

Simulating from Marginal Structural Models with Time-Dependent Confounding

W. G. Havercroft* and V. Didelez*

We discuss why it is not always obvious how to simulate longitudinal data from a general marginal structural model (MSM) for a survival outcome while ensuring that the data exhibit complications due to time-dependent confounding. Based on the relation between a directed acyclic graph (DAG) and an MSM, we suggest a data-generating process that satisfies both these requirements, the general validity of which we prove. Our approach is instructive regarding the interpretation of MSMs, and useful in that it allows one to examine the finite sample performance of methods which claim to adjust for time-dependent confounding. We apply our methodology to design a simulation study which emulates typical longitudinal studies such as the Swiss HIV Cohort Study, so that competing methods of adjusting for time-dependent covariates can be compared. Copyright © 2011 John Wiley & Sons, Ltd.

Keywords: Causal inference; Time-dependent confounding; Marginal structural models; Survival analysis; Longitudinal data; Simulation; Inverse probability weights

1. Introduction

Marginal structural models (MSMs), together with the method of inverse probability of treatment weighting (IPTW) used to fit these models, have become a popular tool for investigating causal effects of time-varying treatments (or exposures) in longitudinal studies [1, 2, 3, 4, 5]. A typical application is the Swiss HIV Cohort Study investigating the effect of highly active antiretroviral therapy (HAART) on survival [3]; other potential examples could deal with the adjustment of anticoagulant dosage for stroke patients [6, 7]. It is well known that when investigating the effect of time-varying treatments from longitudinal studies time-dependent confounding has to be taken into account [8]. For a recent overview see Daniel et al. [9].

Here we are concerned with the question of how to simulate data from a given MSM, such that the data-generating process exhibits time-dependent confounding with the added complication of current treatment affecting future covariates. In short, the problem occurs when the *conditional* distributions one might use to draw the simulated data from are not compatible with the desired properties of the *marginal* model, e.g. because of non-collapsibility. However, it is clear that being able to simulate data from a known true model is crucial to investigate for instance finite sample properties or efficiency of different methods. We propose a general approach that allows such simulation from a given MSM, as well

as a specific algorithm for a discrete-time hazard model, which emulates longitudinal HIV studies where antiretroviral treatment may be started at different times, depending on changes in CD4 count.

Our proposal extends and gives a formal justification to work by Bryan et al. [10] and Young et al. [11, 12].

The paper is structured as follows. In Section 2 we introduce the basic concepts, including a definition of causal effect, representation of relationships between random variables using directed acyclic graphs (DAGs), structural assumptions, time-dependent confounding, MSMs and IPTW. In Section 3, we present and apply the sampling algorithm. Section 3.1 details the basic mechanism by which the presented algorithm generates data from the desired MSM with time-dependent confounding, a proof is given in Appendix B; then Section 3.2 gives the details for a specific case. In Section 3.3 we compare our algorithm to others suggested in the literature. Section 3.4 illustrates the use of the algorithm by carrying out a small simulation study comparing the performance of different approaches dealing with time-dependent confounding. Finally we discuss possible extensions and future work.

2. Basic Concepts

As mentioned in the introduction, we focus on the application of MSMs to typical longitudinal studies, for example the Swiss HIV Cohort Study [3]. Let $\mathcal{T} = \{0, \dots, T\}$ denote the finite sequence of discrete observation time points. We denote the treatment (action) process by $\{A_t\}_{t \in \mathcal{T}}$ — this might simply be a sequence of binary variables to indicate whether the patient is on treatment or not; it could also be a general discrete or continuous variable, e.g. the changing dosage of an anticoagulant for stroke patients. Further $\{L_t\}_{t \in \mathcal{T}}$ will typically denote one (or a set of) repeatedly observed covariate(s) such as CD4 count for HIV patients. There may be latent variables which we denote by $\{U_t\}_{t \in \mathcal{T}}$. Finally we consider either a single outcome variable Y measured after time T , or a process $\{Y_t\}_{t \in \mathcal{T}^+}$, where $\mathcal{T}^+ = \{1, \dots, T + 1\}$. Typically Y_t is an indicator for the occurrence of an event, e.g. death, in which case there will be an equivalence between $\{Y_t\}_{t \in \mathcal{T}^+}$ and $\tilde{Y} = \min\{t : Y_t = 1\}$.

Throughout, we denote the history of a stochastic process $\{X_t\}_{t \in \mathcal{T}}$ up to time t as \bar{X}_t , while $\bar{X} = \bar{X}_T$ denotes the complete history.

2.1. Causal Effects and Marginal Structural Models

First we consider the single time point case, to keep the exposition simple. We define the causal effect of one variable A on another Y in terms of interventions. Following the approach of Pearl [13], we denote the distribution of Y given an intervention by which A is fixed at the value a as $P(Y | do(A = a))$ (abbreviated to $P(Y | do(a))$ when the context is clear). Note the distinction between this, the distribution of the random variable ‘ Y given an intervention a ’, and $P(Y | A = a)$, the distribution of the random variable ‘ Y given A is observed to be a ’. It is also common to define causal effects in terms of potential outcomes, where $Y(a)$ denotes the outcome that would be observed if A were (possibly counter to the fact) set to a [14, 15, 8]. For the present purpose we can regard $P(Y(a))$ and $P(Y | do(a))$ as equivalent, but see discussion by Dawid [16, 17].

Broadly, we say that there is a causal effect if there is some contrast between the distribution $P(Y | do(a))$ for different values of a . Contrasts under different interventions for other quantities, such as expectations, odds ratios or hazard functions, can be used to define causal effects, depending on the context of the problem. More specifically, a model that parameterises $P(Y | do(a))$ is a *marginal structural model* (MSM). The model is *marginal* over any covariates, while the term *structural* distinguishes it as an interventional model, versus some observational marginal model $P(Y | A = a)$.

The case of a sequence of treatment decisions or actions $\bar{A} = \{A_0, \dots, A_T\}$ is analogous. For instance, the causal effect of any sequence \bar{a} of actions, which we also call a *strategy*, can be modeled by an MSM which parameterises

$$P(Y | do(\bar{a})) = f(\bar{a}, \gamma),$$

where γ is a vector of parameters. The exact form of the model depends on the structure of the data, such that it may be defined for example in terms of $E(Y | do(\bar{a}))$, or for instance the corresponding hazard function when Y is a survival outcome.

Note that an MSM as used here models the effect of a given strategy \bar{a} . For example in the case of binary treatment decisions, an MSM may model the effect of strategies such as $\bar{a} = (0, \dots, 0)$, ‘never treat’ and $\bar{a} = (1, \dots, 1)$, ‘always treat’. As this means that all treatment decisions are fixed in advance one could say that an MSM specifies the effects of earlier treatments not mediated by later treatments, e.g. if we were to compare $\bar{a} = (0, \dots, 0)$ with $\bar{a} = (1, 0, \dots, 0)$. For the moment, we do not consider dynamic strategies which make a_t a function of time-varying covariates (see discussion in Section 4).

2.2. Directed Acyclic Graphs (DAGs)

We will use directed acyclic graphs (DAGs) to illustrate key aspects of causal structures. Such a DAG is given by a set of nodes V representing variables, and a set of directed edges. We call V_i a ‘parent’ of V_j if $V_i \rightarrow V_j$, and an ‘ancestor’ of V_k if $V_i \rightarrow \dots \rightarrow V_k$. A DAG represents assumptions of conditional independencies by the absence of edges between nodes [18]. For example Figure 1 includes the independence $Y \perp\!\!\!\perp \{L, W\} | \{A, U\}$ (where $B \perp\!\!\!\perp C | D$ denotes conditional independence between B and C given D [19]).

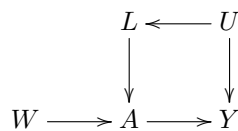


Figure 1. Example of a directed acyclic graph (DAG).

A DAG imposes a factorisation on the joint distribution of the set of random variables V , which is given by the product of the conditional distributions of each variable in V given all its parents. The DAG can be used to further encode *causal* assumptions in the following way. We define that the DAG is *causal* with respect to the variable $S \in V$ if we can modify the non-interventional joint distribution for the effect of an intervention which fixes S at the value s , simply by [13, p. 24] [20]:

1. Replacing the conditional distribution of S given its parents with the identity function $\mathbb{I}_{S=s}$, and
2. Substituting the intervention $S = s$ into the remaining conditional distributions.

It is common (but often not necessary nor appropriate) to assume the whole DAG is *causal*, meaning the above holds for any subset of variables $S \subset V$.

For the DAG of Figure 1, we have the factorisation

$$P(U, L, W, A, Y) = P(U)P(W)P(L|U)P(A|L, W)P(Y|U, A). \tag{1}$$

Demanding that the DAG be *causal* with respect to the variable A means that the joint intervention distribution of the remaining variables is given by

$$P(U, L, W, Y | do(a)) = P(U)P(W)P(L|U)P(Y|U, A = a)\mathbb{I}_{A=a}. \tag{2}$$

In words, the DAG will be *causal* with respect to A if it describes the system in sufficient detail so that we can believe that an intervention in A does not alter the remaining conditional distributions, other than through the value of $A = a$. This

has been called the *stability* property by Dawid and Didelez [21], and can also be made graphically explicit with influence diagrams [22, 23].

The general version of (2) is known as a *truncated factorisation formula* as it is obtained by ‘dropping’ certain factors from (1) [13, p. 72] (see also the manipulation theorem of Spirtes et al. [24, p. 51]). This is in turn equivalent to dividing the joint distribution by these factors; for example (2) is obtained from (1) upon dividing by $P(A | L, W)$. This is an alternative way of motivating the IPTW method which we describe in more detail in Section 2.6.

2.3. Relation between MSMs and DAGs

The factorisation of the joint distribution immediately suggests a generating process from which it would be straightforward to simulate data that obey a given DAG. For Figure 1, choose specific conditional distributions and sequentially generate first W and U , then L given U , followed by A given $\{L, W\}$ and finally Y given $\{A, U\}$. However, assume that we also wanted this data-generating process to obey a given MSM parameterising $P(Y | do(a))$. By integrating out the remaining variables in (2) we see that

$$\begin{aligned}
 P(Y | do(a)) &= \sum_l P(Y | L = l, A = a)P(L = l) & (3) \\
 &= \sum_u P(Y | U = u, A = a)P(U = u). & (4)
 \end{aligned}$$

(Replace sums by integrals if variables are continuous.) This means that in order to simulate data from a given MSM within the structure of Figure 1, we need either suitable choices of $P(Y | L, A)$ and $P(L)$, or of $P(Y | U, A)$ and $P(U)$. In some cases this is easy to achieve; e.g. if the MSM assumes a linear mean $E(Y | do(a)) = \beta_0 + \beta_1 a$ and we also have $E(Y | L = l, A = a) = \alpha_0 + \alpha_1 a + \alpha_2 l$, then $\sum_l E(Y | L = l, A = a)P(L = l) = \alpha_0 + \alpha_1 a + \alpha_2 E(L)$ so that $\beta_1 = \alpha_1$. As soon as general non-linear models (possibly with interactions) are involved, especially non-collapsible ones such as logistic models, it is no longer obvious how the conditional distribution on the RHS of (3) or (4) can result in a given MSM which specifies the LHS of (3).

For the case of a single time point, our proposal in Section 3.1 works as follows: take U to be an unobservable latent variable, e.g. representing general health, with a $U[0, 1]$ marginal distribution. Hence any distribution $P(Y | do(a))$ can be obtained from (4) by transforming U with the CDF F^a of $P(Y | do(a))$. At the same time, the variable L in Figure 1 can be taken to play the role of an observable covariate used to adjust for confounding as we explain in the next section.

2.4. No Unmeasured Confounding

We say that there exists some confounding of the effect of A on Y if $P(Y | do(a)) \neq P(Y | A = a)$, i.e. if the distribution of Y given an intervention $do(A = a)$ is not equal to the distribution of Y given $A = a$ is observed. In the example of Figure 1, the effect of A on Y is confounded, which we can see formally by noting that neither (3) nor (4) are equal to $P(Y | A = a)$, since the distributions of L and U respectively are not conditional on A in the sum. Another way of putting this is by noting that in the DAG of Figure 1 there is a back-door path from A to Y which creates a non-causal association between them [25].

It may be possible to obtain $P(Y | do(a))$ by suitably making use of additional data on covariates which is often called *adjusting for confounding*; if this is possible then $P(Y | do(a))$ (or a specific parameter of this distribution) is *identifiable*. There are different approaches to adjust for confounding, but all of them require that *sufficient* additional information is available. In the case of a single treatment variable, a set of covariates is sufficient to adjust for confounding if it blocks all back-door paths from A to Y in a DAG that is causal with respect to A [25, 23]. For Figure 1 we find that either L or U is sufficient to adjust for confounding. An analogous criterion can be given for the case of sequential (non-dynamic) treatments: For every A_t , every back-door path to Y must be blocked by $\{\bar{L}_t, \bar{A}_{t-1}\}$ in the modified DAG where all

incoming arrows into any future A_{t+k} are deleted [26, 27, 21]. This property is also often called *sequential randomisation* [8], or *no unmeasured confounding* [2, 12].

Given covariates that are sufficient to adjust for confounding, the most popular method to do so is *regression adjustment*. In the example of Figure 1, with L as the covariate, this amounts to estimating $P(Y | L, A)$ in (3). This is clearly informative about the causal effect, as we have $P(Y | L, A = a) = P(Y | L; do(A = a))$, though it does not explicitly target the marginal effect $P(Y | do(a))$ for which we require the additional standardisation with $P(L)$. In the case of sequential treatments, regression adjustment is more problematic as $P(Y | \bar{L}, \bar{A} = \bar{a})$ is *not* generally equal to $P(Y | \bar{L}; do(\bar{A} = \bar{a}))$ when some of the time-varying covariates L_t are measured after and affected by some of the treatments \bar{A}_{t-1} . This is discussed in the next section.

2.5. Time-Dependent Confounding

Now, we turn to the case of more than one treatment variable and the problem of time-dependent confounding. We consider a very simple case of longitudinal data. For each subject, a treatment decision A_0 is made at time 0, e.g. whether to initiate HAART for HIV patients. The value of a diagnostic covariate L_1 , e.g. CD4 count, measured at time 1, is potentially affected by the earlier A_0 and may itself inform a second treatment decision A_1 , e.g. to start HAART if CD4 is low. The diagnostic covariate L_1 will typically also depend on an underlying latent general health process which we represent by U , and which in turn is predictive of the outcome of interest Y , e.g. survival.

We may represent the above structure in the DAG in Figure 2, which we assume is causal with respect to \bar{A} . We see that there is confounding of the relationship between A_1 and Y , since L_1 is a common parent (and U a common ancestor) of the two; hence $P(Y | do(\bar{a})) \neq P(Y | \bar{A} = \bar{a})$ in general.

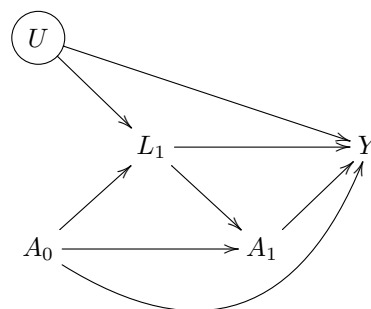


Figure 2. DAG depicting simple time-dependent confounding. Circled node denotes latent variable.

Recall that in a DAG, the absence of edges encodes crucial independencies; here we assume that the treatment decisions A_t are conditionally independent of U given $\{\bar{L}_t, \bar{A}_{t-1}\}$ as shown by the absence of an edge from U into A_t . Such an assumption could be justified if it is plausible that any information which clinicians in the study use to make their treatment decisions is contained in \bar{L}_t . This structure ensures that in Figure 2, L_1 is sufficient to adjust for confounding, i.e. as long as we have data on L_1 , there is no unobserved confounding.

In analogy to the time-fixed case one may consider regression adjustment, i.e. estimation of the observational conditional distribution $P(Y | L_1, A_0, A_1)$. However in this time-dependent case, adjusting for L_1 in this way will introduce other sources of bias. Firstly, the time-dependent data structure permits the existence of causal pathways from treatment to outcome which are *mediated* via the covariates. In this case, we have the mediated causal pathway $A_0 \rightarrow L_1 \rightarrow Y$. Regression adjustment for L_1 could block such pathways, resulting in distorted estimates of the overall effects of earlier treatment decisions on the final outcome. Secondly, conditioning on L_1 , which is graphically known as a *collider* on the path $A_0 \rightarrow L_1 \leftarrow U \rightarrow Y$, induces an association between its co-parents A_0 and U , and hence between A_0 and Y . This is known as *selection bias* [28, 29, 30]. Selection bias often gets overlooked, but its effect can be clearly demonstrated if we consider the special case where neither A_0 nor A_1 have any causal effect on Y at all, as in Figure 3. There are no

causal pathways from \bar{A} to Y , yet conditioning on L_1 results in a conditional association between A_0 and Y , which will show up in a regression model for $P(Y | L_1, A_0, A_1)$ as a non-zero coefficient of A_0 .

In the general case with more time points, selection bias, as well as bias due to blocked mediated effects, will almost automatically become a problem whenever *summary* measures of treatment and/or covariates, such as *cumulative* or average treatment or CD4 count, are used [31]. It can then also happen that a mediation pathway like $A_t \rightarrow L_{t+1} \rightarrow A_{t+1} \rightarrow Y$ becomes problematic.

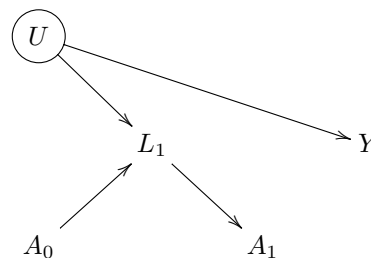


Figure 3. DAG representing the case of no causal effect of \bar{A} , but selection bias.

Thus we have three potential sources of bias in a longitudinal setting: time-dependent covariates may confound the effect of future treatments on the outcome; but if we condition on these covariates we could potentially induce selection bias as well as obscuring mediation effects with respect to earlier treatments. These issues characterise *time-dependent confounding* [2], and are problematic in that there is potential for bias whether or not we use $\{L_t\}$ in regression adjustment. The following section addresses one possible method which takes time-dependent confounding into account without introducing bias.

2.6. Inverse Probability of Treatment Weighting

Inverse probability of treatment weighting (IPTW) [1, 2] is a method for correcting for time-dependent confounding when fitting an MSM. This is achieved by weighting each subject's data based on their observed treatment history. The basic inverse probability of treatment weight for subject i at time t takes the form

$$W_{t,i} = \frac{1}{\prod_{\tau=0}^t p_{\tau}(A_{\tau,i} | \bar{A}_{\tau-1,i}, \bar{L}_{\tau,i})}$$

where $p_t(x|\bar{z}, \bar{w}) = P(A_t = x | \bar{A}_{t-1} = \bar{z}, \bar{L}_t = \bar{w})$, so that the denominator is the probability that the subject received the particular treatment history they were observed to receive up to time t , given prior observed treatment and covariate histories. Applying the terminal weights $W_{T,i}$ to each subject in the sample results in a pseudo-population in which treatment A_t is no longer affected by past covariates \bar{L}_t , breaking the confounding; but crucially, the causal effect of \bar{A} on Y remains unchanged. Then the parameters of the MSM coincide with those of the re-weighted observational marginal model, which may be estimated using standard methods on the re-weighted data. The resulting estimates are consistent under the following assumptions: the MSM is correctly specified, the denominators are non-zero as well as being correctly specified and fitted using a consistent method, and given $\{L_t\}$ there is no unobserved confounding.

In practice, we use the stabilised form of the IPTW:

$$\widetilde{W}_{t,i} = \frac{\prod_{\tau=0}^t p_{\tau}(A_{\tau,i} | \bar{A}_{\tau-1,i})}{\prod_{\tau=0}^t p_{\tau}(A_{\tau,i} | \bar{A}_{\tau-1,i}, \bar{L}_{\tau,i})} \tag{5}$$

where $p_t(x|\bar{z}) = P(A_t = x | \bar{A}_{t-1} = \bar{z})$. It can be shown that replacing the numerator of the weight by any function of \bar{A}_t affects the efficiency of the IPTW estimator of γ , but not its consistency. Furthermore, $\widetilde{W}_{t,i}$ generally gives a more efficient estimator than $W_{t,i}$ [1, 2]. It also has the property that the weights reduce to 1 when there is no confounding, so that the weighted analysis is asymptotically equivalent to the unweighted analysis (see an empirical example of this in Section 3.4.3).

2.7. Extension to Time-to-Event Outcomes

Sections 2.5–2.6 focus on a simple form of longitudinal data in which follow-up proceeds for a fixed amount of time, and the outcome Y is measured at the end. This extends naturally to the case of time-to-event data (in discrete as well as continuous time) with finite follow-up, in which Y is treated as a binary indicator process $\{Y_t\}_{t \in \mathcal{T}^+}$ [11, 32, 33, 34, 35]. Here, $Y_t = 1$ when the subject has failed (e.g. died) before time t , $Y_t = 0$ otherwise, and the outcome of interest is the survival time, \widetilde{Y} . The key differences in methodology for this case are as follows: Instead of applying terminal IPTWs $W_{T,i}$ to each subject, the time-varying weights $W_{t,i}$ are applied at each time point t for each subject; and often the hazard function, or the probability of failure in a particular interval, is of primary interest, so that it is necessary to condition on $\bar{Y}_t = 0$, i.e. ‘survival’ up to the present.

3. Simulating Data from a Given MSM with Time-Dependent Confounding

We now address the issue of how to simulate longitudinal data with a survival outcome from a given MSM, so that the underlying data-generating process exhibits time-dependent confounding. For the purpose of evaluating methods which correct for time-dependent confounding, we require that the true marginal structural model for \bar{Y} given $do(\bar{a})$ takes a known form.

If, instead of a survival analysis, we were examining a data structure with a fixed-time outcome, with linear relationships between variables (e.g. multivariate Gaussian), the requirement would be straightforward to achieve. We could simply use the structure of Figure 2, in which \bar{Y} is affected by \bar{A} and \bar{L} . This would give a known *conditional* model for \bar{Y} given \bar{A} and \bar{L} , which, due to linearity, we could *collapse* over \bar{L} to derive the marginal model for \bar{Y} given \bar{A} alone. However this is not as simple in the case of a survival outcome, where hazard functions, or their discrete equivalents based on generalised linear models, are typically non-collapsible.

3.1. General Idea

The general idea is essentially based on a longitudinal version of Figure 1 and the issues discussed around (4). We take a uniformly distributed U_0 , which represents a latent general health variable; we generate $\{L_t, A_t\}$ depending on U_0 , as well as transforming U_0 to obtain Y_t with the desired MSM’s survival function $S^{\bar{a}}(t)$ once the actual treatments $\bar{A}_t = \bar{a}_t$ are known. Note that the uniform distribution of U_0 is not crucial, as long as Y_t given U_0 can be generated from the desired $S^{\bar{a}}(t)$. We can ensure that $\{L_t\}$ is sufficient to adjust for confounding by choosing conditional distributions for A_t that are independent of U_0 given $(\bar{L}_t, \bar{A}_{t-1})$. We can further easily choose conditional distributions for L_t that depend on \bar{A}_{t-1} and U_0 so that past treatment decisions affect future covariates. (Note that in the specific algorithm given below, we induce noise in the dependence of L_t on U_0 by including further latent variables U_1, \dots, U_T .)

The procedure suggested above can simply be regarded as one way of matching the LHS of (3) as given by the MSM with a suitable choice of the factors on the RHS of (4). However, it could also be interpreted in terms of potential outcomes as follows. The MSM determines a survival function (or CDF) $S^{\bar{a}}(t)$ for *any* treatment strategy $do(\bar{a})$. So, in principle we could generate the set of all counterfactual outcomes of an individual i , $\{\bar{Y}_i(\bar{a}) : \bar{a} \in \mathcal{A}\}$, at the start, and then generate the remaining processes in a way which ensures $\{L_t\}$ is sufficient to adjust for confounding. Once the observed treatment process $\bar{A} = \bar{a}^*$ is known, the corresponding potential outcome process $\bar{Y}_i(\bar{a}^*)$ is picked and becomes the observed

outcome \bar{Y}_i . The procedure suggested above instead provides a shortcut in that it waits until the treatment process is known and then generates the potential outcome, which is far more computationally efficient.

An obvious disadvantage of the suggested approach is the way in which the dependence of $\{Y_t\}$ on $\{L_t\}$ can only implicitly be specified via their respective dependencies on U_0 . This also means that $\{L_t\}$ is only a very weak mediator of past treatment decisions in the sense that it may influence survival via future treatment decisions but not via other pathways. If we wanted to model a more explicit dependence we would either need to use collapsible models (or maybe exploit conditions for near-collapsibility [32]), or make the RHS of (3) match with the LHS as given by the MSM. This approach will generally be more challenging and intricate, especially in the general longitudinal case.

3.2. Data-Generating Algorithm

We now present a specific algorithm based on the above proposal. We wish to emulate a simplified version of the data structure of the Swiss HIV Cohort Study [3], for the effect of HAART versus no treatment on the ‘survival’ time of HIV patients (time to AIDS-related illness or death; we assume death from here onwards for simplicity), with time-dependent confounding by CD4 count and no other measured covariates.

3.2.1. Data Structure. Our data structure operates in discrete time $t \in \mathcal{T}$. The four processes included in the model are as follows.

A binary treatment indicator process $\{A_t\}_{t \in \mathcal{T}}$, where $A_t = 1$ when the subject is on treatment (e.g. HAART) at time t , $A_t = 0$ otherwise. Once a subject has started treatment, they remain on treatment until failure or end of follow-up.

A covariate process $\{L_t\}_{t \in \mathcal{T}}$ represents one or more diagnostic variables. In the context of the Swiss HIV Cohort Study, L_t would be the subject’s CD4 count in cells/ μL at time t . A lower value of L_t indicates more severe illness.

The binary failure process $\{Y_t\}_{t \in \mathcal{T}^+}$ indicates when death has occurred before time t , $Y_t = 1$, and otherwise $Y_t = 0$.

As before $\{U_t\}_{t \in \mathcal{T}}$ represents a latent ‘general health’ process, where $U_t \in [0, 1]$. A value of U_t close to 0 indicates poor health, U_t close to 1 indicates good health.

‘Check-up’ times, at which L_t is measured and A_t is chosen, only occur at every k th time point (the first check-up is at $t = 0$), while U_t is updated and Y_t is measured at each t . This permits a greater granularity in the failure times; for a fixed ratio t/k , as $t \rightarrow \infty$, the failure time approaches a continuous process. Between check-ups L_t and A_t are taken to be constant.

3.2.2. MSM. As we assume discrete-time data, we will not take a continuous proportional hazards approach to the survival analysis. Instead we work with the probability of failure in a single interval $(t, t + 1]$ given survival up to t , as the discrete equivalent of the hazard function. Let $\lambda_t^{\bar{a}}$ denote $\lambda_t(\text{do}(\bar{a})) = P(Y_{t+1} = 1 \mid Y_t = 0; \text{do}(\bar{a}))$. Note that, in the present case, any intervention strategy \bar{a} is fully specified by the time t^* of treatment initiation so that we can reformulate our model in terms of t^* . The algorithm given below generates data from the following MSM:

$$\begin{aligned} \lambda_t^{\bar{a}} &= \text{logit}^{-1}(\gamma_0 + \gamma_1[(1 - a_t)t + a_t t^*] + \gamma_2 a_t + \gamma_3 a_t(t - t^*)) \\ &= \text{logit}^{-1}(\gamma_0 + \gamma_1 d_{1,t} + \gamma_2 a_t + \gamma_3 d_{3,t}), \end{aligned} \tag{6}$$

where $d_{1,t} = \min\{t, t^*\}$ and $d_{3,t} = \max\{t - t^*, 0\}$ represent the time elapsed before and after treatment initiation respectively. Thus the structural model for the effect of treatment on survival can represent situations in which, for example, initiating treatment improves survival chance ($\gamma_2 < 0$), but also increases the rate at which survival chance decreases with time ($\gamma_3 > \gamma_1 > 0$). In this case, the problem of determining the optimal time to initiate treatment is non-trivial. Note that often the hazard function is chose to depend only on the current treatment, e.g. $\text{logit}^{-1}(\alpha_0 + \alpha_1 a_t)$; the above choice

means that the hazard depends on a summary of the treatment history \bar{a}_t by including the duration without treatment as well as the duration on treatment.

3.2.3. The Algorithm. The following algorithm generates data under the structure of Section 3.2.1 which follow the MSM of Section 3.2.2, exhibiting time-dependent confounding, and ensuring that data on \bar{L}_t are sufficient to adjust for confounding. For each subject:

- Draw $U_0 \sim U[0, 1]$
- Compute the baseline CD4 L_0 as a transformation of U_0 by the inverse CDF of the $\Gamma(k = 3, \theta = 154)$ distribution, plus a $N(0, 20)$ error term. This inverse CDF transform was chosen to heuristically approximate the sample distribution of baseline CD4 counts in the Swiss HIV Cohort Study [3].
- Draw treatment decision $A_0 \sim \text{Bernoulli}(p_{A_0})$, where $p_{A_0} = \text{logit}^{-1}(\theta_0 + \theta_2(L_0 - 500))$, where clinical consensus holds that a CD4 count of 500 is a conservative lower ‘healthy’ bound.
- If $A_0 = 1$, i.e. treatment starts at this check-up, set treatment start time $T^* \leftarrow 0$.
- Compute $\lambda_0 \leftarrow \text{logit}^{-1}(\gamma_0 + \gamma_2 A_0)$, the probability of failure in the interval $(0, 1]$ conditional on survival up to time 0.
- If $\lambda_0 \geq U_0$, the subject failed in the interval $(0, 1]$ and set $Y_1 \leftarrow 1$. Else the subject survived and set $Y_1 \leftarrow 0$.
- For $t = 1, \dots, T$, while the subject is still alive,
 - Draw $U_t \leftarrow \min\{1, \max\{0, U_{t-1} + N(0, 0.05)\}\}$, a perturbation of U_{t-1} restricted to $[0, 1]$.
 - If $t \not\equiv 0 \pmod{k}$, i.e. t is not a check-up time, $L_t \leftarrow L_{t-1}$ and $A_t \leftarrow A_{t-1}$. Else:
 - * If treatment started at the previous check-up, $L_t \leftarrow \max(0, L_{t-1} + 150 + N(100(U_t - 2), 50))$. Else $L_t \leftarrow \max(0, L_{t-1} + N(100(U_t - 2), 50))$. The Gaussian drift term is included such that the closer U_t is to 0, the stronger the negative drift in CD4.
 - * If treatment has not already started, draw treatment decision $A_t \sim \text{Bernoulli}(p_{A_t})$, where $p_{A_t} = \text{logit}^{-1}(\theta_0 + \theta_1 t + \theta_2(L_t - 500))$. Else $A_t \leftarrow 1$ (subject remains on treatment once initiated.)
 - * If $A_t = 1$ and $A_{t-k} = 0$, i.e. treatment starts at this check-up, set treatment start time $T^* \leftarrow t$.
 - Compute $\lambda_t \leftarrow \text{logit}^{-1}(\gamma_0 + \gamma_1[(1 - A_t)t + A_t T^*] + \gamma_2 A_t + \gamma_3 A_t(t - T^*))$, the probability of failure in the interval $(t, t + 1]$ conditional on survival up to time t .
 - Compute $S(t) = \prod_{\tau=0}^t (1 - \lambda_\tau)$, the probability of survival up to time $t + 1$. If $1 - S(t) \geq U_0$, the subject failed in the interval $(t, t + 1]$ and set $Y_{t+1} = 1$. Else the subject survived and set $Y_{t+1} = 0$.

The complete sampling algorithm for the modified model is given in Appendix A, and Appendix B proves that the generated data follow our desired MSM (6). A DAG representing this model for the case $T = 1, k = 1$ is given in Figure 4. Dot arrowheads represent deterministic relationships, circled nodes represent latent variables. Arrows from Y_t to future variables, representing dependence on prior survival, are omitted for clarity.

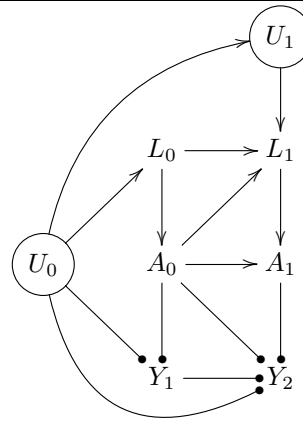


Figure 4. DAG corresponding to the data-generating algorithm for the case $T = 1, k = 1$.

From Figure 4, it is clear that the process exhibits time-dependent confounding. There is confounding due to U_0 being a common ancestor of \bar{A} (via \bar{L}) and \bar{Y} . There is potential for selection bias conditioning on L_1 induces an association between U_1 and A_0 , and hence between \bar{Y} and A_0 . Finally the only mediating pathways are such as $A_0 \rightarrow L_1 \rightarrow A_1 \rightarrow Y_2$. Furthermore, \bar{L} is sufficient to adjust for confounding, since A_t is independent of \bar{U}_t given $(\bar{L}_t, \bar{A}_{t-1})$.

The key mechanism by which this sampling algorithm generates confounded data from this MSM is the way in which \bar{A} and \bar{Y} depend on U_0 . The $U_0 \rightarrow \bar{Y}$ relationship is due to U_0 acting as a probability threshold, against which $S(t)$, the probability of surviving the time interval $(0, t + 1]$, is compared at each t . Hence, the subject fails at the point at which $1 - S(t)$ exceeds U_0 . The mechanism makes sense intuitively; a value of U_0 close to 1 indicates good health, and means it is less likely that $1 - S(t)$ will exceed U_0 . Thus failure is unlikely when baseline general health is good.

It is straightforward to demonstrate that using U_0 in this way does indeed result in the algorithm generating survival processes with interventional failure probabilities $\lambda_t^{\bar{a}}$. One such demonstration is as follows: Given a strategy \bar{a} , we note that λ_t in the algorithm is equal to $\lambda_t^{\bar{a}}$ from (6) when $\bar{A}_t = \bar{a}_t$. Thus, we have

$$\{Y_{t+1} = 1, \bar{Y}_t = 0\} \Leftrightarrow \left\{ 1 - \prod_{\tau=0}^t (1 - \lambda_{\tau}^{\bar{a}}) \geq U_0 \geq 1 - \prod_{\tau=0}^{t-1} (1 - \lambda_{\tau}^{\bar{a}}) \right\}.$$

This event occurs with probability

$$\begin{aligned} & P \left(U_0 \leq 1 - \prod_{\tau=0}^t (1 - \lambda_{\tau}^{\bar{a}}) \right) - P \left(U_0 \leq 1 - \prod_{\tau=0}^{t-1} (1 - \lambda_{\tau}^{\bar{a}}) \right) \\ &= 1 - \prod_{\tau=0}^t (1 - \lambda_{\tau}^{\bar{a}}) - \left(1 - \prod_{\tau=0}^{t-1} (1 - \lambda_{\tau}^{\bar{a}}) \right) \\ &= \lambda_t^{\bar{a}} \prod_{\tau=0}^{t-1} (1 - \lambda_{\tau}^{\bar{a}}), \end{aligned}$$

which is the probability of surviving the first $t - 1$ intervals, and failing in the t -th, as required.

Note that the corresponding *observational hazard* $\lambda_t^{obs}(\bar{A}_t = \bar{a}_t) = P(Y_{t+1} = 1 | Y_t = 0, \bar{A}_t = \bar{a}_t)$ cannot be so easily obtained. It would require one to evaluate the above probabilities with respect to the *conditional* distribution of U_0 given $\bar{A}_t = \bar{a}_t$, which is different from the marginal distribution of U_0 , since U_0 and \bar{A} are not independent (see also Appendix B for more details).

3.3. Comparison with Other Approaches

The above data-generating algorithm is inspired by the one used in Bryan et al. [10] and the chosen MSM hazard (6) is the same. However, their algorithm fixes the whole sequence $\{L_t\}$ at the start as a deterministic function of U_0 . Hence present treatment A_t cannot affect future L_{t+k} , which is unrealistic and means that we cannot examine whether a given method appropriately deals with the problems due to treatment affecting future covariates. Our algorithm extends the approach of Bryan et al. in this crucial aspect.

Furthermore, Young et al. [11, 12] propose an algorithm to simulate data from an MSM that allows a continuous survival outcome. The two papers focus on describing and proving a set of conditions under which a continuous-time proportional hazards MSM, a Structural Nested Accelerated Failure Time Model (SNAFTM) and a Structural Nested Cumulative Failure Time Model (SNCFTM) approximately coincide, so that the performance of IPTW can be compared with g-estimation. The MSM we have chosen here (6) is more complex than that in Young et al. and does not satisfy the conditions under which the three models coincide, for instance the effect of treatment start is allowed to depend on the time delay. The data-generating approach of Young et al., however, is essentially the same as ours; in particular they allow future covariates to depend on both past treatments and a latent variable representing the counterfactual baseline survival time T_0 under the strategy of ‘never treat’, which obviously predicts the outcome variable; this is analogous to our U_0 . The authors do not explicitly prove their algorithm, but our proof in Appendix B can easily be adapted, by changing the final step to demonstrate the validity of the respective MSMs.

The simulations carried out by Daniel et al. [9] used linear models, which are collapsible and therefore do not have the problem of incompatibility between the marginal model and the conditional distributions used to simulate the data. The approach of Xiao et al. [32] considers survival outcomes but does not actually simulate data from an MSM as it is based on a *conditional* model which includes a L_t -term in the hazard. They argue that choosing a very small rate of events means that the model is near-collapsible, i.e. that the conditional and marginal parameter values are almost identical. However, this further relies also on the absence of a latent process $\{U_t\}$ which means that selection bias cannot occur and hence cannot be investigated with their simulations.

3.4. Application to a Simulation Study

With our data-generating algorithm established, we wish to evaluate and compare the performance of the following methods for a finite sample size. Note that throughout we only consider the case where censoring occurs at end of follow-up, for simplicity. We apply four different pooled logistic regressions to the data. (Pooled logistic regression is the discrete analogue of proportional hazards regression, in which all subject–time points are pooled into one set of observations [36].) These are:

1. An unweighted regression which nominally fits model (6) using data on $\{A_t, Y_t\}$ but without any adjustment for confounding, i.e. not using $\{L_t\}$.
2. An unweighted, adjusted regression which nominally fits model (6) but adds an extra term $\gamma_4 \frac{L_t}{100}$ in an attempt to adjust for confounding by $\{L_t\}$.
3. A second unweighted, adjusted regression which nominally fits model (6) but adds an extra term $\gamma_4 \frac{L^{\text{av}}}{100}$, where L^{av} is the average over all (past and future) observed values of $\{L_t\}$. The idea here is that L^{av} may be regarded as a summary measure for the CD4 counts, and it should be biased as it conditions on future covariates.
4. A weighted regression which fits (6) and uses the stabilised weights of (5).

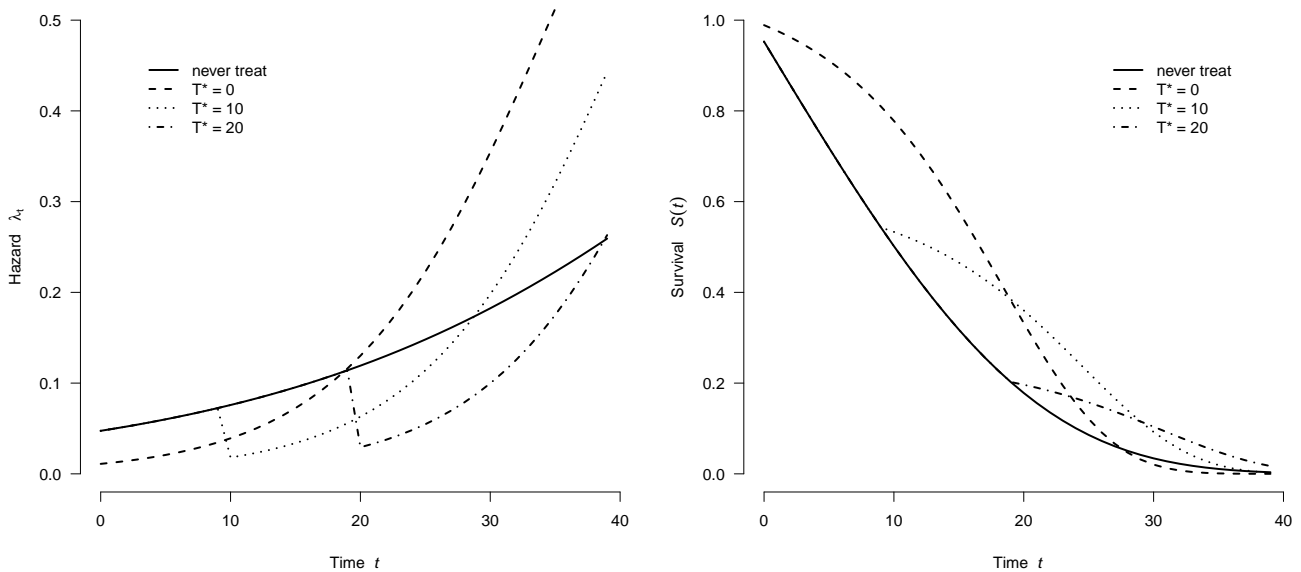


Figure 5. Structural hazard and survival curves for four possible treatment regimes.

3.4.1. *Simulation Settings.* We simulate data from our generating model with the following parameter values.

- Follow-up occurs over $T = 40$ time points, with check-ups every ($k = 5$)th point. These values match those used in Bryan et al. [10]
- The parameters for the conditional distributions of treatment are

$$(\theta_0, \theta_1, \theta_2) = (-0.405, 0.0205, -0.00405).$$

These values calibrate the logistic function such that $P(A_0 = 1 | L_0 = 500) = 0.4$, $P(A_0 = 1 | L_0 = 400) = 0.5$ and $P(A_{10} = 1 | L_{10} = 500) = 0.45$. This calibration generally ensures that treatment assignment probabilities close to 0 or 1 do not occur. Such a requirement could be relaxed when we wish to specifically investigate the stability of IPTW when there are rare treatment decisions.

- The parameters of the MSM hazard $\lambda_t^{\bar{a}}$ are

$$(\gamma_0, \gamma_1, \gamma_2, \gamma_3) = (-3, 0.05, -1.5, 0.1).$$

These values satisfy the non-trivial case described in Section 3.2.3, in which initiating treatment provides a one-time decrease in λ_t , but also increases the rate at which λ_t increases with time.

The above settings result in structural hazard and survival functions as shown in Figure 5, and the rate of failure events between check-up times is around 12% (by changing the intercept γ_0 this can be reduced; however, even with a rate of at most 4% we found qualitatively analogous results as given below).

3.4.2. *Results.* We generated 100 replications each of a study with $n = 1000$ subjects, a moderate sample size (for comparison, the Swiss HIV Cohort Study has a sample size of 3245). For each study, we fit the four regressions. The mean and standard deviation of each regression's estimates of γ are taken over all replications. The results (omitting the intercept γ_0) are given in Table 1. Note that the coefficient γ_4 is defined such that one unit corresponds to 100 cells/ μ L.

As expected, all unweighted estimates are biased, while the IPTW estimate appears to perform well even for this moderate sample size. We see that the unweighted unadjusted estimates, especially of γ_1, γ_2 , are biased towards the null. This is plausible because patients with low CD4 counts are more likely to start treatment, but also more likely to have shorter survival, hence treatment will appear less effective. Also, as expected the regression that adjusts for L^{av} is more biased than the one which adjusts for L_t . We note the non-zero regression coefficients γ_4 which indicate that the covariate process $\{L_t\}$, and moreso L^{av} , is indeed predictive of the outcome.

When interpreting the results of the two regression adjustments, we have to keep in mind that our target parameter is a *marginal* causal effect, while these two regressions are estimating conditional coefficients (given the covariate L_t or L^{av}). As the model is not collapsible these cannot be expected to be the same, even if $\{L_t\}$ and $\{A_t\}$ were independent processes, as we will see below.

3.4.3. Isolating Sources of Bias. In order to isolate separate potential sources of bias, we now consider the following special cases of the above data-generating algorithm obtained by removing certain dependencies corresponding to edges in the underlying DAG.

- (i) We can break the dependence of L_t on \bar{A}_{t-1} by drawing $L_t \leftarrow \max\{0, L_{t-1} + N(100(U_t - 2), 50)\}$. For the HIV scenario this would mean that treatment no longer has a beneficial (or any other) effect on CD4 count, but nevertheless improves survival chances. Hence, U_0 is still confounding A_t and survival via L_t , but as the latter is not affected by past treatment there will be no selection bias nor mediation.
- (ii) We can break the dependence of L_t on \bar{U}_t by drawing $L_t \leftarrow \max\{0, L_{t-1} + 150A_{t-k}(1 - A_{t-k-1}) + N(-150, 50)\}$. For the HIV scenario this would mean that a subject's CD4 count at a particular time does not depend on their latent general health. Now, there is no confounding at all, and selection bias cannot occur either.
- (iii) We can break the dependence of A_t on \bar{L}_t by drawing $A_t \sim \text{Bernoulli}(0.5)$ whenever $A_{t-1} = 0$. This clearly corresponds to performing a randomised controlled trial, in which each subject who has not already initiated treatment at a particular time is randomly put on treatment with probability 0.5. Again, there is no confounding in this scenario, but there is potential for selection bias.
- (iv) Finally we combine (i) and (iii) of the above. Graphically this removes any arrows from A_{t-1} to L_t , as well as arrows from L_t to A_t . In other words the processes $\{A_t\}$ on the one hand and $\{L_t, U_t\}$ on the other hand are entirely independent.

With these modified generating processes established, we may once again apply the same four pooled logistic regressions: unweighted, L_t -adjusted, L^{av} -adjusted and IPTW-weighted, each applied to 100 replications of 1000 subjects of each data-generating process. Means and standard deviations are taken across the replications. The results are given in Tables 3–5. We note that the results for the IPTW approach are very good throughout, as can be expected from its theoretical properties; we know that it is consistent in all scenarios and the finite, rather small, sample size does not appear to cause a problem.

Let us first consider the performance for the last data-generating process, see Table 3. As we would expect, neither the unweighted unadjusted nor the IPTW estimators are biased; in fact their standard deviations are very similar. This is because the stabilised weights will all be close to 1, numerator and denominator being essentially identical except for sampling variability. The unweighted unadjusted estimators can be regarded as using the true weights. It is known that using estimated weights even when the true weights are known yields more efficient estimators, which explains why the variability of IPTW appears even smaller than of the unweighted unadjusted approach [1, 10]. While there is no structural bias that can possibly affect the two regression-adjusted estimators for these simulated data, the estimates appear biased due to the lack of collapsibility and ensuing misspecification of the regression model. As expected, higher CD4 counts are associated with longer survival, however all other effects appear biased away from null.

In Table 2 we see the performance on data where treatment does not affect future covariates. The unweighted unadjusted regression exhibits the exact same bias as in Table 1; the same method is not biased in Tables 4 and 5 which confirms

that this method is only affected by confounding. The two regression adjustments are also affected by the remaining confounding in Table 2, even though they are trying to adjust for it. These two methods neither seem to estimate the true marginal parameters, nor the conditional ones of Table 3.

The L_t -adjusted regression appears consistent for data where $\{U_t\}$ does not affect $\{L_t\}$ in Table 4, but is further affected by selection bias in Table 5. In the former case, we have that L_t does not predict the outcome given A_t as can be seen from the zero γ_4 coefficient. Hence the model is correctly specified and there is no source of bias, so that the estimates are consistent. If we were able to simulate data from an MSM such that L_t predicts Y other than through future treatment we would expect to see a bias here.

In contrast, the L^{av} -adjusted regression struggles in all cases. Including L^{av} as covariate at any given time t implies conditioning on future covariates which predict Y via future treatment, and we can see from Table 4 that this leads to bias of the treatment effect towards the null, i.e. some of the effect of starting treatment is wrongly attributed to L^{av} instead. In fact the sign of γ_4 means that large average CD4 count is associated with shorter survival which can only be explained by its predicting that future treatment is less likely.

4. Discussion and Outlook

We have discussed a general approach for simulating longitudinal data from a given MSM while ensuring that the data exhibit time-dependent confounding, including the potential for bias due to past treatment affecting future covariates. This is especially relevant in cases where the desired model is not collapsible, such as in typical survival analyses using proportional hazards or pooled logistic regressions, because the data-generating process suggested by a corresponding DAG relies on conditional model specification which may not be compatible with the desired marginal model.

The sampling algorithm presented here is quite flexible; there is scope to examine and develop the finer details of it. This may involve altering aspects such as the distribution of U_t given \bar{U}_{t-1} , with the aim of generating data which more closely match that of the Swiss HIV Cohort Study, or perhaps examining different longitudinal data structures. However the key mechanism by which time-dependent confounding is included must remain intact.

As MSMs and IPTW are gaining popularity, especially in the context of survival or time-to-event analyses, it is important to be able to simulate data from a known true model so as to evaluate finite-sample properties and compare competing approaches. We have illustrated the shortcomings of unweighted and regression-adjusted methods with a small simulation study. To address weaknesses of IPTW, future simulation studies could consider problems due to large or unstable weights, censoring, misspecification of the structural model and/or the treatment model, especially in situations when the set of covariates to be adjusted-for is high-dimensional, as well as mild violations of the no unmeasured confounding assumption. Our data-generating algorithm can easily be extended to such cases. A shortcoming however is that the dependence of the outcome on any covariates can only indirectly be specified, and hence mediation as well as effect modification are difficult to model.

The problems addressed in this paper may lead some researchers to conclude that MSMs are not a natural model class as it appears more intuitive and flexible to generate simulated data from the conditional distributions of a DAG factorisation. In such a case we find it still valuable and important to consider the induced $P(Y|do(\bar{a}))$ which can be obtained from the truncated factorisation formula, integrating out the time-dependent covariates and latent processes; it may then even be possible to work out what the misspecified MSM is estimating. In cases where this is analytically difficult it is relatively straightforward to simulate the corresponding data similar to [37].

Finally, MSMs are now becoming more popular for analysing the effect of *dynamic* treatment strategies by artificial censoring [38, 5, 39, 40]. The aim is to determine the optimal strategy within a given set of possible ones. A dynamic strategy takes the treatment decisions as functions of the time-varying covariate, e.g. 'start treatment with HAART when CD4 drops below 600'. The time at which a given patient will receive treatment is not then predetermined, as it might

happen that their CD4 drops below 600 very early on, later, or never. In future work we will consider how to adapt our data-generating algorithm to MSMs for such dynamic strategies.

Acknowledgement

We thank Rhian Daniel for helpful discussions, and the Associate Editor and two anonymous referees for insightful comments and questions on the original submission. William Havercroft acknowledges funding by the Engineering and Physical Sciences Research Council (EPSRC). Vanessa Didelez acknowledges financial support from the Leverhulme Trust (grant RF-2011-320).

References

1. Robins J. Marginal structural models versus structural nested models as tools for causal inference. *Statistical Models in Epidemiology, the Environment and Clinical Trials*, Halloran E, Berry D (eds.). Springer, 1999; 95.
2. Robins J, Hernán M, Brumback B. Marginal structural models and causal inference in epidemiology. *Epidemiology* 2000; **11**(5):550.
3. Sterne J, Hernán M, Ledergerber B, Tilling K, Weber R, Sendi P, Rickenbach M, Robins J, Egger M. Long-term effectiveness of potent antiretroviral therapy in preventing AIDS and death: a prospective cohort study. *The Lancet* 2005; **366**(9483):378–384.
4. Robins J, Orellana L, Rotnitzky A. Estimation and extrapolation of optimal treatment and testing strategies. *Statistics in Medicine* 2008; **27**(23):4678–4721.
5. Sterne J, May M, Costagliola D, de Wolf F, Phillips A, Harris R, Funk M, Geskus R, Gill J, Dabis F, *et al.*. Timing of initiation of antiretroviral therapy in AIDS-free HIV-1-infected patients: a collaborative analysis of 18 HIV cohort studies. *The Lancet* 2009; **373**(9672):1352–1363.
6. Rosthøj S, Fullwood C, Henderson R, Stewart S. Estimation of optimal dynamic anticoagulation regimes from observational data: a regret-based approach. *Statistics in Medicine* 2006; **25**(24):4197–4215.
7. Henderson R, Ansell P, Alshibani D. Regret-regression for optimal dynamic treatment regimes. *Biometrics* 2009; **66**(4):1192–1201.
8. Robins J. A new approach to causal inference in mortality studies with a sustained exposure period—application to control of the healthy worker survivor effect. *Mathematical Modelling* 1986; **7**(9-12):1393–1512.
9. Daniel R, Cousens S, De Stavola B, Kenward M, Sterne J. Tutorial in biostatistics: methods for dealing with time-varying confounding 2011; (submitted).
10. Bryan J, Yu Z, Van Der Laan M. Analysis of longitudinal marginal structural models. *Biostatistics* 2004; **5**(3):361.
11. Young J, Hernán M, Picciotto S, Robins J. Simulation from structural survival models under complex time-varying data structures. *JSM Proceedings, Section on Statistics in Epidemiology, Denver, CO: American Statistical Association* 2008; .
12. Young J, Hernán M, Picciotto S, Robins J. Relation between three classes of structural models for the effect of a time-varying exposure on survival. *Lifetime data analysis* 2010; **16**(1):71–84.
13. Pearl J. *Causality: Models, Reasoning and Inference*. 2nd edn., Cambridge University Press, 2000.
14. Rubin D. Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of Educational Psychology* 1974; **66**:688–701.
15. Rubin D. Bayesian inference for causal effects: the role of randomization. *Annals of Statistics* 1978; **6**(1):34–68.
16. Dawid A. Causal inference without counterfactuals (with Discussion). *Journal of the American Statistical Association* 2000; **95**:407–448.
17. Dawid A. Counterfactuals, hypotheticals and potential responses: a philosophical examination of statistical causality. *Causality and Probability in the Sciences*, Russo F, Williamson J (eds.). College Publications, Texts In Philosophy Series Vol. 5: London, 2007; 503532.
18. Lauritzen S. *Graphical Models*. Oxford University Press, 1996.
19. Dawid A. Conditional independence in statistical theory (with Discussion). *Journal of the Royal Statistical Society* 1979; **41**:1–31.
20. Lauritzen S. Causal inference from graphical models. *Complex Stochastic Systems*, Barndorff-Nielsen O, Cox D, Klüppelberg C (eds.). chap. 2, CRC Press: London, 2000; 63–107.
21. Dawid A, Didelez V. Identifying the consequences of dynamic treatment strategies: A decision-theoretic overview. *Statistics Surveys* 2010; **4**:184–231.
22. Pearl J. Graphical models, causality and interventions. *Statistical Science* 1993; **8**:266269.
23. Dawid A. Influence diagrams for causal modelling and inference. *International Statistical Review* 2002; **70**:161–189. Corrigenda, *ibid.*, 437.
24. Spirtes P, Glymour C, Scheines R. *Causation, Prediction and Search*. Springer: New York, 1993.
25. Pearl J. Causal diagrams for empirical research. *Biometrika* 1995; **82**(4):669–710.
26. Robins J. Causal inference from complex longitudinal data. *Latent Variable Modeling and Applications to Causality, Lecture Notes in Statistics*, vol. 120, Berkane M (ed.). Springer: New York, 1997; 69–117.

27. Pearl J, Robins J. Probabilistic evaluation of sequential plans from causal models with hidden variables. *Proceedings of the 11th Conference on Uncertainty in Artificial Intelligence*, Besnard P, Hanks S (eds.), Morgan Kaufmann: San Francisco, 1995; 444–453.
28. Hernán M, Hernández-Díaz S, Robins J. A structural approach to selection bias. *Epidemiology* 2004; **15**(5):615.
29. Cole S, Platt R, Schisterman E, Chu H, Westreich D, Richardson D, Poole C. Illustrating bias due to conditioning on a collider. *International Journal of Epidemiology* 2009; **39**(2):417–420.
30. Greenland S. Quantifying biases in causal models: classical confounding vs collider-stratification bias. *Epidemiology* 2003; **14**(3):300.
31. Robins J, Blevins D, Ritter G, Wulfsohn M. G-estimation of the effect of prophylaxis therapy for pneumocystis carinii pneumonia on the survival of aids patients. *Epidemiology* 1992; **3**(4):319–336.
32. Xiao Y, Abrahamowicz M, Moodie E. Accuracy of conventional and marginal structural Cox model estimators: A simulation study. *The International Journal of Biostatistics* 2010; **6**(2):13.
33. Naimi A, Cole S, Westreich D, Richardson D. A comparison of methods to estimate the hazard ratio under conditions of time-varying confounding and nonpositivity. *Epidemiology* 2011; **22**(5):718.
34. Westreich D, Cole S, Tien P, Chmiel J, Kingsley L, Funk M, Anastos K, Jacobson L. Time scale and adjusted survival curves for marginal structural cox models. *American journal of epidemiology* 2010; **171**(6):691–700.
35. Røysland K. A martingale approach to continuous-time marginal structural models. *Bernoulli* 2011; **17**(3):895–915.
36. D'Agostino R, Lee M, Belanger A, Cupples L, Anderson K, Kannel W. Relation of pooled logistic regression to time dependent Cox regression analysis: the Framingham Heart Study. *Statistics in Medicine* 1990; **9**(12):1501–1515.
37. Daniel R, De Stavola B, Cousens S. gformula: Estimating causal effects in the presence of time-varying confounding or mediation using the g-computation formula. *The Stata Journal* 2011; (to appear).
38. Hernán M, Lanoy E, Costagliola D, Robins J. Comparison of dynamic treatment regimes via inverse probability weighting. *Basic & Clinical Pharmacology & Toxicology* 2006; **98**(3):237–242.
39. Orellana L, Rotnitzky A, Robins J. Dynamic regime marginal structural mean models for estimation of optimal dynamic treatment regimes, part i: Main content. *The International Journal of Biostatistics* 2010; **6**(2).
40. Cain L, Robins J, Lanoy E, Logan R, Costagliola D, Hernán MA. When to start treatment? a systematic approach to the comparison of dynamic regimes using observational data. *The International Journal of Biostatistics* 2010; **6**(2).

A. The Full Sampling Algorithm

for subject $i = 1, \dots, n$ **do**

$$U_{0,i} \sim U[0, 1]$$

$$\epsilon_{0,i} \sim N(0, 20)$$

$$L_{0,i} \leftarrow F_{\Gamma(3,154)}^{-1}(U_{0,i}) + \epsilon_{0,i}$$

$$A_{-1,i} \leftarrow 0$$

$$A_{0,i} \sim \text{Bernoulli}(\text{logit}^{-1}(\theta_0 + \theta_2(L_{0,i} - 500)))$$

if $A_{0,i} = 1$ **then**

$$T^* \leftarrow 0$$

end if

$$\lambda_{0,i} \leftarrow \text{logit}^{-1}(\gamma_0 + \gamma_2 A_{0,i})$$

if $\lambda_{0,i} \geq U_{0,i}$ **then**

$$Y_{1,i} \leftarrow 1$$

else

$$Y_{1,i} \leftarrow 0$$

end if

for $t = 1, \dots, T$ **do**

while $Y_{t,i} = 0$ **do**

$$\Delta_{t,i} \sim N(0, 0.05)$$

$$U_{t,i} \leftarrow \min\{1, \max\{0, U_{t-1,i} + \Delta_{t,i}\}\}$$

if $t \not\equiv 0 \pmod{k}$ **then**


```

 $L_{t,i} \leftarrow L_{t-1,i}$ 
 $A_{t,i} \leftarrow A_{t-1,i}$ 
else
 $\epsilon_{t,i} \sim N(100(U_{t,i} - 2), 50)$ 
 $L_{t,i} \leftarrow \max\{0, L_{t-1,i} + 150A_{t-k,i}(1 - A_{t-k-1,i}) + \epsilon_{t,i}\}$ 
if  $A_{t-1,i} = 0$  then
 $A_{t,i} \sim \text{Bernoulli}(\text{logit}^{-1}(\theta_0 + \theta_1 t + \theta_2(L_{t,i} - 500)))$ 
else
 $A_{t,i} \leftarrow 1$ 
end if
if  $A_{t,i} = 1$  and  $A_{t-k,i} = 0$  then
 $T^* \leftarrow t$ 
end if
end if
 $\lambda_{t,i} \leftarrow \text{logit}^{-1}(\gamma_0 + \gamma_1[(1 - A_{t,i})t + A_{t,i}T^*] + \gamma_2 A_{t,i} + \gamma_3 A_{t,i}(t - T^*))$ 
if  $1 - \prod_{\tau=0}^t (1 - \lambda_{\tau,i}) \geq U_{0,i}$  then
 $Y_{t+1,i} = 1$ 
else
 $Y_{t+1,i} = 0$ 
end if
end while
end for
end for

```

B. Sketch Proof that the Generating Model Simulates from the Desired MSM

We begin by considering the factorisation of the joint distribution of all the variables in the model. Then, as in Section 2.2 we obtain the truncated factorisation formula under an intervention following treatment plan \bar{a} , and show that it results in the MSM (6). The reasoning is essentially the same as for the example with a single time point, which yields (4).

Note that the set of parents of each L_t and A_t depends on whether or not t is a check-up time. If t is *not* a check-up time, the only parent of L_t is L_{t-1} . Otherwise, L_t has the additional parents A_{t-k} , A_{t-k-1} and U_t . Similarly if t is not a check-up time, the only parent of A_t is A_{t-1} , otherwise it has the additional parent L_t . However since the removal of parents represents the addition of conditional independencies, we may proceed without loss of generality by assuming the larger set of parents for all L_t and A_t .

In the following we have that all vectors \bar{Y} consist of a sequence of zeros followed by a one. Let $\tilde{Y} = \min\{t : Y_t = 1\}$ be the time where the event (e.g. death) occurs. All remaining variables are undefined for $t \geq \tilde{Y}$. The joint distribution factorises as follows (where variables with negative index are defined to be the empty set).

$$P(\bar{L}, \bar{A}, \bar{Y}, \bar{U}) = \prod_{\tau=0}^{\tilde{Y}-1} \{P(U_\tau | U_{\tau-1}, Y_\tau = 0) \times P(L_\tau | L_{\tau-1}, A_{\tau-k}, A_{\tau-k-1}, U_\tau, Y_\tau = 0) \\ \times P(A_\tau | A_{\tau-1}, L_\tau, Y_\tau = 0) \times P(Y_{\tau+1} | \bar{A}_\tau, U_0, Y_\tau = 0)\}.$$

Intervening to fix treatment \bar{A} at \bar{a} yields the truncated factorisation formula

$$P(\bar{L}, \bar{Y}, \bar{U} | do(\bar{a})) = \prod_{\tau=0}^{\tilde{Y}-1} \{P(U_\tau | U_{\tau-1}, Y_\tau = 0) \times P(L_\tau | L_{\tau-1}, A_{\tau-k} = a_{\tau-k}, A_{\tau-k-1} = a_{\tau-k-1}, U_\tau, Y_\tau = 0) \times P(Y_{\tau+1} | \bar{A}_\tau = \bar{a}_\tau, U_0, Y_\tau = 0)\}.$$

Integrating out L_0, \dots, L_T and U_1, \dots, U_T yields

$$P(\bar{Y}, U_0 | do(\bar{a})) = P(U_0) \prod_{\tau=0}^{\tilde{Y}-1} P(Y_{\tau+1} | \bar{A}_\tau = \bar{a}_\tau, U_0, Y_\tau = 0).$$

Note that equivalent operations can be derived if we replace \bar{Y} with the events $\tilde{Y} > t$, so that equivalently we have

$$P(\tilde{Y} > t, U_0 | do(\bar{a})) = P(U_0) \prod_{\tau=0}^T P(\tilde{Y} > \tau | \bar{A}_\tau = \bar{a}_\tau, U_0, \tilde{Y} \geq \tau).$$

Let $\lambda_\tau^{\bar{a}} = \lambda_\tau(do(\bar{a}))$ as in (6). Then, in our algorithm $\{\tilde{Y} = t + 1\} \Leftrightarrow \{\prod^{t-1} (1 - \lambda_\tau^{\bar{a}}) \geq 1 - U_0 \geq \prod^t (1 - \lambda_\tau^{\bar{a}})\}$, and $U_0 \sim U[0, 1]$, so we obtain that

$$P(\tilde{Y} > t | do(\bar{a})) = \prod_{\tau=0}^t (1 - \lambda_\tau^{\bar{a}})$$

which is the survival function of our desired MSM under an intervention $do(\bar{a})$ as required.

Note that apart from the very last step, the above proof is for general conditional distributions. It therefore justifies more general versions of the algorithm, where for instance we could have non-binary treatment, or treatment that can repeatedly be switched on or off, as well as multivariate covariate processes etc. The validity of the algorithm of Young et al. [12] can be shown just as above, but replacing the final step so that U_0 is transformed into an exponentially distributed variable, which is then further transformed according to an accelerated failure time model depending on the actual treatment $\bar{A}_t = \bar{a}_t$.

Tables

	True	Unweighted (s.d.)	L_t -adjusted (s.d.)	L^{av} -adjusted (s.d.)	Weighted (s.d.)
γ_1	0.050	0.009 (0.006)	0.072 (0.008)	0.230 (0.012)	0.050 (0.010)
γ_2	-1.500	-0.402 (0.111)	-1.917 (0.143)	-2.374 (0.173)	-1.526 (0.134)
γ_3	0.100	0.113 (0.005)	0.097 (0.007)	0.287 (0.011)	0.103 (0.009)
γ_4	—	—	-0.978 (0.032)	-1.375 (0.056)	—

Table 1. Monte Carlo means and standard deviations of estimates produced by different regression methods.

	True	Unweighted (s.d.)	L_t -adjusted (s.d.)	L^{av} -adjusted (s.d.)	Weighted (s.d.)
γ_1	0.050	0.009 (0.006)	0.066 (0.008)	0.247 (0.012)	0.051 (0.011)
γ_2	-1.500	-0.399 (0.109)	-3.019 (0.166)	-2.411 (0.153)	-1.516 (0.133)
γ_3	0.100	0.113 (0.004)	0.077 (0.007)	0.245 (0.008)	0.102 (0.010)
γ_4	—	—	-1.184 (0.040)	-1.320 (0.046)	—

Table 2. Pooled logistic regression results for data simulated according to (i), where L_t does not depend on \bar{A}_{t-1} .

	True	Unweighted (s.d.)	L_t -adjusted (s.d.)	L^{av} -adjusted (s.d.)	Weighted (s.d.)
γ_1	0.050	0.049 (0.008)	0.068 (0.008)	0.249 (0.011)	0.050 (0.006)
γ_2	-1.500	-1.496 (0.116)	-3.557 (0.178)	-3.654 (0.206)	-1.503 (0.098)
γ_3	0.100	0.100 (0.004)	0.103 (0.007)	0.278 (0.009)	0.100 (0.004)
γ_4	—	—	-1.127 (0.043)	-1.272 (0.047)	—

Table 3. Pooled logistic regression results for data simulated according to (iv), where $\{A_t\}$ and $\{L_t, U_t\}$ are independent.

	True	Unweighted (s.d.)	L_t -adjusted (s.d.)	L^{av} -adjusted (s.d.)	Weighted (s.d.)
γ_1	0.050	0.051 (0.009)	0.051 (0.009)	0.033 (0.009)	0.052 (0.017)
γ_2	-1.500	-1.523 (0.104)	-1.523 (0.110)	-1.184 (0.111)	-1.525 (0.157)
γ_3	0.100	0.100 (0.005)	0.100 (0.005)	0.107 (0.005)	0.101 (0.006)
γ_4	—	—	-0.000 (0.017)	0.167 (0.017)	—

Table 4. Pooled logistic regression results for data simulated according to (ii), where L_t does not depend on \bar{U}_t .

	True	Unweighted (s.d.)	L_t -adjusted (s.d.)	L^{av} -adjusted (s.d.)	Weighted (s.d.)
γ_1	0.050	0.051 (0.008)	0.074 (0.009)	0.232 (0.013)	0.051 (0.007)
γ_2	-1.500	-1.504 (0.111)	-2.541 (0.143)	-3.480 (0.195)	-1.502 (0.097)
γ_3	0.100	0.101 (0.004)	0.119 (0.007)	0.321 (0.012)	0.100 (0.004)
γ_4	—	—	-0.943 (0.029)	-1.335 (0.053)	—

Table 5. Pooled logistic regression results for data simulated according to (iii), where A_t does not depend on \bar{L}_t .