 0 \\ \Delta_2 \downarrow 0 \\ \Delta_3 \downarrow 0}} \lim_{\delta \downarrow 0} \frac{P(Z_{t_1} \in [z_1, z_1 + \Delta_1], Z_{t_2+\delta} \in [z_2, z_2 + \Delta_2], Z_{t_3} \in [z_3, z_3 + \Delta_3])}{\Delta_1 \Delta_2 \Delta_3} \\
&= \lim_{\substack{\Delta_1 \downarrow 0 \\ \Delta_2 \downarrow 0 \\ \Delta_3 \downarrow 0}} \frac{P(Z_{t_1} \in [z_1, z_1 + \Delta_1], Z_{t_2} \in [z_2, z_2 + \Delta_2], Z_{t_3} \in [z_3, z_3 + \Delta_3])}{\Delta_1 \Delta_2 \Delta_3} \\
&= f(Z_{t_1} = z_1, Z_{t_2} = z_2, Z_{t_3} = z_3).
\end{aligned}$$

The interchanging of limits in the second equality is valid because of the third regularity condition. The third equality follows from the fact that probabilities are expectations of indicator functions and that the dominated convergence theorem applies.

We introduce the following lemma for technical convenience.

LEMMA 8. *If the càdlàg process Z_t follows the finite-time sequential randomization as defined in Definition 3, then the following version of FTSR is also true:*

$$\begin{aligned}
(18) \quad &P(A_{t_n} | \bar{L}_{t_{n-1}}, L_{t_n-}, \bar{A}_{t_{n-1}}, \bar{Y}_{t_{n-1}}^0, Y_{t_n-}^0, \underline{Y}_{t_n+}^0) \\
&= P(A_{t_n} | \bar{L}_{t_{n-1}}, L_{t_n-}, \bar{A}_{t_{n-1}}, \bar{Y}_{t_{n-1}}^0, Y_{t_n-}^0),
\end{aligned}$$

where $\bar{L}_{t_{n-1}} = (L_{t_1}, L_{t_2}, \dots, L_{t_{n-1}})$, $\bar{A}_{t_{n-1}} = (A_{t_1}, A_{t_2}, \dots, A_{t_{n-1}})$, $\bar{Y}_{t_{n-1}}^0 = (Y_{t_1}^0, Y_{t_2}^0, \dots, Y_{t_{n-1}}^0)$ and $\underline{Y}_{t_n+}^0 = (Y_{t_{n+1}}^0, Y_{t_{n+2}}^0, \dots, Y_{t_{n+t}}^0)$.

REMARK 9. The difference between (18) and the original definition of FTSR is that in (18), most L 's and Y^0 's are stated in their concurrent values, while in Definition 3, they are all stated in their left limits. Lemma 8 is only for technical convenience.

PROOF OF LEMMA 8. The result follows directly from the definition of a càdlàg process. \square

We now consider a discrete-time property.

LEMMA 10. *Suppose FTSR holds true. If we define*

$$\begin{aligned}
\mathcal{F} &= \sigma(Z_{t_1}, \dots, Z_{t_{n-1}}, Z_{t-}, Y_{t_{n+1}}^0, \dots, Y_{t_{n+t}}^0), \\
\mathcal{H} &= \sigma(Z_{t_1}, \dots, Z_{t_{n-1}}, Z_{t-}),
\end{aligned}$$

then we have for every t that

$$(19) \quad \lim_{\delta \downarrow 0} \frac{E[r_t(\delta) | \mathcal{F}]}{\delta} = \lim_{\delta \downarrow 0} \frac{E[r_t(\delta) | \mathcal{H}]}{\delta} \quad a.s.$$

PROOF. First, we note that the limits on both sides of equation (19) exist and are finite. This fact follows from the regularity condition 1. Take $\lim_{\delta \downarrow 0} \frac{E[r_t(\delta)|\mathcal{F}]}$, for example:

$$\begin{aligned} \lim_{\delta \downarrow 0} \frac{E[r_t(\delta)|\mathcal{F}]}{\delta} &= \lim_{\delta \downarrow 0} \frac{E[E[r_t(\delta)|\sigma(\bar{Z}_{t-}, \underline{Y}_t^0)]|\mathcal{F}]}{\delta} \\ &= E\left[\lim_{\delta \downarrow 0} \frac{E[r_t(\delta)|\sigma(\bar{Z}_{t-}, \underline{Y}_t^0)]}{\delta} \middle| \mathcal{F}\right] \\ &= E[\eta_t|\mathcal{F}]. \end{aligned}$$

The interchange of limit and expectation is guaranteed by the assumption in regularity condition 1 that $\frac{E[r_t(\delta)|\sigma(\bar{Z}_{t-}, \underline{Y}_t^0)]}{\delta}$ is bounded. The existence is then guaranteed by the dominated convergence theorem and $E[\eta_t|\mathcal{F}]$ is obviously finite.

Given equation (10) and Lemma 8, we always have

$$(20) \quad E[I_{A_t \neq A_{t_n}}[\bar{L}_{t-}, \bar{A}_{t_n}, \bar{Y}_{t-}^0, \underline{Y}_{t+}^0]] = E[I_{A_t \neq A_{t_n}}[\bar{L}_{t-}, \bar{A}_{t_n}, \bar{Y}_{t-}^0]] \quad \text{a.s.},$$

where $\bar{L}_{t-} = (L_{t_1}, L_{t_2}, \dots, L_{t_{n-1}}, L_{t_n}, L_{t-})^T$, $\bar{A}_{t_n} = (A_{t_1}, A_{t_2}, \dots, A_{t_n})^T$, $\bar{Y}_{t-}^0 = (Y_{t_1}^0, Y_{t_2}^0, \dots, Y_{t_n}^0, Y_{t-}^0)^T$ and $\underline{Y}_{t+}^0 = (Y_{t_{n+1}}^0, Y_{t_{n+2}}^0, \dots, Y_{t_{n+l}}^0)^T$.

In the regularity conditions, since we assumed the existence of joint density, the usual definition of conditional probability is a version of the conditional expectation defined using σ -fields. In our case, we have

$$\begin{aligned} &\lim_{\delta \downarrow 0} \frac{E[r_t(\delta)|\mathcal{F}]}{\delta} \\ &= \lim_{\delta \downarrow 0} \frac{P(A_{t+\delta} \neq A_{t-} | Z_{t_1}, \dots, Z_{t_{n-1}}, Z_{t-}, Y_{t_{n+1}}^0, \dots, Y_{t_{n+l}}^0)}{\delta} \\ &= \lim_{\delta \downarrow 0} \lim_{t_n \uparrow t-} \frac{P(A_{t+\delta} \neq A_{t_n} | Z_{t_1}, \dots, Z_{t_{n-1}}, Z_{t_n}, L_{t-}, Y_{t-}^0, Y_{t_{n+1}}^0, \dots, Y_{t_{n+l}}^0)}{\delta + (t - t_n)} \\ &= \lim_{t_n \uparrow t-} \lim_{\delta \downarrow 0} \frac{P(A_{t+\delta} \neq A_{t_n} | Z_{t_1}, \dots, Z_{t_{n-1}}, Z_{t_n}, L_{t-}, Y_{t-}^0, Y_{t_{n+1}}^0, \dots, Y_{t_{n+l}}^0)}{\delta + (t - t_n)} \\ &= \lim_{t_n \uparrow t-} \frac{P(A_t \neq A_{t_n} | Z_{t_1}, \dots, Z_{t_{n-1}}, Z_{t_n}, L_{t-}, Y_{t-}^0, Y_{t_{n+1}}^0, \dots, Y_{t_{n+l}}^0)}{t - t_n} \\ &= \lim_{t_n \uparrow t-} \frac{E[I_{A_t \neq A_{t_n}}[\bar{L}_{t-}, \bar{A}_{t_n}, \bar{Y}_{t-}^0, \underline{Y}_{t+}^0]]}{t - t_n}. \end{aligned}$$

The second equality is guaranteed by the third regularity condition. By Remark 7, we can show that the conditional density in the third line converges to the second line as A_{t_n} and Z_{t_n} converges to A_{t-} and Z_{t-} . The $(t - t_n)$ term in the denominator is not needed for the second equality, but is crucial for the interchangeability of

limits in the third equality. The interchangeability of limits is guaranteed by the fourth regularity condition. By the fourth regularity condition, the following limit

$$\begin{aligned} \lim_{\delta \downarrow 0} \frac{P(A_{t+\delta} \neq A_{t_n} | Z_{t_1}, \dots, Z_{t_{n-1}}, Z_{t_n}, Z_{t-}, Z_{t_{n+1}}, \dots, Z_{t_{n+l}})}{\delta + (t - t_n)} \\ = \frac{P(A_t \neq A_{t_n} | Z_{t_1}, \dots, Z_{t_{n-1}}, Z_{t_n}, Z_{t-}, Z_{t_{n+1}}, \dots, Z_{t_{n+l}})}{t - t_n} \end{aligned}$$

is uniform in t_n .

If we integrate out some extra variables, we can get that

$$\begin{aligned} \lim_{\delta \downarrow 0} \frac{P(A_{t+\delta} \neq A_{t_n} | Z_{t_1}, \dots, Z_{t_{n-1}}, Z_{t_n}, L_{t-}, Y_{t-}^0, Y_{t_{n+1}}^0, \dots, Y_{t_{n+l}}^0)}{\delta + (t - t_n)} \\ = \frac{P(A_t \neq A_{t_n} | Z_{t_1}, \dots, Z_{t_{n-1}}, Z_{t_n}, L_{t-}, Y_{t-}^0, Y_{t_{n+1}}^0, \dots, Y_{t_{n+l}}^0)}{t - t_n} \end{aligned}$$

is uniform in t_n .

Therefore, we can interchange the limits in the third equality.

Similarly, we can prove that

$$\lim_{\delta \downarrow 0} \frac{E[r_t(\delta) | \mathcal{H}] }{\delta} = \lim_{t_n \uparrow t-} \frac{E[I_{A_t \neq A_{t_n}} | \bar{L}_{t-}, \bar{A}_{t_n}, \bar{Y}_{t-}^0]}{t - t_n}.$$

Therefore, we have

$$\begin{aligned} \lim_{\delta \downarrow 0} \frac{E[r_t(\delta) | \mathcal{F}]}{\delta} &= \lim_{t_n \uparrow t-} \frac{E[I_{A_t \neq A_{t_n}} | \bar{L}_{t-}, \bar{A}_{t_n}, \bar{Y}_{t-}^0, \underline{Y}_{t+}]}{t - t_n} \\ &= \lim_{t_n \uparrow t-} \frac{E[I_{A_t \neq A_{t_n}} | \bar{L}_{t-}, \bar{A}_{t_n}, \bar{Y}_{t-}^0]}{t - t_n} \\ &= \lim_{\delta \downarrow 0} \frac{E[r_t(\delta) | \mathcal{H}]}{\delta}. \end{aligned}$$

The second equality comes from (20). \square

We now prove the final key lemma.

LEMMA 11. *Given FTSR, η_t is \mathcal{H}_{t-} -measurable.*

PROOF. We prove the result by using the definition of a measurable function with respect to a σ -field.

For any $a \in \mathcal{R}$, consider the following set:

$$B \equiv \{\omega : \eta_t \leq a\}.$$

Since η_t is measurable with respect to $\mathcal{F}_{t-,t+}$, $B \in \mathcal{F}_{t-,t+}$.

By Lemma 25.9 of Rogers and Williams (1994), B is a σ -cylinder and it can be decided by variables from countably many time points. Suppose the collection of these countably many time points is S . $S = S_1 \cup S_2$, where $t_{1,i} < t$ for $t_{1,i} \in S_1$ and $t_{2,j} > t$ for $t_{2,j} \in S_2$.

Let \mathcal{F}_S denote the σ -field generated by $(Z_{t_{1,i}}, i \in \mathcal{N}; Z_{t-}; Y_{t_{2,j}}^0, j \in \mathcal{N})$. We have augmented the σ -field generated by variables from S with Z_{t-} .

Next, define the following series of σ -fields:

$$\mathcal{F}_1 \equiv \sigma(Z_{t_{1,1}}, Z_{t-}, Y_{t_{2,1}}^0),$$

$$\mathcal{F}_2 \equiv \sigma(\mathcal{F}_1, Z_{t_{1,2}}, Y_{t_{2,1}}^0),$$

...

$$\mathcal{F}_\infty \equiv \mathcal{F}_S.$$

Considering the following sets:

$$B_1 \equiv \{\omega : E[\eta_t | \mathcal{F}_1] \leq a\},$$

$$B_2 \equiv \{\omega : E[\eta_t | \mathcal{F}_2] \leq a\},$$

...

$$B_S \equiv B_\infty \equiv \{\omega : E[\eta_t | \mathcal{F}_S] \leq a\}.$$

We have $B_k \in \mathcal{F}_k$.

It is easy to see that

$$B_1 \supset B_2 \supset \dots \supset B_S$$

because

$$E[E[\eta_t | \mathcal{F}_k] | \mathcal{F}_{k-1}] = E[\eta_t | \mathcal{F}_{k-1}]$$

and taking conditional expectation preserves the direction of inequality.

Also, with the above definitions, $\mathcal{F}_k \uparrow \mathcal{F}_S$. Therefore, by Theorem 5.7 from Durrett [(2005), Chapter 4], we know that

$$E[\eta_t | \mathcal{F}_k] \rightarrow E[\eta_t | \mathcal{F}_S] \quad \text{a.s.}$$

It is then easy to see that $I_{B_1} \rightarrow I_{B_S}$ a.s. and that

$$B_S = \bigcap_{i=1}^{\infty} B_i$$

with difference up to a null set.

We now claim that

$$(21) \quad B_S = B$$

with difference up to a null set.

Obviously, $B \subset B_S$. Suppose that $P(B_S - B) > 0$. Since $B_S - B \in \mathcal{F}_S$, we have

$$\int_{B_S - B} \eta_t P(d\omega) = \int_{B_S - B} E[\eta_t | \mathcal{F}_S] P(d\omega).$$

Then

$$LHS > aP(B_S - B)$$

and

$$RHS \leq aP(B_S - B).$$

This is a contradiction.

Therefore, $B = \bigcap_{i=1}^{\infty} B_i$ with difference up to a null set.

Next, we define

$$\begin{aligned} \mathcal{H}_1 &= \sigma(Z_{t_{1,1}}, Z_{t-}), \\ \mathcal{H}_2 &= \sigma(\mathcal{H}_1, Z_{t_{1,2}}), \\ &\dots \end{aligned}$$

Given FTSR, by Lemma 10, we have

$$E[\eta_t | \mathcal{F}_k] = E[\eta_t | \mathcal{H}_k].$$

Therefore, every $B_k \in \mathcal{H}_k$ and thus $B_k \in \mathcal{H}_{t-}$.

Since $B = \bigcap_{i=1}^{\infty} B_i$, $B \in \mathcal{H}_{t-}$ as well. By the definition of a measurable function, η_t is measurable with respect to \mathcal{H}_{t-} . \square

Combining all of the results in this appendix, we have proven Theorem 4.

APPENDIX D: PROOF OF THEOREM 5

Let $\mathcal{G}_t = \sigma(Y_{t-}^0, L_{t-}, A_{t-})$. Recall the definition of $r_t(\delta) = (1 - A_{t-})A_{t+\delta} + A_{t-}(1 - A_{t+\delta})$ and that $Z_t = (Y_t^0, L_t, A_t)^T$. By the Markovian property and the càdlàg property, it is easy to show that

$$E[r_t(\delta) | \sigma(\bar{Z}_{t-})] = E[r_t(\delta) | \mathcal{G}_t]$$

and that

$$E[r_t(\delta) | \sigma(\bar{Z}_{t-}, Y_{t+s}^0)] = E[r_t(\delta) | \sigma(\mathcal{G}_t, Y_{t+s}^0)].$$

Note that, without loss of generality, we only consider Y_{t+s}^0 in the proof, rather than \underline{Y}_{t+}^0 .

Therefore, we have a reduced form of *continuous-time sequential randomization*:

$$\begin{aligned} \lim_{\delta \downarrow 0} \frac{E[r_t(\delta)|\sigma(\mathcal{G}_t, Y_{t+s}^0)]}{\delta} &= \lim_{\delta \downarrow 0} \frac{E[r_t(\delta)|\sigma(\bar{Z}_{t-}, Y_{t+s}^0)]}{\delta} \\ &= \lim_{\delta \downarrow 0} \frac{E[r_t(\delta)|\sigma(\bar{Z}_{t-})]}{\delta} \\ &= \lim_{\delta \downarrow 0} \frac{E[r_t(\delta)|\mathcal{G}_t]}{\delta}. \end{aligned}$$

First, we note that if we can prove

$$(22) \quad \begin{aligned} f(Y_{t+s}^0, A_{t-}|Y_{t-}^0, L_{t-})P(A_t|Y_{t-}^0, L_{t-}) \\ = f(Y_{t+s}^0, A_t|Y_{t-}^0, L_{t-})P(A_{t-}|Y_{t-}^0, L_{t-}), \end{aligned}$$

then we can conclude (11). The reason is as follows: assuming (22) to be true, we integrate A_{t-} out on both sides of the equation. We will get

$$f(Y_{t+s}^0|Y_{t-}^0, L_{t-})P(A_t|Y_{t-}^0, L_{t-}) = f(Y_{t+s}^0, A_t|Y_{t-}^0, L_{t-}).$$

Dividing the above equation by $f(Y_{t+s}^0|Y_{t-}^0, L_{t-})$, we obtain (11).

Consider

$$g(\delta_1, \delta_2) \equiv f(Y_{t+s}^0|A_{t+\delta_1} = a_1, A_{t-\delta_2} = a_2, Y_{t-}^0, L_{t-}),$$

where $\delta_1 > 0$ and $\delta_2 > 0$.

We observe that

$$\begin{aligned} &\lim_{\delta_1 \downarrow 0} \lim_{\delta_2 \downarrow 0} g(\delta_1, \delta_2) \\ &= \lim_{\delta_1 \downarrow 0} f(Y_{t+s}^0|A_{t+\delta_1} = a_1, A_{t-} = a_2, Y_{t-}^0, L_{t-}) \\ &= \lim_{\delta_1 \downarrow 0} \frac{f(Y_{t+s}^0, A_{t+\delta_1} = a_1|A_{t-} = a_2, Y_{t-}^0, L_{t-})}{P(A_{t+\delta_1}|A_{t-} = a_2, Y_{t-}^0, L_{t-})} \\ &= f(Y_{t+s}^0|A_{t-} = a_2, Y_{t-}^0, L_{t-}) \\ &\quad \times \lim_{\delta_1 \downarrow 0} \frac{P(A_{t+\delta_1} = a_1|Y_{t+s}^0, A_{t-} = a_2, Y_{t-}^0, L_{t-})}{P(A_{t+\delta_1} = a_1|A_{t-} = a_2, Y_{t-}^0, L_{t-})} \end{aligned}$$

$$\begin{aligned}
&= \begin{cases} f(Y_{t+s}^0 | A_{t-} = a_2, Y_{t-}^0, L_{t-}) \\ \quad \times \lim_{\delta_1 \downarrow 0} \frac{1 - P(A_{t+\delta_1} \neq A_{t-} | Y_{t+s}^0, A_{t-} = a_2, Y_{t-}^0, L_{t-})}{1 - P(A_{t+\delta_1} \neq A_{t-} | A_{t-} = a_2, Y_{t-}^0, L_{t-})}, \\ \text{if } a_1 = a_2, \\ f(Y_{t+s}^0 | A_{t-} = a_2, Y_{t-}^0, L_{t-}) \\ \quad \times \lim_{\delta_1 \downarrow 0} \frac{P(A_{t+\delta_1} \neq A_{t-} | Y_{t+s}^0, A_{t-} = a_2, Y_{t-}^0, L_{t-})/\delta_1}{P(A_{t+\delta_1} \neq A_{t-} | A_{t-} = a_2, Y_{t-}^0, L_{t-})/\delta_1}, \\ \text{if } a_1 \neq a_2. \end{cases} \\
&= f(Y_{t+s}^0 | A_{t-} = a_2, Y_{t-}^0, L_{t-}).
\end{aligned}$$

Here, the validity of taking the limit inside the density is guaranteed by the third regularity condition, and the last equality follows because of the continuous-time sequential randomization assumption.

We also observe that

$$\begin{aligned}
\lim_{\delta_2 \downarrow 0} \lim_{\delta_1 \downarrow 0} g(\delta_1, \delta_2) &= \lim_{\delta_2 \downarrow 0} f(Y_{t+s}^0 | A_{t-\delta_2}, A_t, Y_{t-}^0, L_{t-}) \\
&= \lim_{\delta_2 \downarrow 0} f(Y_{t+s}^0 | A_t, Y_{t-}^0, L_{t-}) = f(Y_{t+s}^0 | A_t, Y_{t-}^0, L_{t-}).
\end{aligned}$$

The second equality uses the Markov property.

If we can interchange the limits, then we have

$$f(Y_{t+s}^0 | A_{t-}, Y_{t-}^0, L_{t-}) = f(Y_{t+s}^0 | A_t, Y_{t-}^0, L_{t-}).$$

Equation (22) follows from the definition of conditional density.

We now establish the fact that

$$\lim_{\delta_2 \downarrow 0} \lim_{\delta_1 \downarrow 0} g(\delta_1, \delta_2) = \lim_{\delta_1 \downarrow 0} \lim_{\delta_2 \downarrow 0} g(\delta_1, \delta_2)$$

by showing that $\lim_{\delta_1 \downarrow 0} g(\delta_1, \delta_2)$ is uniform in δ_2 .

If we define $g_1(\delta_2) = \lim_{\delta_1 \downarrow 0} g(\delta_1, \delta_2)$, then

$$\begin{aligned}
|g(\delta_1, \delta_2) - g_1(\delta_2)| &= \left| \frac{f(Y_{t+s}^0, A_{t+\delta_1} = a_1 | A_{t-\delta_2} = a_2, Y_{t-}^0, L_{t-})}{P(A_{t+\delta_1} = a_1 | A_{t-\delta_2} = a_2, Y_{t-}^0, L_{t-})} \right. \\
&\quad \left. - \frac{f(Y_{t+s}^0, A_t = a_1 | A_{t-\delta_2} = a_2, Y_{t-}^0, L_{t-})}{P(A_t = a_1 | A_{t-\delta_2} = a_2, Y_{t-}^0, L_{t-})} \right| \\
&= f(Y_{t+s}^0 | A_{t-\delta_2} = a_2, Y_{t-}^0, L_{t-}) \\
&\quad \times \left| \frac{P(A_{t+\delta_1} = a_1 | A_{t-\delta_2} = a_2, Y_{t-}^0, L_{t-}, Y_{t+s}^0)}{P(A_{t+\delta_1} = a_1 | A_{t-\delta_2} = a_2, Y_{t-}^0, L_{t-})} \right. \\
&\quad \left. - \frac{P(A_t = a_1 | A_{t-\delta_2} = a_2, Y_{t-}^0, L_{t-}, Y_{t+s}^0)}{P(A_t = a_1 | A_{t-\delta_2} = a_2, Y_{t-}^0, L_{t-})} \right|.
\end{aligned}$$

Consider the ratio $\frac{P(A_{t+\delta_1}=a_1|A_{t-\delta_2}=a_2, Y_{t-}^0, L_{t-}, Y_{t+s}^0)}{P(A_{t+\delta_1}=a_1|A_{t-\delta_2}=a_2, Y_{t-}^0, L_{t-})}$. We claim that it converges to $\frac{P(A_t=a_1|A_{t-\delta_2}=a_2, Y_{t-}^0, L_{t-}, Y_{t+s}^0)}{P(A_t=a_1|A_{t-\delta_2}=a_2, Y_{t-}^0, L_{t-})}$ uniformly in δ_2 .

If $a_1 = a_2$, then the density $P(A_{t+\delta_1} = a_1|A_{t-\delta_2} = a_2, Y_{t-}^0, L_{t-})$ is bounded from below by a positive number. By the fourth regularity condition,

$$\begin{aligned} P(A_{t+\delta_1} = a_1|A_{t-\delta_2} = a_2, Y_{t-}^0, L_{t-}, Y_{t+s}^0) \\ \rightarrow P(A_t = a_1|A_{t-\delta_2} = a_2, Y_{t-}^0, L_{t-}, Y_{t+s}^0) \end{aligned}$$

and

$$P(A_{t+\delta_1} = a_1|A_{t-\delta_2} = a_2, Y_{t-}^0, L_{t-}) \rightarrow P(A_t = a_1|A_{t-\delta_2} = a_2, Y_{t-}^0, L_{t-}),$$

uniformly in δ_2 , as $\delta_1 \downarrow 0$. When the denominators are bounded from below by a positive number, the ratio also converges uniformly.

If $a_1 \neq a_2$, then, by the fourth regularity condition, we have

$$\begin{aligned} \frac{P(A_{t+\delta_1} = a_1|A_{t-\delta_2} = a_2, Y_{t-}^0, L_{t-}, Y_{t+s}^0)}{\delta_1 + \delta_2} \\ \rightarrow \frac{P(A_t = a_1|A_{t-\delta_2} = a_2, Y_{t-}^0, L_{t-}, Y_{t+s}^0)}{\delta_2} \end{aligned}$$

and

$$\begin{aligned} \frac{P(A_{t+\delta_1} = a_1|A_{t-\delta_2} = a_2, Y_{t-}^0, L_{t-})}{\delta_1 + \delta_2} \\ \rightarrow \frac{P(A_t = a_1|A_{t-\delta_2} = a_2, Y_{t-}^0, L_{t-})}{\delta_2}, \end{aligned}$$

uniformly in δ_2 , as $\delta_1 \downarrow 0$. Also, the denominator $\frac{P(A_{t+\delta_1}=a_1|A_{t-\delta_2}=a_2, Y_{t-}^0, L_{t-})}{\delta_1+\delta_2}$ is bounded from below by a positive number. Hence, we establish the uniform convergence of the ratio.

Combining the two cases above, $|g(\delta_1, \delta_2) - g_1(\delta_2)|$ is bounded by $O(\delta_1)$, which does not depend on δ_2 , so $g(\delta_1, \delta_2) \rightarrow g_1(\delta_2)$ uniformly in δ_2 . Therefore,

$$\lim_{\delta_2 \downarrow 0} \lim_{\delta_1 \downarrow 0} g(\delta_1, \delta_2) = \lim_{\delta_1 \downarrow 0} \lim_{\delta_2 \downarrow 0} g(\delta_1, \delta_2).$$

By the argument at the beginning of the proof, we have proven the first part of the theorem.

To show that (11) implies FTSR, without of loss of generality, we consider

$$\begin{aligned} P(A_t|L_{t-}, Y_{t-}^0, A_{t-m}, L_{(t-m)-}, Y_{(t-m)-}^0, Y_{t+s}^0) \\ = \frac{f(A_t, L_{t-}, Y_{t-}^0, A_{t-m}, L_{(t-m)-}, Y_{(t-m)-}^0, Y_{t+s}^0)}{\sum_{i=0,1} f(A_t = i, L_{t-}, Y_{t-}^0, A_{t-m}, L_{(t-m)-}, Y_{(t-m)-}^0, Y_{t+s}^0)} \end{aligned}$$

$$\begin{aligned}
&= (f(Y_{t+s}^0 | A_t, L_{t-}, Y_{t-}^0) f(A_t, L_{t-}, Y_{t-}^0, A_{t-m}, L_{(t-m)-}, Y_{(t-m)-}^0)) \\
&\quad / \left(\sum_{i=0,1} f(Y_{t+s}^0 | A_t = i, L_{t-}, Y_{t-}^0) \right. \\
&\quad \quad \left. \times f(A_t = i, L_{t-}, Y_{t-}^0, A_{t-m}, L_{(t-m)-}, Y_{(t-m)-}^0) \right) \\
&= \frac{f(Y_{t+s}^0 | L_{t-}, Y_{t-}^0) f(A_t, L_{t-}, Y_{t-}^0, A_{t-m}, L_{(t-m)-}, Y_{(t-m)-}^0)}{\sum_{i=0,1} f(Y_{t+s}^0 | L_{t-}, Y_{t-}^0) f(A_t = i, L_{t-}, Y_{t-}^0, A_{t-m}, L_{(t-m)-}, Y_{(t-m)-}^0)} \\
&= \frac{f(A_t, L_{t-}, Y_{t-}^0, A_{t-m}, L_{(t-m)-}, Y_{(t-m)-}^0)}{\sum_{i=0,1} f(A_t = i, L_{t-}, Y_{t-}^0, A_{t-m}, L_{(t-m)-}, Y_{(t-m)-}^0)} \\
&= P(A_t | L_{t-}, Y_{t-}^0, A_{t-m}, L_{(t-m)-}, Y_{(t-m)-}^0).
\end{aligned}$$

The second equality follows because of the Markov property. The third equality uses equation (11). We have thus proven the second half of the theorem.

APPENDIX E: SIMULATION PARAMETERS

In all simulation models from M1 to M4, we specify the parameters as follows:

- let $g(V, t) = C$, a constant; let $C = 100$;
- for M1 (also for M3 and M4), let $\theta = 0.2$ and $\sigma = 1$;
- for M2, let $m = 2$, $\theta_1 = 0.2$, $\sigma_1 = 1$ and $\theta_2 = 1$, $\sigma_2 = 0.5$. The transition probability of J_t would be $P(t) = e^{At}$, where $A = \begin{pmatrix} -1 & 1 \\ 1 & -1 \end{pmatrix}$;
- for initial value, e_0 is generated from $N(0, \frac{\sigma}{\sqrt{2\theta}})$;
- the causal parameter $\Psi = 1$;
- in M1, M2 and M3, $s(A_t, Y_t) = e^{\alpha_0 + \alpha_1 A_t + \alpha_2 Y_t + \alpha_3 A_t Y_t}$; let $\alpha_1 = -0.3$, $\alpha_2 = -0.005$, $\alpha_3 = 0.007$ and $\alpha_0 = -0.2$;
- in M4, A_t is generated as follows: if $Y_{t-0.5} > 101$ and $Y_t > 101$, $s(A_t = 1, L_t^*) = 2.8$; if $Y_{t-0.5} < 99$ and $Y_t < 99$, $s(A_t = 0, L_t^*) = 2.8$; otherwise, A_t is generated following a model similar to that in M1, except that $s(A_t, L_t^*) = e^{\alpha_0 + \alpha_1 A_t + \alpha_2 L_t^* + \alpha_3 A_t L_t^*}$; the values of the α 's are the same as before;
- in M4, η_t follows an Ornstein–Uhlenbeck process with parameters $\theta = 0.2$ and $\sigma = 1$;
- for initial value, A_0 is generated from Bernoulli($\text{expit}(\alpha_0 + \alpha_2 Y_0)$);
- $K = 5$ is the number of periods;
- number of subjects $n = 5000$.

Acknowledgments. The authors would like to thank the causal inference reading group at the University of Pennsylvania and Judith Lok for helpful discussions.

REFERENCES

- ANDERSEN, P. K., BORGAN, Ø., GILL, R. and KEIDING, N. (1992). *Statistical Models Based on Counting Processes*. Springer, New York.
- BRUMBACK, B. A., HERNÁN, M. A., HANEUSE, S. J. and ROBINS, J. M. (2004). Sensitivity analyses for unmeasured confounding assuming a marginal structural model for repeated measures. *Stat. Med.* **23** 749–767.
- COCHRAN, W. G. (1965). The planning of observational studies of human populations (with discussion). *J. Roy. Statist. Soc. Ser. A* **128** 134–155.
- DEL NINNO, C., DOROSH, P., SMITH, L. and ROY, D. (2001). The 1998 floods in Bangladesh: Disaster impacts, household coping strategies and response. Research Report 122, International Food Policy Research Institute, Washington, DC.
- DEL NINNO, C. and LUNDBERG, M. (2005). Treading water: The long-term impact of the 1998 flood on nutrition in Bangladesh. *Economics and Human Biology* **3** 67–96.
- DURRETT, R. (2005). *Probability: Theory and Examples*, 3rd ed. Brooks/Cole–Thomson Learning, Belmont, CA.
- GASSER, T., MULLER, H., KOHLER, W., MOLINARI, L. and PRADER, A. (1984). Nonparametric regression analysis of growth curves. *Ann. Statist.* **12** 210–229. [MR0733509](#)
- GUERRANT, R. L., KOSEK, M., LIMA, A. A. M., LORNTZ, B. and GUYATT, H. L. (2002). Updating the Dalys for diarrhoeal disease. *Trends in Parasitology* **18** 191–193.
- HERNÁN, M. A., BRUMBACK, B. and ROBINS, J. M. (2002). Estimating the causal effect of zidovudine on CD4 count with a marginal structural model for repeated measures. *Stat. Med.* **21** 1689–1709.
- JOFFE, M. M. and ROBINS, J. M. (2009). Controlling the future: Revised assumptions and methods for causal inference with repeated measures outcomes. Working paper.
- KASLOW, R. A., OSTROW, D. G., DETELS, R., PHAIR, J. P., POLK, B. F. and RINALDO, C. R., Jr. (1987). The Multicenter AIDS Cohort Study: Rationale, organization, and selected characteristics of the participants. *American Journal of Epidemiology* **126** 310–318.
- KOSEK, M., BERN, C. and GUERRANT, R. L. (2003). The global burden of diarrhoeal disease as estimated from studies published between 1992 and 2000. *Bulletin of the World Health Organization* **81** 197–204.
- LOK, J. J. (2008). Statistical modeling of causal effects in continuous time. *Ann. Statist.* **36** 1464–1507. [MR2418664](#)
- MARTORELL, R. and HO, T. J. (1984). Malnutrition, morbidity and mortality. *Population and Development Review* **10** Supplement: Child Survival: Strategies for Research.
- MOORE, S. R., LIMA, A. A. M., CONAWAY, M. R., SCHORLING, J. B., SOARES, A. M. and GUERRANT, R. L. (2001). Early childhood diarrhea and helminthiases associate with long-term linear growth faltering. *International Journal of Epidemiology* **30** 1457–1464.
- ROBINS, J. M. (1986). A new approach to causal inference in mortality studies with sustained exposure periods—Application to control of the healthy worker survivor effect. *Math. Model.* **7** 1393–1512 [errata (1987) **14** 917–921].
- ROBINS, J. M. (1992). Estimation of the time-dependent accelerated failure time model in the presence of confounding factors. *Biometrika* **79** 321–334. [MR1185134](#)
- ROBINS, J. M. (1994). Correcting for non-compliance in randomized trials using structural nested mean models. *Comm. Statist.* **23** 2379–2412. [MR1293185](#)
- ROBINS, J. M. (1997). Causal inference from complex longitudinal data. In *Latent Variable Modelling and Applications to Causality* (M. Berkane, ed.). *Lecture Notes in Statist.* **120** 69–117. Springer, New York. [MR1601279](#)
- ROBINS, J. M. (1998). Marginal structural models. In *1997 Proceedings of the Section on Bayesian Statistical Science* 1–10. Amer. Statist. Assoc., Alexandria, VA.

- ROBINS, J. M. (2000). Marginal structural models versus structural nested models as tools for causal inference. In *Statistical Models in Epidemiology, the Environment and Clinical Trials* (M. E. Halloran and D. Berry, eds.). *The IMA Volumes in Mathematics and Its Applications* **116** 95–133. Springer, New York. [MR1731682](#)
- ROBINS, J. M. (2008). Causal models for estimating the effects of weight gain on mortality. *International Journal of Obesity* **32** S15–S41.
- ROGERS, L. C. G. and WILLIAMS, D. (1994). *Diffusions, Markov Processes, and Martingales* **1**. Wiley, Chichester. [MR1331599](#)
- ROSENBAUM, P. R. (1984). The consequences of adjustment for a concomitant variable that has been affected by the treatment. *J. Roy. Statist. Soc. Ser. A* **147** 656–666.
- SINGER, B. (1981). Estimation of nonstationary Markov chains from panel data. *Sociological Methodology* **12** 319–337.
- TSIATIS, A. (2006). *Semiparametric Theory and Missing Data*. Springer, New York. [MR2233926](#)
- VAN DER VAART, A. W. (2000). *Asymptotic Statistics*. Cambridge Univ. Press.

M. ZHANG
D. S. SMALL
DEPARTMENT OF STATISTICS
THE WHARTON SCHOOL
UNIVERSITY OF PENNSYLVANIA
PHILADELPHIA, PENNSYLVANIA 19104
USA
E-MAIL: zhangmi@wharton.upenn.edu
dsmall@wharton.upenn.edu

M. M. JOFFE
DEPARTMENT OF BIostatISTICS
& EPIDEMIOLOGY
UNIVERSITY OF PENNSYLVANIA
PHILADELPHIA, PENNSYLVANIA 19104
USA
E-MAIL: mjoffe@mail.med.upenn.edu