

# Marginal Structural Models versus Structural Nested Models as Tools for Causal Inference

James M. Robins

Professor of Epidemiology and Biostatistics  
Harvard School of Public Health  
677 Huntington Avenue  
Boston, MA 02115  
robins@hsph.harvard.edu

Abstract: Robins (1993, 1994, 1997, 1998ab) has developed a set of causal or counterfactual models, the structural nested models (SNMs). This paper describes an alternative new class of causal models—the (non-nested) marginal structural models (MSMs). We will then describe a class of semiparametric estimators for the parameters of these new models under a sequential randomization (i.e., ignorability) assumption. We then compare the strengths and weaknesses of MSMs versus SNMs for causal inference from complex longitudinal data with time-dependent treatments and confounders. Our results provide an extension to continuous treatments of propensity score estimators of an average treatment effect.

## 1 Introduction

Robins (1993, 1994, 1997, 1998ab) has developed a set of causal or counterfactual models, the structural nested models (SNMs). Robins (1998abcd) has recently described an alternative new class of causal models – the (non-nested) marginal structural models (MSMs). We describe a class of semiparametric estimators for the parameters of these new models under a sequential randomization (i.e., ignorability) assumption. We then compare the strengths and weaknesses of MSMs versus SNMs for causal inference from complex longitudinal data with time-dependent treatments and confounders. Two major strengths of MSMs compared to SNMs are as follows.

- MSMs can be used to provide semiparametric estimates of the causal effect of a time-dependent treatment on a binary outcome using models (e.g. logistic models) which naturally respect the fact that probabilities lie in the interval  $[0, 1]$ .
- MSMs cohere much more closely than do SNMs with models for the analysis of time-dependent treatments that are standardly used in the absence of time-dependent confounders. For example, in the absence of time-dependent confounders, a time-dependent Cox proportional hazards model for the effect of time-dependent treatment on a time-to-event (survival time) outcome is commonly employed. The MSMs provide a natural extension of the time-dependent proportional hazards model. Unlike the usual time-dependent Cox model, the marginal structural time-dependent Cox model can be used to obtain valid causal inferences for the effect of a time-varying treatment in the presence of time-varying confounding factors. [We remind the reader that, as discussed in Robins (1986), one cannot estimate the effect of a time-dependent treatment on survival in the presence of time-dependent confounding factors by using an ordinary time-dependent Cox model that adjusts for the time-dependent confounding factors since, in general, these time-dependent confounding factors will be both determinants of later treatment and affected by earlier treatment. Marginal structural Cox models overcome this deficiency.] Disadvantages of MSMs are discussed later. The

relationship of our approach to the propensity score approach of Rosenbaum and Rubin (1983) is considered in Section 4.1.

We now give a somewhat informal introduction to marginal structural models, and we report the results of a preliminary data analysis of AIDS Clinical Trial Group (ACTG) Trial 002 using MSMs. We begin with the following simple setting. Consider a study of AIDS patients. Let  $A(t)$  be the dose of a treatment of interest, say AZT, at time  $t$  with time measured as days since start of follow-up. Let  $Y$  be an outcome of interest measured at end-of-follow-up at time  $K + 1$ . Our goal is to estimate the causal effect of the time-dependent treatment  $A(t)$  on the mean of  $Y$ . Let  $\bar{A}(t) = \{A(u); 0 \leq u \leq t\}$  be treatment history through  $t$  and let  $\bar{L}(t) = \{L(u); 0 \leq u \leq t\}$  be the history through  $t$  of a vector of relevant prognostic factors  $L(u)$  for (i.e., predictors of)  $Y$ , such as CD4 lymphocyte count, white blood count (WBC), hematocrit, age, gender, etc. Suppose  $Y$  is a dichotomous outcome (e.g.,  $Y = 1$  if HIV RNA is detectable in the blood and zero otherwise), and we entertain a model that says the mean of  $Y$  given AZT history,  $\bar{A} \equiv \bar{A}(K + 1)$ , is a linear logistic function of a subject's cumulative AZT dose. We write the model

$$E[Y | \bar{A}] = g(\bar{A}; \gamma)$$

where

$$g(\bar{A}; \gamma) = [1 + \exp\{-\gamma_1 - \gamma_2 cum(\bar{A})\}]^{-1}$$

and  $cum(\bar{A}) = \int_0^{K+1} A(t) dt$  is the subject's cumulative treatment. The maximum likelihood estimator (MLE) of  $\gamma$  can then be computed from the observed data  $O_i = (\bar{L}_i, \bar{A}_i, Y_i), i = 1, \dots, n$ , on the  $n$  study subjects using standard logistic regression software with  $Y$  as the Bernoulli outcome variable and  $cum(\bar{A})$  as the regressor. That is, the MLE of  $\gamma = (\gamma_1, \gamma_2)'$  maximizes  $\prod_{i=1}^n Lik_i(\gamma)$  with  $Lik_i(\gamma) = g(\bar{A}_i; \gamma)^{Y_i} [1 - g(\bar{A}_i; \gamma)]^{1-Y_i}$  being the likelihood contribution for a single subject. [Note  $Lik_i(\gamma)$  does not depend on the patient's prognostic factor history  $\bar{L}_i \equiv \bar{L}_i(K + 1)$ .] Alternatively, we could have used Bayesian methods by specifying a prior distribution for  $\gamma$  and then estimating  $\gamma$  by its posterior mean given the data. In reasonable large samples, the MLE and Bayes estimate will closely approximate one another.

**Causal Interpretation of Regression Parameters:** The question then is when does  $\gamma_2$  have an interpretation as the causal effect of treatment history on the mean of  $Y$ ? To approach this question, imagine that the decision to administer treatment at each time  $t$  were made totally at random by the treating physician. In that hypothetical case, giving treatment at time  $t$  is not expected to be associated with any measured or unmeasured prognostic factors (i.e., there would be no "confounding") and therefore  $\gamma_2$  would intuitively have a causal interpretation. Similarly,  $\gamma_2$  would keep its causal interpretation if the physician's decision were based only on the history of treatment prior to  $t$ . Whenever the conditional probability of receiving treatment on day  $t$  given past treatment and prognostic factors history (measured and unmeasured) depends only on past treatment history, we say the process is a "causally exogenous or ancillary treatment process". (A more formal mathematical definition is provided below.) It is well-recognized in the social sciences, econometrics, epidemiologic, and biostatistical literature that  $\gamma_2$  will have a causal interpretation if  $A(t)$  is a causally exogenous (or ancillary) covariate process. Randomized treatments like the one described above are causally exogenous treatments.

We say that a treatment  $A(t)$  is a "statistically exogenous or ancillary process" if the probability of receiving treatment at time  $t$  does not depend on the history of measured time-dependent prognostic

factors  $\bar{L}(t)$  up to  $t$  conditional on treatment history prior to  $t$ , i.e.,

$$\bar{L}(t) \coprod A(t) \mid \bar{A}(t-1),$$

where  $A \coprod B \mid C$  means that  $A$  is independent of  $B$  given  $C$ .

Note that a necessary condition for  $A(t)$  to be “causally exogenous” is for it to be “statistically exogenous.” However, that a process is “statistically exogenous” does not imply it is “causally exogenous,” because there may be unmeasured prognostic factors (i.e., confounders) that predict the probability of treatment  $A(t)$  at time  $t$  given past treatment history. We can test from the data whether  $A(t)$  is statistically exogenous but are unable to test whether a statistically exogenous process is causally exogenous.

Suppose  $A(t)$  is discrete and we can correctly model the probability  $f[a(t) \mid \bar{\ell}(t), \bar{a}(t-1)]$  of receiving treatment  $a(t)$  on day  $t$  as a function of past treatment  $\bar{a}(t-1)$  and measured prognostic factor history  $\bar{\ell}(t)$ . We could then measure the degree to which the treatment process is statistically non-exogenous through day  $t$  by the random quantity

$$\mathcal{W}(t) = \prod_{k=0}^t f[A(k) \mid \bar{A}(k-1), \bar{L}(k)] / f[A(k) \mid \bar{A}(k-1)] .$$

The numerator in each term in  $\mathcal{W}(t)$  is the probability that a subject received his own observed treatment at time  $k$ ,  $A(k)$ , given his past treatment and prognostic factor history. The denominator is the probability that a subject received his observed treatment conditional on his past treatment history but not further adjusting for his past prognostic factor history. Note that the treatment process is statistically exogenous just in the case that  $\mathcal{W}(t) = 1$  for all  $t$ . Of course,  $\mathcal{W}(t)$  is unknown and will have to be estimated from the data but, for pedagogic purposes, assume for the moment that it were known.

When  $A(t)$  is a statistically endogenous process, we shall consider estimating  $\gamma$  by a weighted logistic regression in which a subject is given the weight  $\mathcal{W}^{-1} \equiv [\mathcal{W}(K)]^{-1}$ . The weighted logistic regression estimator maximizes  $\prod_{i=1}^n [Lik_i(\gamma)]^{\mathcal{W}_i^{-1}}$ . This weighted logistic regression would agree with the usual unweighted analysis described above just in the case in which  $A(t)$  were exogenous. The somewhat surprising result described in detail below is that, if the vector of prognostic factors recorded in  $L(t)$  constitutes all relevant time-dependent prognostic factors (i.e., confounders), then, whether or not the treatment process is statistically exogenous, the weighted logistic regression estimator of  $\gamma_2$  will converge to a quantity  $\beta_2$  that can be appropriately interpreted as the causal effect of treatment history on the mean of  $Y$ . In contrast, when  $A(t)$  is statistically endogenous, the usual logistic regression estimator will still converge to the parameter  $\gamma_2$ , but now  $\gamma_2$  will have no causal interpretation.

To prove such a claim, we need to give a formal mathematical meaning to the informal concept of the causal effect of treatment history on the mean of  $Y$ . To do so, we first introduce some notational conventions. We use capital letters to represent random variables and lower case letters to represent possible realizations (values) of random variables. For example,  $O_i$  is the random observed data for the  $i^{\text{th}}$  study subject and  $o$  is a possible realization (value) of  $O_i$ . Further, we assume that the random vector  $O_i$  for each subject is drawn independently from a distribution common to all subjects. Because the  $O_i$  have the same distribution, we often suppress the  $i$  subscript.

**Counterfactual Outcomes:** Now we introduce counterfactual or potential outcomes. For any fixed non-random treatment history  $\bar{a} = \{a(u); 0 \leq u \leq K+1\}$ , let  $Y_{\bar{a}}$  be the random variable representing a subject’s outcome had, possibly contrary to fact, the subject been treated with history  $\bar{a}$  rather than his observed history  $\bar{A}$ . Note the  $\bar{a}$ ’s are possible realizations of the random variable  $\bar{A}$ . For each possible history  $\bar{a}$ , we are assuming a

subject's response  $Y_{\bar{a}}$  is well defined (although generally unobserved). Indeed we only observe  $Y_{\bar{a}}$  for that treatment history  $\bar{a}$  equal to a subject's actual treatment history  $\bar{A}$ , i.e.,  $Y = Y_{\bar{A}}$ . Then formally our statement that the effect of treatment history on the mean of  $Y$  is a linear logistic function of cumulative treatment is the statement that, for each  $\bar{a}$ ,

$$E[Y_{\bar{a}}] = g(\bar{a}; \beta) \text{ where } g(\bar{a}; \beta) = [1 + \exp\{-\beta_1 - \beta_2 \text{ cum}(\bar{a})\}]^{-1},$$

which we refer to as a MSM for the effect of treatment on the mean of  $Y$ . The model for  $E[Y_{\bar{a}}]$  is a marginal structural model since it is a model for the marginal distribution of counterfactual variables and, in the econometric and social science literature, causal models (i.e., models for counterfactual variables) are often referred to as structural.

The parameter  $\beta_2$  of our MSM is of important policy interest. To see why, consider a new subject exchangeable with (i.e., drawn from the same distribution as) the  $n$  study subjects. We must decide which treatment history  $\bar{a}$  to administer to the new subject. We would like to provide the treatment that minimizes the subject's probability of having HIV RNA in his blood at end of follow-up. That is, we want to find  $\bar{a}$  that minimizes  $E[Y_{\bar{a}}]$ . Thus, for example, if the parameter  $\beta_2$  of our causal model is positive, we will withhold AZT treatment from our subject (i.e., we will give him the treatment history  $\bar{a} \equiv 0$ ), since positive  $\beta_2$  indicates that the probability of having HIV RNA in one's blood at the end of follow-up increases with increasing cumulative AZT dose. In contrast to  $\beta_2$ , the parameter  $\gamma_2$  of our association model  $E[Y | \bar{A}] = g(\bar{A}; \gamma)$  may have no causal interpretation. For example, suppose physicians preferentially started AZT on subjects who, as indicated by their prognostic factor history, were doing poorly and that AZT has no causal effect on the mean of  $Y$  (i.e.,  $\beta_2 = 0$ ). Nonetheless, the mean of  $Y$  will increase with cumulative AZT doses and thus  $\gamma_2$  will be positive. In this setting, we say that the parameter  $\gamma_2$  of the association model lacks a causal interpretation because it is confounded by the association of the prognostic factors  $\bar{L}(u)$  with the treatment  $A(u)$ .

Formally, in terms of counterfactuals, we say that the  $A(t)$  process is "causally exogenous" if, for all histories  $\bar{a}$ ,

$$Y_{\bar{a}} \perp\!\!\!\perp A(t) \mid \bar{A}(t-1)$$

which is equivalent to

$$Y_{\bar{a}} \perp\!\!\!\perp \bar{A}.$$

Given the covariates recorded in  $L(t)$ , we say there are no unmeasured confounders if for each  $\bar{a}$

$$Y_{\bar{a}} \perp\!\!\!\perp A(t) \mid \bar{L}(t), \bar{A}(t-1).$$

With these formalizations, it can then be shown mathematically, that when there are no unmeasured confounders, (i) statistical exogeneity [i.e.,  $\bar{L}(t) \perp\!\!\!\perp A(t) \mid \bar{A}(t-1)$ ] implies that the  $A(t)$  process is "causally exogenous," (ii) the weighted logistic estimator converges to the parameter  $\beta_2$  of the marginal structural model for  $E[Y_{\bar{a}}]$ , and (iii) the limit  $\gamma_2$  of the usual logistic estimator generally differs from the causal parameter  $\beta_2$  of the MSM unless the treatment process is statistically exogenous.

We shall also refer to the assumption of no unmeasured confounders as the assumption that treatment  $A(t)$  is sequentially ignorable or randomized given the past. The assumption states that, conditional on AZT history and the history of all recorded covariates prior to  $t$ , increments in AZT dosage rate at  $t$  are independent of the counterfactual random variables  $Y_{\bar{a}}$ . This assumption will be true if all prognostic factors for, i.e., predictors of,  $Y_{\bar{a}}$  that are used by patients and physicians to determine the dosage of

AZT at  $t$  are recorded in  $\bar{L}(t)$  and  $\bar{A}(t-1)$ . For example, since physicians tend to withhold AZT from subjects with low white blood count, and in untreated subjects, low white blood count is a predictor of HIV RNA, the assumption of no unmeasured confounders would be false if  $\bar{L}(t)$  does not contain WBC history. It is the primary goal of the epidemiologists conducting an observational study to collect data on a sufficient number of covariates to ensure that the assumption of no unmeasured confounders will be at least approximately true.

The assumption of no unmeasured confounders is the fundamental condition that will allow us to draw causal inferences from observational data. It is precisely because it cannot be guaranteed to hold in an observational study and is not empirically testable that it is so very hazardous to draw causal inferences from observational data. Note that if, as in a sequentially randomized trial, at each time  $t$ , the dose of AZT was chosen at random by the flip of a coin, then the assumption of no unmeasured confounders would be true even if the probability that the coin landed heads depended on past measured covariate and AZT-history. It is because physical randomization guarantees the assumption that most people accept that valid causal inferences can be obtained from a randomized trial. See Rubin (1978), Robins (1986) and Holland (1986) for further discussion. Robins (1997, 1998b) and Robins et al. (1999) discuss how the consequences of violations of the assumption of no unmeasured confounders can be explored through sensitivity analysis. Also see Appendix 3 below.

Given the assumption of no unmeasured confounders, Robins (1987) shows the mean of the dichotomous variable  $Y_{\bar{a}}$  is non-parametrically identified from the joint distribution  $F_O$  of the observed data  $O$  by the g-computation algorithm formula of Robins (1986). Specifically,  $E(Y_{\bar{a}}) = b(\bar{a})$  where

$$b(\bar{a}) \equiv \int \cdots \int E(Y | \bar{\ell}_K, \bar{a}_K) \prod_{k=0}^K f(\ell_k | \bar{\ell}_{k-1}, \bar{a}_{k-1}) d\mu(\ell_k) \quad (*)$$

and for notational convenience we have written  $\bar{z}(k)$  as  $\bar{z}_k$  and  $z(k)$  as  $z_k$ .

The g-computation algorithm functional  $b(\bar{a})$  is the marginal mean of  $Y$  in the manipulated subgraph of the directed acyclic graph (DAG)  $G$  representing the observed data  $O$  in which all arrows into the treatment variables  $\bar{A} = (A_1, \dots, A_K)$  have been removed and  $\bar{A}$  is set to  $\bar{a}$  with probability 1 (Spirtes et al., 1993). More specifically, let DAG  $G$  be the complete DAG with temporally ordered vertex set  $O = \{L_0, A_0, L_1, A_1, \dots, A_K, Y\}$  and let DAG  $G_{\bar{a}}$  be the subgraph of  $G$  in which all arrows into the  $A_k, k = 0, \dots, K$  have been cut. Then  $b(\bar{a})$  is the marginal mean of  $Y$  based on a distribution for  $O$  represented by DAG  $G_{\bar{a}}$  in which  $f(A_k | \bar{A}_{k-1}, \bar{L}_k)$  is replaced by a degenerate density that takes the value  $a_k$  with probability 1 while the conditional density of each other variable in the set  $O$  given its parents remains as in  $F_O$ .

We say that the distribution of  $O = \{L_0, A_0, L_1, A_1, \dots, A_K, Y\}$  is standardly parameterized if, for each variable in  $O$ , we have specified a parametric or semiparametric model for the conditional distribution of that variable given its temporal predecessors (the past) and the parameters of each conditional model are variation-independent of those of any other conditional model. When our goal is to estimate the effect of a sequential (time-dependent) treatment  $\bar{A}$  on an outcome  $Y$ , Lemma 1 and Theorem 2 of Robins and Wasserman (1997) imply that inference procedures based on the standard parameterization will fail. Specifically, they prove that common choices for the parametric families in a standard parameterization often lead to joint densities such that the g-computation formula for  $E(Y_{\bar{a}})$  can never satisfy the causal null hypothesis that  $E(Y_{\bar{a}})$  is the same for all  $\bar{a}$ . In particular, the causal null hypothesis does not imply that  $Y \perp\!\!\!\perp \bar{A}_K | \bar{L}_K$ . As a consequence, in large samples, the causal null hypothesis, even when true, will be falsely rejected regardless of the data. Robins and Wasserman propose reparameterizing the distribution of  $O$  using structural nested models. MSMs represent an alternative reparameterization that

also overcomes the fatal deficiencies of the standard parameterization.

**Theory of Inverse-Probability-of-Treatment-Weighting:** We now explain why weighting by  $\mathcal{W}^{-1}$  corrects our logistic regression estimator for the “confounding” due to the prognostic factors in  $L(t)$ . The first point to note is that in the definition of  $\mathcal{W}(t)$  we could have replaced the denominator  $pr [A(t) | \bar{A}(t-1)]$  by any other function of  $\bar{A}(t)$  without influencing the consistency of our weighted logistic estimator of the parameter  $\beta_2$  of the MSM; only the efficiency (variance) of our estimator would be influenced. However, our estimator would be inconsistent if we replaced the numerator by any other function of  $\bar{A}(t)$  and  $\bar{L}(t)$ . Thus one can view weighting by  $\mathcal{W}^{-1}$  as weighting by the inverse of a subject’s probability of having his own observed treatment history. Now view each person as a member of a pseudo- or ghost population consisting of themselves and  $\mathcal{W}^{-1} - 1$  ghosts (copies) of themselves who have been added by weighting. In this new ghost or pseudo population, it is easy to show that  $\bar{L}(t)$  does not predict treatment at  $t$  given past treatment history, and thus we have created a pseudo-population in which treatment is exogenous. Furthermore, the causal effect of  $\bar{A}$  on  $Y$  in the ghost population is the same as in the original population. That is, if  $E[Y_{\bar{a}}] = g(\bar{a}; \beta)$  in the true population, the same will be true of the ghost population. Hence, we would like to do ordinary logistic regression in the pseudo-population. That is essentially what our weighted logistic regression estimator is doing, since the weights create, as required,  $\mathcal{W}^{-1} - 1$  additional copies of each subject.

A formal, mathematical explanation of why weighting by  $\mathcal{W}^{-1}$  corrects our logistic regression estimator for “confounding” is given in the following lemma characterizing the g-computation algorithm functional  $b(\bar{a})$  defined in (\*) above.

**Lemma 1.1:**  *$b(\bar{a})$  defined in (\*) is the unique function  $c(\bar{a})$  of  $\bar{a}$  such that  $E[q(\bar{A})(Y - c(\bar{A})) / \mathcal{W}] = 0$  for all functions  $q(\bar{A})$  for which the expectation exists.*

Lemma 1.1 has the following corollary.

**Lemma 1.2:** *Under sequential randomization,  $E(Y_{\bar{a}})$  is unique function  $c(\bar{a})$  of  $\bar{a}$  such that  $E[q(\bar{A})(Y - c(\bar{A})) / \mathcal{W}] = 0$  for all functions  $q(\bar{A})$  where the expectation exists.*

Consistency of our weighted estimator then follows from the fact that the probability limit of our weighted score equation is  $E[q(\bar{A})(Y - c(\bar{A})) / \mathcal{W}] = 0$  with  $q(\bar{A}) = (1, cum(\bar{A}))'$  and  $c(\bar{A}) = g(\bar{A}, \beta)$ .

Under a mild strengthening of our assumption of sequential randomization (no unmeasured confounders), a simple, quite revealing, purely “causal” proof of Lemma 1.2 can be obtained that does not use the fact that  $E(Y_{\bar{a}})$  is given by the g-computation algorithm formula  $b(\bar{a})$  of Eq. (\*). Let  $Y_{\bar{\mathcal{A}}} = \{Y_{\bar{a}}; \bar{a} \in \bar{\mathcal{A}}\}$  where  $\bar{\mathcal{A}}$  is the support of the random variable  $\bar{A}$ . Suppose we strengthen our assumption of no unmeasured confounders to

$$Y_{\bar{\mathcal{A}}} \amalg A_k | \bar{L}_k, \bar{A}_{k-1} .$$

Denote the factual and counterfactual data by  $Z = (Y_{\bar{\mathcal{A}}}, \bar{A}, \bar{L})$  and the observed data by  $O = (Y \equiv Y_{\bar{\mathcal{A}}}, \bar{A}, \bar{L})$ . We can factor the true joint density of  $Z$  that generated the data as

$$f(Z) = f(Y_{\bar{\mathcal{A}}}) \prod_{k=0}^K f(L_k | \bar{L}_{k-1}, \bar{A}_{k-1}, Y_{\bar{\mathcal{A}}}) \prod_{k=0}^K f(A_k | \bar{L}_k, \bar{A}_{k-1}) .$$

Now let  $f^*(A_k | \bar{A}_{k-1})$  be a density for  $A_k$  given  $\bar{A}_{k-1}$ . It need not equal the true density  $f(A_k | \bar{A}_{k-1})$ .

Let  $f^*(Z)$  be a joint density for  $Z$  that differs from the true joint density  $f(Z)$  only in that  $f^*(A_k | \bar{L}_k, \bar{A}_{k-1}) = f^*(A_k | \bar{A}_{k-1})$  so that  $A_k$  is strictly exogenous. Thus,

$$f^*(Z) = f(Y_{\bar{\mathcal{A}}}) \prod_{k=0}^K f(L_k | \bar{L}_{k-1}, \bar{A}_{k-1}, Y_{\bar{\mathcal{A}}}) \prod_{k=0}^K f^*(A_k | \bar{A}_{k-1}) .$$

Now  $E(Y_{\bar{a}}) = E^*(Y_{\bar{a}})$  since  $f(Z)$  and  $f^*(Z)$  have the same marginal law for  $Y_{\bar{a}}$ . Second, since  $\bar{A}$  is causally exogenous under  $f^*(z)$  [i.e.,  $Y_{\bar{a}} \perp\!\!\!\perp \bar{A}$ ], we have that  $E^*[Y_{\bar{a}}] = E^*[Y_{\bar{a}} | \bar{A} = \bar{a}] = E^*[Y_{\bar{A}} | \bar{A} = \bar{a}] = E^*[Y | \bar{A} = \bar{a}]$ . That is, by  $\bar{A}$  causally exogenous, the mean of  $Y_{\bar{a}}$  is given by the regression function  $E^*[Y | \bar{A} = \bar{a}]$  of  $Y$  on  $\bar{A} = \bar{a}$ . Now it is a standard result that the regression function  $E^*(Y | \bar{A} = \bar{a})$  is characterized as the unique function  $c(\bar{a})$  solving  $E^*\{q(\bar{A})[y - c(\bar{A})]\} \equiv \int q(\bar{A})(y - c(\bar{A}))f^*(Z)d\mu(Z) = 0$  for all  $q(\bar{A})$  where  $\mu$  is a dominating measure. But,  $\int q(A)(Y - c(\bar{A}))f^*(Z)d\mu(Z) = \int q(A)(Y - c(\bar{A}))\frac{f^*(Z)}{f(Z)}f(Z)d\mu(Z) = E\left[q(A)(Y - c(\bar{A}))\frac{f^*(Z)}{f(Z)}\right]$ . But, by definition,  $\frac{f^*(Z)}{f(Z)} = \mathcal{W}^{-1}$  when  $f^*(A_k | \bar{A}_{k-1}) = f(A_k | \bar{A}_{k-1})$ . Lemma 1.2 then follows, since  $E^*(Y | A = \bar{a}) = E(Y_{\bar{a}})$ . The proof also makes clear that consistency of our weighted estimator does not require that we choose  $f^*(A_k | \bar{A}_{k-1}) = f(A_k | \bar{A}_{k-1})$ .

## Data Analyses: Marginal Structural Mean Model for a Repeated Measures Outcome:

To give a better picture of the meaning and use of MSMs, we report preliminary results of two data analyses. Full details will be published elsewhere. We estimate the joint effect in ACTG Randomized Trial 002 of AZT treatment arm and aerosolized pentamidine (AP) on the evolution of CD4 count in the first analysis and on mortality in the second analysis. This trial was designed to compare the effect of high-dose AZT with low-dose AZT on survival. However, over fifty percent of the subjects failed to comply with the assigned treatment protocol and initiated treatment with a non-randomized therapy, AP, during the course of the trial. The joint effects of AP and AZT treatment arm on survival have been previously estimated using structural nested failure time models by Robins and Greenland (1994). We first consider a MSM model for the effect of AP therapy on the mean of the log transformed CD4 count history while adjusting for baseline variables. CD4 count measurements were obtained at weeks 8, 16, 24, and 32 measured in days. Specifically, we consider the MSM

$$E[Y_{\bar{a}}(m) | V^\dagger] = g_m[\bar{a}(m), V^\dagger, \beta]$$

where  $Y(m) = \log[CD4(m) + 2]$ ,  $CD4(m)$  is the CD4 count on day  $m$ ,  $Y_{\bar{a}}(m)$  is the counterfactual version of  $Y(m)$  under the AP history  $\bar{a}$ ,  $V^\dagger = (1, m, R, Y(0), \log WBC(0))'$  is the vector of baseline regressors with  $R = 1$  denoting the high AZT treatment arm and  $R = 0$  the low AZT treatment arm,  $Y(0)$  is defined above, and  $WBC(0)$  is baseline white blood count. We modeled the regression function as  $g_m[\bar{a}(m), V^\dagger, \beta] = \beta_1'V^\dagger + \beta_2 cum(m, \bar{a})$  with  $cum(m, \bar{a}) = \int_0^m a(t) dt$  being cumulative AP treatment up to day  $m$ . We make the assumption of no unmeasured confounders with  $L(t)$  being white blood count, the number of episodes of pneumocystis pneumonia (PCP), an AIDS-related pneumonia, up to day  $t$ , and CD4 count on day  $t$ . Furthermore, the baseline covariates  $V^\dagger$  are included in  $L(0)$ . Arguing as above, it can be shown a consistent estimator of  $\beta$  is obtained by fitting the model  $E[Y(m) | V^\dagger, \bar{A}] = g_m[\bar{A}(m), V^\dagger, \beta]$  using the generalized estimating equation (GEE) option of Proc genmod in the SAS software package under a working independence matrix and weighting the observation  $Y(m)$  for a subject by  $\{\mathcal{W}(m)\}^{-1}$ . The GEE option of Proc genmod is simply a program that fits the above model by weighted least squares to obtain an estimate of  $\beta$ . In estimating  $\beta$ , the program treats each individual at each of the four times  $m$  as four separate observations when computing the least squares estimator. However, the program outputs a robust variance estimator that appropriately accounts for the fact that the four observations on a given subject are correlated.

Since  $\mathcal{W}(m)$  was unknown, it was estimated from the data by fitting logistic models for  $pr[A(k) | \bar{L}(k), \bar{A}(k-1)]$  and  $pr[A(k) | \bar{A}(k-1), V^\dagger]$ . [Note that when our MSM conditions on baseline variables  $V^\dagger$ , they should be included in the denominator of  $\mathcal{W}(m)$ .] Specifically, we fit the model

$$\text{logit } pr[A(k) = 0 | \bar{L}(k), \bar{A}(k-1) \equiv 0] = \alpha'Q(k)$$

where  $Q(k) = (1, \log k, \log [WBC(k-1)], Y(k-1), PCP\ bouts(k-1), Y(0), \log [WBC(0)], R)$ . Here,  $PCP\ bouts(k-1)$  is the number of episodes (bouts) of pneumocystis pneumonia through  $k-1$ . In fitting the model, we treated each subject at each day  $k, k = 0, 1, \dots, 224$  as an independent observation (which is justified by the conditional martingale structure of the model). We note that since in the 002 data file, any subject starting AP remained on it thereafter, it was only necessary to fit a model for  $A(k)$  for subjects who had yet to begin AP [i.e.,  $\bar{A}(k-1) \equiv 0$ ]. We estimated  $pr[A(k) | \bar{A}(k-1) \equiv 0, V^\dagger]$  by fitting the above model after eliminating the random time-dependent terms that were functions of  $(k-1)$ . The 95 percent Wald intervals computed using the robust variance outputted by the GEE program are conservative (i.e., they are guaranteed to cover the true  $\beta$  at least 95 percent of the time in large samples) because they do not account for estimation of the weights  $\mathcal{W}(m)$ . It is interesting that estimating the weights shrinks the variance of our estimator of  $\beta$ , so that our intervals (which do not account for the fact that the weights are estimated) are conservative. Note that the elements of the vectors  $\alpha$  multiplying the time-dependent covariates  $\log [WBC(k-1)], Y(k-1)$  and  $PCP\ bouts(k-1)$  will all be zero if and only if the AP treatment process is statistically ancillary. A three degree of freedom likelihood ratio test of the hypothesis that all three components of  $\alpha$  were zero rejected at the  $p < .01$  level. As a consequence, we rejected the hypothesis of statistical exogeneity.

The analysis just described assumes that there is no drop-out or censoring by end of follow-up. To correct for this, we defined a subject as censored (i.e., permanently missing) the first time he missed one of his scheduled visits or was censored by end of follow-up. Under the assumption of ignorable drop-out given the time-dependent factors  $L(t)$  and treatment  $A(t)$ , we still obtain consistent estimators of  $\beta$  in the presence of drop-out if we weight a subject uncensored at day  $m$  by  $\{\mathcal{W}(m)\mathcal{W}^\dagger(m)\}^{-1}$  where  $\mathcal{W}^\dagger(m) = \prod_{k=0}^m \{p[R(k)=0 | \bar{R}(k-1)=0, \bar{L}(k-1), \bar{A}(k-1)] / pr[R(k)=0 | \bar{R}(k-1)=0, \bar{A}(k-1), V^\dagger]\}$  is the ratio of a subject's probability of remaining uncensored up to day  $m$  divided by that probability calculated as if there had been no time-dependent determinants of drop-out except past treatment history. Here,  $R(m) = 0$  if a subject has not dropped out or reached end to follow-up by day  $m$ . Since  $\mathcal{W}^\dagger(m)$  is unknown, it was estimated from the data in a manner completely analogous to the estimation of  $\mathcal{W}(m)$  except with  $A(k)$  replaced by  $R(k)$  as the outcome variable and with  $A(k-1)$  added as an additional regressor. Furthermore, in the presence of censoring, it is necessary when estimating  $pr[A(k)=0 | \cdot]$  to add the event  $R(k) = 0$  to the conditioning event.

We fit the above models and obtained an estimate  $\hat{\beta}_2 = .001$  and conservative 95 percent confidence interval  $(-.026, .028)$  for the parameter  $\beta_2$  representing the causal effect of cumulative AP dose on CD4 count. Furthermore, since in trial 002 the assignment to AZT treatment arm was at random with probability 1/2, the AZT treatment arm indicator is exogenous. It follows that the component  $\beta_{1R}$  of  $\beta_1$  multiplying the AZT treatment indicator  $R$  has the interpretation as the direct effect of AZT treatment arm on the evolution of mean CD4 count that is not through AP history. We obtained an estimate of .0196 with a conservative 95 percent confidence interval of  $(-.123, .163)$  for  $\beta_{1R}$ .

**Marginal Structural Cox Proportional Hazards Model:** We next estimated the joint effects of AP therapy and AZT treatment arm on survival by specifying a marginal structural Cox proportional hazards model

$$\lambda_{T_{\bar{a}}}(t | V^\dagger) = \lambda_0(t) \exp[\beta'_1 V^\dagger + \beta_2 a(t)]$$

where  $T_{\bar{a}}$  is the subject's time to death if he had followed AP history  $\bar{a}$ ,  $\lambda_{T_{\bar{a}}}(t | V^\dagger)$  is the hazard of  $T_{\bar{a}}$  at  $t$  given  $V^\dagger$ ,  $\lambda_0(t)$  is an unspecified baseline hazard function, and  $V^\dagger = (R, Y(0), \log WBC(0))$ . Note this model specifies that the hazard of failure at time  $t$  depends on current AP status rather than on



cumulative AP history. Arguing as in the previous subsection, we can obtain consistent estimates of the unknown parameter  $\beta = (\beta_1', \beta_2)'$  by fitting the ordinary time-dependent Cox model  $\lambda_T(t | \bar{A}(t), V^\dagger) = \lambda_0(t) \exp[\beta_1' V^\dagger + \beta_2 A(t)]$  except that the contribution of subject to a calculation performed on subjects at risk at time  $t$  is weighted by  $\widehat{W}(t)^{-1} \widehat{W}^\dagger(t)^{-1}$ . Note that now when we model  $pr[R_k = 0 | \cdot]$  and  $pr[A_k = 0 | \cdot]$  we must include the event  $T > k$  among the conditioning events. Here we have adopted the convention that on any day  $k$  censoring occurs at the end of the day. Note since the subject-specific weights change with time, one either needs to write a special program or trick a standard time-dependent Cox model that allows weights into allowing for time-varying weights by a clever use of the time-varying stratum option available in many off-the-shelf Cox programs. To obtain conservative 95 percent confidence intervals for  $\beta$ , one needs to compute the so-called robust variance of Lin et al. (1989). Implementing the above procedure, we obtained an estimate  $\widehat{\beta}_2 = -.1362$  with 95 percent conservative confidence interval  $(-.35, .09)$  for  $\beta_2$ . For the component  $\beta_{1R}$  of  $\beta_1$  representing the direct effect of AZT treatment arm on survival, we obtained an estimate of .1890 with a 95 percent confidence interval of  $(-.01, .21)$  indicating borderline statistically significant evidence for a beneficial effect of the low-dose AZT arm. Both the results obtained for the prophylaxis effect and for the AZT effect were consistent with those obtained by Robins and Greenland (1994) using SNFTMs.

**Philosophical Interlude:** We pause to comment briefly on the definition and nature of the counterfactual random variables  $T_{\bar{a}}$ . Following Lewis (1973), we consider  $T_{\bar{a}}$  to be a subject's death time in the closest possible world to this in which, possibly contrary to fact, the subject was treated with the AP history  $\bar{a}$ . Consider a subject  $i$  who, in the 002 trial, was assigned to the high-dose AZT arm, received AP from week 10 to 40, took the assigned 1500 mg. of AZT daily until week 12 but then stopped all further AZT therapy, and finally died in week 40. If AP had been withheld, it is quite conceivable that subject  $i$  would have continued to be assigned 1500 mg. of AZT daily past week 12 if either (1) AP potentiated the toxic effects of AZT, precipitating a life-threatening toxic episode in week 12, or (2) the subject, although not toxic, had stopped AZT at week 12 because he felt himself to be adequately protected by the AP treatment. To be concrete, say, in the closest possible world, subject  $i$  would have continued to take 1500 mg. of AZT daily through week 14 and none thereafter if AP had been withheld. We now consider the meaning of the counterfactual  $T_0 \equiv T_{\bar{a}=0}$  in which AP was always withheld. Then, by its definition,  $T_0$  would be equal to subject  $i$ 's failure time when the subject was assigned to the high-dose arm, never received AP, and took AZT daily through week 14 (rather than through week 12). Thus,  $T_0$  might well differ from the counterfactual variable, say  $T_0^*$ , representing a subject's survival time in the closest possible world in which AP was withheld but, as in this world, AZT was stopped after week 12. Several comments are in order.

First,  $T_{0i}$  is conceptually rather well-defined, even if we do not observe what the particular subject  $i$  would have done about his AZT dose after week 10 in the absence of AP therapy, as  $T_{0i}$  is just subject  $i$ 's outcome in the closest possible world to this one in which all AP therapy is withheld and all consequences which flow from that (including possibly taking AZT in week 12-14) are all allowed to occur. Second, it may be quite reasonable to make (at least to a good approximation) the assumption of no unmeasured confounders for  $T_0$ , in which case its distribution is non-parametrically identified by inverse-probability-of-treatment-weighting, (equivalently, by the g-computation algorithm formula) from the distribution of the observed data. The intuitive reason for this successful identification is that for a subset of the population (i.e., those who never did take AP), we do observe  $T_0$ , and under the assumption of no unmeasured confounders, we can appropriately reweight them by  $\mathcal{W}^{-1}$  to construct a ghost population whose distribution of  $T_0 = T$  is the same as that of  $T_0$  in the true study population. Third, from a public

health point of view, it is much more important to identify the distribution of  $T_0$  than  $T_0^*$  since it is the distribution of  $T_0$  that would result if we made the public policy decision to withhold AP therapy. Fourth,  $T_0^*$  may be more relevant than  $T_0$  in a legal case against the manufacturers of AP. For example, the manufacturers would argue that they should not be held responsible for any damages if a subject’s observed death time  $T$  equalled their counterfactual death time  $T_0^*$  (even if  $T$  differed from  $T_0$  due to differences in the amount of AZT taken). Fifth, the distribution of  $T_0^*$  is not identified even in an experiment in which both AP and AZT are randomly assigned. Thus, no amount of data evidence will ever determine the distribution of  $T_0^*$ , even in a randomized experiment (Robins and Greenland, 1989).

**Comparison with SNMs:** Marginal structural models are an alternative to structural nested models. A SNM is a model for the magnitude of the causal effect of a final brief blip of a time-dependent treatment at time  $t$  as a function of past time-dependent treatment and prognostic factor history. The causal parameter of a structural nested model is identified under the assumption of no unmeasured confounders. The essential difference between MSMs and SNMs is that SNMs model the magnitude of the effect of a treatment given at  $t$  as a function of the prognostic factor history up to  $t$ . In contrast, MSMs model the causal effect of treatment given at  $t$  only as a function of baseline prognostic factors. Sec. 5 below is devoted to describing what is known about the advantages and disadvantages of MSMs versus SNMs. Some of the advantages of MSMs were discussed above. Possible disadvantages include the following. (i) Inability to easily estimate the effects of dynamic treatment regimes (i.e., treatment plans where a subject’s covariate history up to time  $k$  determines the treatment to be taken at  $k$ ). Actual medical treatments are usually dynamic, since if a subject becomes toxic to a drug, the drug must be stopped. (ii) The inability to directly test the null hypothesis of no effect of any treatment regime (dynamic or non-dynamic) on outcome. (iii) The difficulty in performing likelihood-based inference for MSMs, since the likelihood is a computational nightmare. (iv) Lack of identifiability of the MSM model parameters when sequential ignorability holds for a so-called “instrumental variable” but not for the actual treatment of interest. (v) MSMs, in contrast to SNMs, cannot be used if there exists a value of  $\ell_k$ , say  $\ell_k = 0$ , such that for all but one value of  $a_k$ ,  $f[a_k | \bar{\ell}_{k-1}, \ell_k = 0, \bar{a}_{k-1}] = 0$ . An example would be a study of the effect of an occupational exposure on mortality with  $\ell_k = 0$  if a subject is off work at time  $k$ ,  $\ell_k = 1$  otherwise, and subjects off work can only receive exposure level  $a_k = 0$ . We note that SNMs do not suffer from any of these five deficiencies.

## 2 Advantages of MSMs with Continuous $Y$ or with Failure Time: A Formal Definition of MSMs

### 2.1 The Data

Consider a study where we observe  $n$  i.i.d. copies of data  $O = (\bar{A}(C), \bar{L}(C))$ , where  $C$  is an administrative end of follow-up time,  $\bar{A}(C)$  is a treatment process,  $\bar{L}(C)$  is an outcome or response process and, for any  $Z(u)$ ,  $\bar{Z}(t) \equiv \{Z(u); 0 \leq u \leq t\}$ . We assume  $C$  is an element of  $L(0)$  since it is assumed known at time 0.

For purposes of causal inference, we assume the existence of an underlying treatment process  $\bar{A} = \{A(u); 0 \leq u < \infty\}$  with  $A(u)$  taking values in a set  $\mathcal{A}(u)$  and the existence of underlying counterfactual random variables

$$\{\bar{L}_{\bar{a}}; \bar{a} \in \bar{\mathcal{A}}\} \tag{1}$$

where  $\bar{L}_{\bar{a}} = \{L_{\bar{a}}(u); 0 \leq u < \infty\}$ ,  $\bar{a} = a(\cdot) = \{a(t); 0 \leq t < \infty \text{ and } a(t) \in \mathcal{A}(t)\}$  is a treatment plan (equivalently, regime or function) lying in a set of functions  $\bar{\mathcal{A}}$ . Given a regime  $\bar{a}$ , let  $\bar{L}_{\bar{a}(u),0}$  be counterfactual history under a regime  $\bar{a}^*$  that agrees with  $\bar{a}$  through time  $u$  and is 0 thereafter, where 0 is the baseline value of  $a(t)$ . Then we assume that the  $\bar{L}_{\bar{a}}$  satisfy the following consistency assumption with probability 1:

$$\bar{L}_{\bar{a}(u),0}(u) = \bar{L}_{\bar{a}(t),0}(u) = \bar{L}_{\bar{a}}(u) = \bar{L}_{\bar{a}^\dagger}(u) \quad (2)$$

for all  $t > u$  and all  $\bar{a}^\dagger$  with  $\bar{a}^\dagger(u) = \bar{a}(u)$ . This assumption essentially says that the future does not determine the past. The observed data are linked to the counterfactual data by

$$\bar{L}(C) = \bar{L}_{\bar{A}(C),0}(C) . \quad (3)$$

Eq. (3) states that a subject's observed outcome history through end of follow-up is equal to their counterfactual outcome history corresponding to the treatment they did indeed receive. We assume  $\bar{L}_{\bar{a}} = (\bar{Y}_{\bar{a}}, \bar{V}_{\bar{a}})$  where  $\bar{Y}_{\bar{a}}$  is an outcome process of interest and  $\bar{V}_{\bar{a}}$  is the process of other recorded variables. Further, we shall make the sequential randomization (i.e., ignorable treatment assignment) assumption that for all  $t$  and  $\bar{a} \in \bar{\mathcal{A}}$ ,

$$\underline{Y}_{\bar{a}}(t) \amalg \amalg A(t) \mid \bar{L}(t^-), \bar{A}(t^-) \quad (4)$$

where for any variable  $\underline{Z}(t) = \{Z(u); u \geq t\}$  is the history of that variable from  $t$  onwards. We also refer to (4) as the assumption of no unmeasured confounders given prognostic factors  $L(t)$ . Because of measurability issues, (4) is not well-defined. If the  $A(t)$  process can only jump at discrete non-random times  $t_1, t_2, \dots$  and the  $\bar{L}(t)$  process has left-hand limits, i.e.,  $\bar{L}(t^-) \equiv \lim_{u \uparrow t} \bar{L}(u)$ , (4) is formally, for each  $t_k$ ,

$$f[A(t_k) \mid \bar{L}(t_k^-), \bar{A}(t_k^-), \underline{Y}_{\bar{a}}(t_k)] = f[A(t_k) \mid \bar{L}(t_k^-), \bar{A}(t_k^-)] . \quad (5)$$

where  $f(\cdot \mid \cdot)$  is the conditional density of  $A(t_k)$  with respect to a dominating measure. If  $A(t)$  is a marked point process that can jump in continuous time with CADLAG (continuous from the right with left-hand limits) step-function sample paths, then Eq. (4) is formally that

$$\lambda_A[t \mid \bar{L}(t^-), \bar{A}(t^-), \underline{Y}_{\bar{a}}(t)] = \lambda_A[t \mid \bar{L}(t^-), \bar{A}(t^-)] \quad (6a)$$

and

$$\begin{aligned} f[A(t) \mid \bar{L}(t^-), \bar{A}(t^-), A(t) \neq A(t^-), \underline{Y}_{\bar{a}}(t)] = \\ f[A(t) \mid \bar{L}(t^-), \bar{A}(t^-), A(t) \neq A(t^-)] . \end{aligned} \quad (6b)$$

Here, the intensity process  $\lambda_A(t \mid \cdot)$  is  $\lim_{\delta \rightarrow 0} pr[A(t + \delta t) \neq A(t^-) \mid A(t^-), \cdot] / \delta t$ . Eq. (6a) says that given past treatment and confounder history, the probability that the  $A$  process jumps at  $t$  does not depend on the future counterfactual history of the outcome of interest. Eq. (6b) says that given that the covariate process did jump at  $t$ , the probability it jumped to a particular value of  $A(t)$  does not depend on the future counterfactual history of the outcome of interest.

Following Heitjan and Rubin (1991), we say the data are coarsened at random (CAR) if

$$f[\bar{A}(C) \mid \{\bar{L}_{\bar{a}}; \bar{a} \in \bar{\mathcal{A}}\}] \text{ depends only on } O = (\bar{A}(C), \bar{L}(C)) . \quad (7)$$

Note that we can use ideas from the ‘‘missing data’’ literature because one's treatment history  $\bar{A}(C)$  determines which components of one's counterfactual history  $\{\bar{L}_{\bar{a}}; \bar{a} \in \bar{\mathcal{A}}\}$  one observes. Thus we can

view causal inference as a missing data problem (Rubin, 1976). If, as in all the models we shall consider in this paper, for each  $\bar{\mathcal{A}}^\dagger \subseteq \bar{\mathcal{A}}$  satisfying  $\bar{a}_1(u) \neq \bar{a}_2(u)$  for all  $\bar{a}_1, \bar{a}_2 \in \bar{\mathcal{A}}^\dagger$ , the  $\{L_{\bar{a}}(u); \bar{a} \in \bar{\mathcal{A}}^\dagger\}$  may have a non-degenerate joint distribution, then CAR implies sequential randomization (4) but the converse is not true (Robins et al., 1999). Robins (1997, pg. 83) gives examples where one would expect (4) to be true even when (7) is false. In this paper, we shall only need (4). However, even if (7) is also imposed this, by itself, essentially places no restrictions on the joint distribution of the observable random variables (Gill, van der Laan, Robins, 1997) and, thus, is not subject to empirical test.

## 2.2 MSMs

A MSM for  $\{\bar{Y}_{\bar{a}}; \bar{a} \in \bar{\mathcal{A}}\}$  places restrictions on the marginal distribution of the  $\bar{Y}_{\bar{a}}$  possibly conditional on a baseline variable  $V^\dagger$  in  $V(0)$  (with  $C \in V^\dagger$  if  $C$  is random). Examples of MSMs follow. Each of these examples will be important in our comparison of MSMs with structural nested models below.

**Model 1:** Suppose  $C = K + 1$  w.p.1., the  $\bar{A}(C)$  process jumps only at times  $0, 1, 2, \dots, K$  and the  $\bar{L}_{\bar{a}}$  process jumps only at times  $0^-, 1^-, 2^-, \dots, K^-, K + 1^-$ . In models 1a-1c, we are only concerned with an outcome measured at end of follow-up. Hence, we set  $Y_{\bar{a}}(m) \equiv 0$  with probability 1 for  $m \leq K$  and define  $Y_{\bar{a}} = Y_{\bar{a}}(K + 1)$ . Then we have

**Model 1a – non-linear least squares:**  $E[Y_{\bar{a}} | V^\dagger] = g[\bar{a}(K), V^\dagger, \beta_0]$  where  $g(\cdot, \cdot, \cdot)$  is a known function. This is the logistic regression MSM we discussed in the Introduction.

**Model 1b – semiparametric regression:**  $\eta\{E[Y_{\bar{a}} | V^\dagger]\} = g[\bar{a}(K), V^\dagger, \beta_0] + g^\dagger(V^\dagger)$  where  $\eta(\cdot)$  is a known monotone link function,  $g^\dagger(\cdot)$  is unknown and unrestricted and  $g(\cdot, \cdot, \cdot)$  is a known function satisfying  $g(\mathbf{0}, V^\dagger, \beta) = 0$ . The requirement that  $g(\mathbf{0}, V^\dagger, \beta) = 0$  implies that  $g^\dagger(V^\dagger)$  is the “main effect of  $V^\dagger$ .” Such models are also referred to as partial spline models. They are semiparametric because the main effect of  $V^\dagger$  is modelled non-parametrically.

**Model 1c – stratified transformation model:**  $pr[R(\bar{a}, \beta_0) < t | V^\dagger] = F_0(t | V^\dagger), F_0(t | V^\dagger)$  an unknown distribution function,  $R(\bar{a}, \beta) = r(Y_{\bar{a}}, \bar{a}, V^\dagger, \beta)$  is a known increasing function of  $Y_{\bar{a}}$  satisfying  $r(y, \bar{a}, V^\dagger, \beta) = y$  if  $\bar{a} \equiv \mathbf{0}$  or  $\beta = 0$ . This model says that we know the conditional quantile-quantile function linking the  $Y_{\bar{a}}$ 's given  $V^\dagger$  up to an unknown parameter  $\beta$ . It is the natural generalization of model 1b for mean functions to quantile-quantile functions.

In the following model, we are interested in the outcome at each  $m \geq 1$  so we no longer assume that  $Y_{\bar{a}}(m) \equiv 0$  with probability 1.

**Model 1d – multivariate non-linear least squares:**  $E[Y_{\bar{a}}(m) | V^\dagger] = g_m[\bar{a}(m - 1), V^\dagger, \beta_0]$ ,  $m = 1, \dots, K + 1$  where the  $g_m(\cdot, \cdot, \cdot)$  are known. This is the natural MSM version of longitudinal generalized estimating equation models for marginal means (Liang and Zeger, 1986). It is the model we use to analyze the 002 CD4 count data in the Introduction.

**Model 2:**  $C = \infty$ ,  $Y_{\bar{a}}$  is a failure time process, i.e.,  $Y_{\bar{a}}$  jumps from 0 to 1 at some particular time and stays at 1. Then define the failure time  $T_{\bar{a}}$  by the equation  $Y_{\bar{a}}(T_{\bar{a}}) = 1$  and  $Y_{\bar{a}}(T_{\bar{a}}^-) = 0$ . Let  $\lambda_0(t)$  and  $\lambda_0(t | V^\dagger)$  be unknown non-negative functions of  $t$  and  $(t, V^\dagger)$  respectively and, for any  $Z$ ,  $\lambda_Z(u)$  is the hazard of  $Z$ .

**Model 2a – Cox proportional hazards model:**  $\lambda_{T_{\bar{a}}}[t | V^\dagger] = \lambda_0(t) \exp[r(\bar{a}(t^-), t, V^\dagger; \beta_0)]$  where  $r(\cdot)$  is a known function satisfying  $r(\mathbf{0}, t, 0; \beta) = 0$ . This is the model we use to analyze the 002 mortality data in the Introduction.

**Model 2b – stratified Cox proportional hazards model:**  $\lambda_{T_{\bar{a}}}(t | V^\dagger) = \lambda_0(t | V^\dagger) \exp[r(\bar{a}(t^-), t, V^\dagger; \beta_0)]$  where, now,  $r(\mathbf{0}, t, V^\dagger; \beta) = 0$ .

**Model 2c – stratified time-dependent accelerated failure time model:**  $pr[r(T_{\bar{a}}, \bar{a}, V^\dagger, \beta_0) <$

$t | V^\dagger] =$

$F_0(t | V^\dagger)$  where  $r(u, \bar{a}, V^\dagger, \beta) = r(u, \bar{a}(u), V^\dagger, \beta)$  is a known function increasing in its first argument satisfying  $r(u, \mathbf{0}, V^\dagger, \beta) = u$ . This model can also be written as

$$\lambda_{R(\bar{a}, \beta_0)}(t | V^\dagger) = \lambda_0(t | V^\dagger)$$

for  $\bar{a} \in \bar{\mathcal{A}}$ , where  $R(\bar{a}, \beta) = r(T_{\bar{a}}, \bar{a}, V^\dagger, \beta)$ . This is the extension of model 1c to a failure time variable. It is the model studied by Robins and Tsiatis (1992).

### 3 Estimation

#### 3.1 Ancillary treatment process

In this section, we consider estimation of the parameter  $\beta_0$  of our marginal structural models. In this subsection, we will suppose that  $\bar{A}$  is a causally ancillary (i.e., exogenous) covariate process, i.e.,

$$\bar{A} \amalg \{\bar{L}_{\bar{a}}; \bar{a} \in \bar{\mathcal{A}}\} | V^\dagger. \quad (8)$$

The often unrealistic assumption (8) implies CAR but, in contrast to CAR, places restrictions on the joint distribution of the data. Specifically (8) implies statistical ancillarity

$$A(t) \amalg \bar{L}(t^-) | \bar{A}(t^-), V^\dagger \quad (9)$$

and thus (8) is subject to an empirical test.

Given (8), the restrictions on the observables  $O$  implied by any MSM are (9) and that the restrictions on the distribution of  $\bar{Y}_{\bar{a}}$  given  $V^\dagger$  specified by the MSM hold for the conditional distribution of the observable  $\bar{Y}(C)$  conditional on  $(\bar{A}(C), V^\dagger)$ .

For reasons that will become clear below, we indicate with a “\*” any expectations, probabilities or hazard functions computed under the assumption that (8) and (9) hold. For convenience, denote  $\bar{A}(C)$  as  $\bar{A}$ . Thus, for our MSM models (1a) – (2c), (8) implies the association models

**Model 1a:**  $E^*[Y | V^\dagger, \bar{A}] = g(\bar{A}, V^\dagger, \beta_0)$

**Model 1b:**  $\eta \{E^*[Y | V^\dagger, \bar{A}]\} = g(\bar{A}, V^\dagger, \beta_0) + g^\dagger(V^\dagger)$ .

**Model 1c:**  $R(\beta_0) \amalg^* \bar{A} | V^\dagger$  where  $R(\beta_0) \equiv R(\bar{A}, \beta_0)$ .

**Model 1d:**  $E^*[Y(m) | V^\dagger, \bar{A}] = g_m[\bar{A}(m-1), V^\dagger; \beta_0]$ ,  $m = 1, \dots, K+1$ .

**Model 2a:**  $\lambda_T^*[t | V^\dagger, \bar{A}] = \lambda_T^*[t | V^\dagger, \bar{A}(t^-)] = \lambda_0(t) \exp[r(\bar{A}(t^-), t, V^\dagger, \beta_0)]$ .

**Model 2b:**  $\lambda_T^*[t | V^\dagger, \bar{A}] = \lambda_T^*[t | V^\dagger, \bar{A}(t^-)] = \lambda_0(t | V^\dagger) \exp[r(\bar{A}(t^-), t, V^\dagger, \beta_0)]$ .

**Model 2c:**  $\lambda_{R(\beta_0)}^*[u | V^\dagger, \bar{A}] = \lambda_{R(\beta_0)}^*[u | \bar{A}[r^{-1}(u, \bar{A}, V^\dagger, \beta_0)], V^\dagger] = \lambda_0(u | V^\dagger)$  where  $R(\beta_0) \equiv R(\bar{A}, \beta_0)$  and  $r^{-1}(u, \bar{a}, V^\dagger, \beta) \equiv t$  if  $r(t, \bar{a}, V^\dagger, \beta) = u$ .

We shall now consider estimation of these models for the observables, under assumption (9), and the further assumption that

$$\bar{A}(C) \text{ has a known conditional distribution given } V^\dagger. \quad (10)$$

Semiparametric inference in the association models 1a - 2c without (10) imposed has been examined previously by many authors. Below we use their results to solve the estimation problem in our semiparametric model.

We will show that associated with each MSM model with (9) and (10) imposed is a class of regular asymptotically linear (RAL) estimators  $\{\hat{\beta}^*(h, \phi)\}$  for  $\beta_0$ , indexed by vector functions  $h \in \mathcal{H}$  and

$\phi \in \Phi$  such that the set  $\mathcal{IF}^* = \{IF^*(h, \phi)\}$  of influence functions of the  $\widehat{\beta}^*(h, \phi)$  constitute all the influence functions for the model, in the sense that if  $\widetilde{\beta}^*$  is any other RAL estimator, then the influence function of  $\widetilde{\beta}^*$  equals  $IF^*(h, \phi)$  for some functions  $h \in \mathcal{H}, \phi \in \Phi$ . Recall that an estimator  $\widetilde{\beta}$  of  $\beta_0$  is RAL with influence function  $IF$  if  $n^{\frac{1}{2}}(\widetilde{\beta} - \beta_0) = n^{-\frac{1}{2}} \sum_i IF_i + o_p(1)$ , the  $IF_i$  are i.i.d, and the convergence of  $\widetilde{\beta}$  to  $\beta_0$  is locally uniform. Here  $o_p(1)$  denotes a random variable converging in probability to zero. Thus a RAL estimator is asymptotically equivalent to a sum of the i.i.d random variables  $IF_i$ . We obtain  $\widehat{\beta}^*(h, \phi)$  by solving the estimating equations  $n^{-\frac{1}{2}} \sum_i \widehat{D}_i^*(\beta, h, \phi) = o_p(1)$  described below. [We put  $o_p(1)$  on the right side of the estimating equation to take care of cases (e.g., rank estimators) in which the estimating function  $\widehat{D}^*(\beta, h, \phi)$  is not continuous in  $\beta$  and, thus, the left-hand side of the previous equality may never be exactly zero.] The solution  $\widehat{\beta}^*(h, \phi)$  has influence function  $IF^*(h, \phi) = \{\kappa^*(h)\}^{-1} U^*(\beta_0, h, \phi)$  where  $U_i^*(\beta_0, h, \phi)$  depends only on subject  $i$ 's data,  $\kappa^*(h) = -\partial E^*[U^*(\beta, h, \phi)]/\partial \beta|_{\beta=\beta_0}$  does not depend on  $\phi$ , and  $n^{-\frac{1}{2}} \sum_i \widehat{D}_i^*(\beta_0, h, \phi) = n^{-\frac{1}{2}} \sum_i U_i^*(\beta_0, h, \phi) + o_p(1)$ . Furthermore,  $\Lambda^\perp = \{U^*(\beta_0, h, \phi)\}$  with  $h \in \mathcal{H}, \phi \in \Phi$  is the linear span of  $\mathcal{IF}^*$  and thus is the orthogonal complement to the nuisance tangent space for the model in the Hilbert space induced by the covariance norm. (Here we are quoting a well known result from the theory of semiparametric models. See Robins and Ritov (1997) for discussion.) We refer to  $U^*(\beta_0, h, \phi)$  as the influence function for the estimating function  $\widehat{D}(\beta_0, h, \phi)$ . More specifically,  $U^*(\beta, h, \phi)$  and  $\widehat{D}^*(\beta, h, \phi)$  are each expressed as the sum of the two components, one of which  $U_{tp}^*(\phi) = D_{tp}^*(\phi)$  is independent of the choice of the MSM and follows from the fact that, for the ‘‘treatment process (tp),’’ (9) and (10) are assumed. Specifically, if the  $A(t)$  can jump only at times  $0, 1, 2, \dots$ ,  $U_{tp}^*(\phi) = \sum_{k=0}^{int(C)} \phi(k, \overline{A}(k), \overline{L}(k^-)) - E^*[\phi(k, \overline{A}(k), \overline{L}(k^-)) | \overline{L}(k^-), \overline{A}(k^-)]$  where  $int(C)$  is the greatest integer less than or equal to  $C$ . It is easy to see that  $\{U_{tp}^*(\phi)\}$  is, as  $\phi$  varies, the sum over  $k$  of functions of the observed data  $(\overline{A}(k), \overline{L}(k^-))$  with mean zero given  $(\overline{A}(k^-), \overline{L}(k^-))$ . If  $A(t)$  is a continuous time marked point process, then  $U_{tp}^*(\phi) = \int dM_A^*(u) \phi_1(u, \overline{A}(u^-), \overline{L}(u^-)) + \int dN_A(u) \{\phi_2(u, \overline{A}(u), \overline{L}(u^-)) - E^*[\phi_2(u, \overline{A}(u), \overline{L}(u^-)) | A(u) \neq A(u^-), \overline{L}(u^-), \overline{A}(u^-)]\}$  where  $dM_A^*(u) = dN_A(u) - \lambda_A^*[u | \overline{A}(u^-), \overline{L}(u^-)] du$  and  $dN_A(u) = I\{A(u) \neq A(u^-)\}$  counts jumps in the  $\overline{A}$  process. [In the examples of the Introduction, we chose the function  $\phi$  to be identically zero so that  $\widehat{D}_{tp}^*(\phi)$  was also zero. As we shall see later, the choice  $\phi$  identically zero, although computationally convenient because we can then use standard software, is somewhat inefficient.]

The other structural model-specific component  $\widehat{D}_{sm}^*(\beta, h)$  and  $U_{sm}^*(\beta, h)$  of  $\widehat{D}^*(\beta, h, \phi)$  and  $U^*(\beta, h, \phi)$  are the well-known estimating functions and their associated influence functions for the association models 1a - 2c with neither (9) nor (10) imposed.

**Model 1a:**  $\widehat{D}_{sm}^*(\beta, h) = U_{sm}^*(\beta, h) = h(\overline{A}, V^\dagger) \varepsilon(\beta)$  with  $\varepsilon(\beta) = Y - g(\overline{A}, V^\dagger, \beta)$  and  $h(\overline{A}, V^\dagger)$  is any  $\dim(\beta)$  vector function. In the linear logistic cumulative treatment model of the Introduction,  $\widehat{D}_{sm}^*(\beta, h)$  was the score equation from the logistic model and thus  $h(\overline{A}, V^\dagger)$  was the vector  $(1, cum(\overline{a}))'$ .

**Model 1b:**  $\eta(x) = x : \widehat{D}_{sm}^*(\beta, h) = U_{sm}^*(\beta, h) = \{\varepsilon(\beta) - h_1(\overline{A}, V^\dagger)\} \{h_2(\overline{A}, V^\dagger) - E^*[h_2(\overline{A}, V^\dagger) | V^\dagger]\}$  where  $h_1$  is any real valued function,  $\varepsilon(\beta)$  is as just defined and the range of  $h_2$  is of  $\dim(\beta)$ .

$\eta(x) = \ln[x/(1-x)] : U_{sm}^*(\beta, h) = U^\dagger(h, P(\beta))$  and  $\widehat{D}_{sm}^*(\beta, h) \equiv U^\dagger(h, \widehat{P}(\beta))$ , where  $P(\beta) = expit[g(\overline{A}, V^\dagger, \beta) + g^\dagger(V^\dagger)]$ ,  $\widehat{P}(\beta) = expit[g(\overline{A}, V^\dagger, \beta) + \widehat{g}^\dagger(V^\dagger)]$ ,  $expit(x) = e^x/(1+e^x)$ ,  $\widehat{g}^\dagger(V^\dagger)$  is a  $n^{\frac{1}{4}}$ -consistent estimate of  $g^\dagger(V^\dagger)$ , and  $U^\dagger(h, P(\beta)) \equiv \{Y - P(\beta)\} \{h(\overline{A}, V^\dagger) - E^*[h(\overline{A}, V^\dagger) P(\beta) \{1 - P(\beta)\} | V^\dagger] / E^*[P(\beta) \{1 - P(\beta)\} | V^\dagger]\}$ .

**Model 1c:**  $\widehat{D}_{sm}^*(\beta, h) = U^*(\beta, h) = h[R(\beta), \overline{A}, V^\dagger] - \int h[R(\beta), \overline{a}, V^\dagger] dF^*[\overline{a} | V^\dagger]$ .

**Model 1d:** Let  $\varepsilon(\beta) = \{\varepsilon_1(\beta), \dots, \varepsilon_{K+1}(\beta)\}'$ ,  $\varepsilon_m(\beta) = Y(m) - g_m[\overline{A}(m-1), V^\dagger; \beta]$ . Then  $\widehat{D}_{sm}^*(\beta, h) =$

$U_{sm}^*(\beta, h) = h(\bar{A}, V^\dagger) \varepsilon(\beta)$  where  $h(\bar{A}, V^\dagger)$  is now any  $\dim(\beta) \times (K+1)$  matrix of real valued functions.

**Model 2a:**  $\hat{D}_{sm}^*(\beta, h) = \int_0^\infty dN_T(u) \{h(u, \bar{A}(u), V^\dagger) - \tilde{\mathcal{L}}(h, u, \beta)\}$ , where  $\tilde{\mathcal{L}}(h, u, \beta) = \tilde{J}[h, \beta] / \tilde{J}[\mathbf{1}, \beta]$ ; for any  $h(u, \bar{A}(u), V^\dagger)$ ,  $\tilde{J}(h, \beta) = \tilde{E}[h(u, \bar{A}(u), V^\dagger) I(T > u) \exp\{r[\bar{A}(u), u, V^\dagger, \beta]\}]$ ; for any  $H_i$ ,  $\tilde{E}(H) = \sum_{i=1}^n H_i/n$ ;  $\mathbf{1}$  is the constant function equal to one; and  $N_T(u) = I(T \leq u)$ .  $U_{sm}^*(\beta, h) = \int_0^\infty dM_T(u) \{h(u, \bar{A}(u), V^\dagger) - \mathcal{L}^*(h, u, \beta)\}$  where  $\mathcal{L}^*(h, u, \beta) = J^*[h, \beta] / J^*[\mathbf{1}, \beta]$ ;  $J^*[h, \beta]$  is defined like  $\tilde{J}(h, \beta)$  but with  $E^*$  replacing  $\tilde{E}$ ; and  $dM_T(u) = dN_T(u) - \lambda_T(u | \bar{A}, V^\dagger) I(T > u) du$ .

**Model 2b:**  $U_{sm}^*(\beta, h)$  and  $\hat{D}_{sm}^*(\beta, h)$  are as above except  $J^*(h, \beta) \equiv E^*[h(u, \bar{A}(u), V^\dagger) I(T > u) \exp\{r[\bar{A}(u), u, V^\dagger, \beta]\} | V^\dagger]$  and  $\tilde{J}(h, \beta)$  replaces  $E^*(\cdot | V^\dagger)$  in  $J^*(h, \beta)$  by a  $n^{\frac{1}{4}}$ -consistent estimator  $\hat{E}(\cdot | V^\dagger)$ .

**Model 2c:**  $\hat{D}_{sm}^*(\beta, h) = \int_0^\infty du I[R(\beta) > u] \{H_2(u, \beta) - E^*[H_2(u, \beta) | V^\dagger]\} + \int_0^\infty dN_{R(\beta)}(u) [H_1(u, \beta) - E^*[H_1(u, \beta) | V^\dagger]]$  and, for  $j = 0, 1$ ,  $H_j(u, \beta) = h_j[u, \bar{A}\{r^{-1}(u, \bar{A}, V^\dagger, \beta)\}, V^\dagger]$ .  $U_{sm}^*(\beta, h) = \hat{D}_{sm}^*(\beta, h) - E^*[D_{sm}^*(\beta, h) | \bar{A}, V^\dagger]$ .

**Remark:** Note that in model 2b and in model 1b with  $\eta(x) = \ln[x/(1-x)]$ , smooths are necessary to estimate  $g^\dagger(V^\dagger)$  and  $E^*(\cdot | V^\dagger)$  if  $V^\dagger$  has continuous components. In particular, due to the curse of dimensionality, it is not possible to obtain a reasonable  $n^{\frac{1}{2}}$ -consistent estimator of  $\beta_0$  in these models when  $V^\dagger$  has multiple continuous components. This can be formalized using the concept of curse of dimensionality appropriate (CODA) semiparametric information bounds introduced by Robins and Ritov (1997). Specifically, models 2b and 1b have CODA information bounds of zero, although they have positive ordinary semiparametric information bounds.

### 3.2 Non-ancillary treatment process

In this section, we no longer assume (8) is true. The essential idea of this section (requiring some minor modification) is to reweight  $\hat{D}_{sm}^*(\beta, h)$  by the inverse of a subject's probability of having had his observed treatment history. We continue to assume that

$$f[a(t) | \bar{L}(t^-), \bar{A}(t^-), V^\dagger] \text{ is known for } t \leq C \quad (11)$$

which implies that if  $A(t)$  jumps at non-random times  $0, \dots, K$ ,  $W(k) = f[A(k) | \bar{L}(k^-), \bar{A}(k^-)]$  and  $\bar{W}(k) = \prod_{m=0}^k W(k)$  are known. If  $A(t)$  jumps in continuous time,  $\bar{W}(t) = \exp\left[-\int_0^t \lambda_A[u | \bar{L}(u^-), \bar{A}(u^-)] du\right] \prod_{\{u; A(u) \neq A(u^-), u < t\}} \lambda_A[u | \bar{L}(u^-), \bar{A}(u^-)]$   $f[A(u) | \bar{A}(u^-), \bar{L}(u^-), A(u) \neq A(u^-)]$  is known.

We now come to a subtle but crucial idea. We need to artificially censor a subject at the first time  $C^\dagger$  that the density of receiving his observed treatment  $A(C^\dagger)$  at  $C^\dagger$  was zero for some prognostic factor history  $\bar{\ell}(C^\dagger)$  in order to insure that the reweighted  $\hat{D}_{sm}^*(\beta_0, h)$  still has asymptotic mean zero. We formalize this idea as follows. If  $A(t)$  jumps at non-random times, let  $\overset{\circ}{A}(k, \bar{a}(k^-), v^\dagger) = \{a(k); f[a(k) | \bar{L}(k^-), \bar{A}(k^-) = \bar{a}(k^-), V^\dagger = v^\dagger] \neq 0 \text{ w.p.1}\}$  and set  $C^\dagger = \min\left\{k; A(k) \notin \overset{\circ}{A}(k, \bar{A}(k^-), V^\dagger)\right\}$ . If  $A(t)$  jumps in continuous time, let  $\overset{\circ}{A}(t, \bar{a}(t^-), v^\dagger) = \{a(t); f[a(t) | \bar{L}(t^-), \bar{A}(t^-) = \bar{a}(t^-), A(t) \neq a(t^-), V^\dagger = v^\dagger] \neq 0 \text{ w.p.1} \text{ or } a(t) = A(t^-)\}$  and set  $C^\dagger = \inf\left\{t; A(t) \notin \overset{\circ}{A}(t, \bar{A}(t^-), V^\dagger)\right\}$ . The variable  $C^\dagger$  is crucial because, as indicated in the remark following

Lemma 3.1 below, one can only unbiasedly reweight a function of  $A(t)$  for  $A(t) \in \overset{\circ}{\mathcal{A}}(t, \bar{a}(t^-), v^\dagger)$ .

Let  $f^*(\bar{a} | V^\dagger)$  be a density (chosen by the analyst). Let  $F^*$  denote the joint distribution which differs from the true distribution  $F$  of  $O$  only in that  $f[a(t) | \bar{L}(t^-), \bar{A}(t^-), V^\dagger]$  is replaced by the ancillary density  $f^*[a(t) | \bar{A}(t^-), V^\dagger]$ . Further, define  $\mathcal{W}(t) \equiv \bar{W}(t) / f^*[\bar{A}(t) | V^\dagger]$  and  $O^\dagger = (\bar{L}(C^\dagger), \bar{A}(C^\dagger))$ . Note that in the examples of the Introduction, we chose  $f^*[a(t) | \bar{A}(t^-), V^\dagger]$  to be  $f[a(t) | \bar{A}(t^-), V^\dagger]$ . Note even with this choice, the distribution  $F$  of  $O$  differs from the distribution  $F^*$  if (9) is false. A key result, which follows from direct calculation, is

**Lemma 3.1:** For any  $z(O^\dagger)$ ,  $E[z(O^\dagger) / \mathcal{W}(C^\dagger) | V^\dagger] = E^*[z(O^\dagger) | V^\dagger]$ .

**Remark:** It is false that  $E[z(O) / \mathcal{W}(C) | V^\dagger] = E^*[z(O) | V^\dagger]$ .

Let  $D_{sm}^*(\beta, h)$  be the probability limit under  $F^*$  of  $\hat{D}_{sm}^*(\beta, h)$  and let  $\{U_{sm}(\beta, h)\}$  and  $\{D_{sm}(\beta, h)\}$  be the subsets of  $\{U_{sm}^*(\beta, h)\}$  and  $\{D_{sm}^*(\beta, h)\}$ , respectively, that depend on the data only through  $O^\dagger$ . Set  $\mathcal{W} = \mathcal{W}(C^\dagger)$  and note  $D_{sm}(\beta, h)$  and  $U_{sm}(\beta, h)$  often depend on  $E^*[\cdot | V^\dagger] = E[\cdot / \mathcal{W} | V^\dagger]$  or  $E^*[\cdot] = E[\cdot / \mathcal{W}]$ . Define  $\hat{D}_{sm}(\beta, h)$  and  $\hat{U}_{sm}(\beta, h)$  like  $D_{sm}(\beta, h)$  and  $U_{sm}(\beta, h)$  except replace any unknown expectations  $E[\cdot / \mathcal{W} | V^\dagger]$  and  $E[\cdot / \mathcal{W}]$  with appropriate estimates  $\hat{E}[\cdot / \mathcal{W} | V^\dagger]$  and  $\hat{E}[\cdot / \mathcal{W}]$ .

**Examples:** In model 1b, with  $\eta(x) = x$ ,  $U_{sm}(\beta, h) = \hat{D}_{sm}(\beta, h) = U_{sm}^*(\beta, h)$  has  $h_1(\bar{A}, V^\dagger)$  and  $h_2(\bar{A}, V^\dagger)$  being functions only of  $\{\bar{A}(C^\dagger), V^\dagger\}$ . Note  $E^*[h_2(\bar{A}, V^\dagger) | V^\dagger]$  is known and need not be estimated.

In contrast, in models 2a and 2b,  $\hat{D}_{sm}(\beta, h)$  will be defined like  $\hat{D}_{sm}^*(\beta, h)$  except in defining  $\tilde{J}(h, \beta)$  we replace  $I(T > u)$  by  $I(T > u) / \mathcal{W}$  in order to estimate the unknown expectations. Alternatively, we can replace  $I(T > u)$  by  $I(T > u) / \mathcal{W}(u)$ . This latter choice (i) will in general have better finite sample properties, (ii) tend to increase efficiency unless the estimator with  $I(T > u) / \mathcal{W}$  was already semiparametric efficient, and (iii) was the approach we took in the Introduction. The issues are exactly those discussed in Robins (1993), which the reader may consult for further clarification.

In model 1b with  $\eta(x) = \ln[x / (1 - x)]$ ,  $\hat{g}(V^\dagger) \equiv \hat{g}(V^\dagger, \beta)$  could be chosen to minimize  $\tilde{E}[(Y - \text{expit}\{g(\bar{A}, V^\dagger, \beta) + g^\dagger(V^\dagger)\})^2 / \mathcal{W}]$  over  $g^\dagger(V^\dagger)$  in some class (e.g., splines), whose dimension may increase with sample size.

In the appendix we briefly sketch a proof of the following.

**Theorem 3.1:** Subject to regularity conditions, in the semiparametric model (i) characterized by (4), (11), the data  $O$ , and a MSM, the class  $\{\hat{\beta}(h, \phi)\}$  with  $h \in \mathcal{H}$  and  $\phi \in \Phi$  of estimators which solve  $0 = \sum_i \hat{D}_i(\beta, h, \phi)$  with  $\{\hat{D}(\beta, h, \phi)\} = \{\hat{D}_{sm}(\beta, h) / \mathcal{W} + D_{tp}(\phi)\}$  is a class of RAL estimators with influence functions  $\mathcal{IF} = \{IF(h, \phi)\}$ ,  $IF(h, \phi) = \{\kappa(h)\}^{-1} U(\beta_0, h, \phi)$ ,  $\kappa(h) = -\partial E[U(\beta, h, \phi)] / \partial \beta|_{\beta=\beta_0}$ ,  $U(\beta_0, h, \phi) = U_{sm}(\beta_0, h) / \mathcal{W} + U_{tp}(\phi)$ , where  $U_{tp}(\phi) = D_{tp}(\phi)$  is defined like  $U_{tp}^*(\phi)$  except with the true law  $F$  replacing  $F^*$ . Furthermore,  $\mathcal{IF}$  is the set of all influence functions.

### 3.3 Efficiency for fixed $h$

We now begin to explore efficiency issues. By a projection argument similar to that given in Robins et al. (1994), we have

**Theorem 3.2:** For a given  $h$ , among all estimators  $\hat{\beta}(h, \phi)$ , the most efficient has  $\phi$  equal to  $\phi_{opt} \equiv \phi_{opt}(h)$ : if  $A(t)$  only jumps at non-random times  $0, 1, 2, \dots$  then  $\phi_{opt} \equiv 0$  if  $k > C^\dagger$ , and if  $k \leq C^\dagger$ ,  $\phi_{opt}[k, \bar{a}(k), \bar{\ell}(k^-)] = E[U_{sm}(h) / \mathcal{W} | \bar{A}(k) = \bar{a}(k), \bar{L}(k^-) = \bar{\ell}(k^-)] = \{\bar{W}(k)\}^{-1} \iint d\mu(\underline{a}_{k+1}) f^*(\bar{a} | v^\dagger) E[u_{sm}\{\bar{a}(C^\dagger), \bar{Y}_{\bar{a}}(C^\dagger), V^\dagger, h\} | \bar{L}_{\bar{a}}(k^-) = \bar{\ell}(k^-), \bar{A}(k) = \bar{a}(k)]$  and  $U_{sm}(h) \equiv U_{sm}(\beta_0, h)$ . Furthermore, if CAR holds,  $\bar{A}(k) = \bar{a}(k)$  can be removed from the last conditioning event above. If  $A(t)$  jumps in continuous time,  $\phi_{1,opt} = E[U_{sm}(h) / \mathcal{W} | \bar{L}(u^-), \bar{A}(u^-), A(u) \neq A(u^-)] - E[U_{sm}(h) / \mathcal{W} | \bar{L}(u^-), \bar{A}(u^-)]$



since  $E [U_{sm}(h) / \mathcal{W} | \bar{L}(u^-), \bar{A}(u^-)] = E [U_{sm}(h) / \mathcal{W} | \bar{L}(u^-), \bar{A}(u^-), A(u) = \bar{A}(u^-)]$ , and  $\phi_{2,opt} = E [U_{sm}(h) / \mathcal{W} | \bar{A}(u), \bar{L}(u^-)]$ .

We now relax unrealistic assumption (11) that the conditional density of  $A(t)$  is known. We consider two cases. In the first case the density is completely unknown and in the second the density follows a parametric model.

**Theorem 3.3: a)** The semiparametric model (ii) characterized by (4), data  $O$ , and a MSM (with (11) not imposed), has the set of influence functions  $\{IF(h, \phi_{opt}(h))\}$  with  $h \in \mathcal{H}$ .

**b):** In model (iii) characterized by (4), data  $O$ , a MSM, and a parametric model indexed by parameter  $\alpha$  for  $f[a(t) | \bar{L}(t^-), \bar{A}(t^-)]$ , the set of influence functions for the model is the set

$\left\{ \kappa(h)^{-1} \left[ U(h, \phi) - E[U(h, \phi) S'_\alpha] \{E[S_\alpha S'_\alpha]\}^{-1} S_\alpha \right] \right\}$  of influence functions of  $\left\{ \hat{\beta}(h, \phi, \hat{\alpha}) \right\}$  solving  $o_p(1) = n^{-\frac{1}{2}} \sum_i \hat{D}_i(\beta, h, \phi, \hat{\alpha})$  where  $\hat{\alpha}$  is the MLE of  $\alpha$ ,  $S_\alpha = \partial \log \{f[A(t) | \bar{L}(t^-), \bar{A}(t^-), \alpha] / \partial \alpha\}$  is the subject-specific score for  $\alpha$ , and  $D(\beta, h, \phi, \hat{\alpha})$  is  $D(\beta, h, \phi)$  evaluated at  $\hat{\alpha}$ .

Theorem 3.3b can be extended to semiparametric models for  $f(a(t) | \bar{L}(t^-), \bar{A}(t^-))$  such as a Cox proportional hazard model as in Robins (1993; 1998b, App. 2). In the non- and semi-parametric case, we have to plug an estimator of  $f[a(t) | \bar{L}(t^-), \bar{A}(t^-)]$  into the estimating function  $\hat{D}_i(\beta, h, \phi)$ . In general, the estimator needs to converge to the true density at a rate greater than  $n^{\frac{1}{4}}$  to obtain a RAL estimator of  $\beta$ .

### 3.4 Censoring

No new idea is required to account for and adjust for right censoring. Specifically, let  $Q$  be censoring time in MSMs. Define a censoring process  $A_2(u)$  by  $A_2(u) = 0$  if  $Q > u$  and  $A_2(u) = 1$  otherwise. Let the treatment of interest be  $a_1(u)$  and define  $a(u) = (a_1(u), a_2(u))$  and write  $T_{\bar{a}}$  as  $T_{\bar{a}_1, \bar{a}_2}$ . To want to adjust for censoring is only to say that interest is in the direct effect of  $\bar{a}_1$  when  $\bar{a}_2 \equiv 0$ , i.e., when censoring is abolished. As a concrete example, the Cox model MSM 2a in the presence of censoring would become

$$\lambda_{T_{\bar{a}_1, \bar{a}_2 \equiv 0}}(t | V^\dagger) = \lambda_0(t) \exp \{r[\bar{a}_1(t^-), t, V^\dagger; \beta_0]\}.$$

If  $\bar{A}$  is ancillary (now including the censoring process),  $\hat{D}_{sm}^*(\beta, h)$  and  $U_{sm}^*(\beta, h)$  are as above except that now  $N_T(u) = I[T \leq u] I[T < Q]$  and  $I[T > u]$  is everywhere replaced by  $I[T > u] I[Q > u]$ . Of course  $\bar{W}(t)$  is now the probability that a subject would have his observed treatment and censoring history. This is exactly the approach we took in analyzing the 002 trial data in the Introduction.

## 4 Semiparametric Efficiency

### 4.1 The efficient score

In any semiparametric model, the semiparametric variance bound is the inverse of the variance of the efficient score  $S_{eff}$ . The efficient score in models (i) - (iii) of Theorems 3.1 and 3.2 are the same and, by Theorem 5.3 in Newey and McFadden (1993), equal  $S_{eff} = U(\beta_0, h_{eff}, \phi_{eff})$  where  $\phi_{eff} = \phi_{opt}(h_{eff})$  and  $h_{eff}$  is uniquely characterized by the requirement that for all  $U(\beta_0, h, \phi)$

$$E[U(\beta_0, h, \phi) U(\beta_0, h_{eff}, \phi_{opt}(h_{eff}))'] = \kappa(h)$$

which is equal to

$$E[U_{sm}(\beta_0, h) U(\beta_0, h_{eff}, \phi_{opt}(h_{eff}))'] = \kappa(h) . \tag{12}$$

To show how to use (12) to calculate  $h_{eff}$ , we consider the following simple example.

**Model 1a:** Consider MSM 1a with  $C^\dagger = C = K + 1 = 1$  w.p.1 so  $K = 0$  and the  $\bar{A}(C)$  process only jumps at time zero. So  $\bar{A} = A(0)$  and  $\mathcal{W}^{-1} = f^*[\bar{A} | V^\dagger] / f[\bar{A} | L(0)]$  where  $V^\dagger \subset L(0) = V(0)$ . For the purposes of computing the efficient score, we can choose  $f^*(\bar{A} | V^\dagger) = 1$  w.p.1 without worrying that it is not a density, because it can be absorbed into  $h_{eff}(\bar{A}, V^\dagger)$ . In Appendix 2, we prove the following.

**Theorem 4.1:** With  $f^*(\bar{A} | V^\dagger) = 1$  w.p.1, Eq. 12 implies  $h_{eff}(\bar{A}, V^\dagger)$  is the unique solution to the type two Fredholm equation  $h_{eff}(\bar{A}, V^\dagger) [\int var[\varepsilon | \bar{A}, L(0)] \{f(\bar{A} | L(0))\}^{-1} f(V^\bullet | V^\dagger) d\mu(V^\bullet)] + \int h_{eff}(\bar{a}, V^\dagger) \omega(\bar{a}, \bar{A}, V^\dagger) d\mu(\bar{a}) = \partial g(\bar{A}, V^\dagger, \beta_0) / \partial \beta$  where  $V^\bullet = L(0) \setminus V^\dagger$  and  $\omega(\bar{a}, \bar{A}, V^\dagger) = [\int E[\varepsilon | \bar{a}, L(0)] E[\varepsilon | \bar{A}, L(0)] f(V^\bullet | V^\dagger) d\mu(V^\bullet)]$ . Note that if  $\bar{A}$  has finite support, this is a finite dimensional matrix equation. Our estimators, specialized to this example, are continuous treatment extensions of efficient propensity score estimators of an average treatment effect. By dividing by the propensity score, we eliminate the bias due to within stratum confounding that can occur with subclassification on the propensity score as recommended by Rosenbaum and Rubin (1983)

## 4.2 Efficiency calculations using missing data theory

Given (4), imposing CAR cannot change the efficient score. Thus, it is of interest to rederive the efficient score using the Hilbert space results of van der Vaart (1991) and of Robins et al. (1994) for missing data models under CAR. For convenience, assume  $C$  is non-random and write  $\bar{A} \equiv \bar{A}(C)$ . The full data are  $\bar{L}^F = \{\bar{L}_{\bar{a}}; \bar{a} \in \bar{A}\}$ . Given any  $B = b(\bar{L}^F)$ , the score operator  $\mathbf{s}(B) = E[B | O], O = (\bar{A}, \bar{L}_{\bar{A}})$ . For any  $Q = q(O)$ , the non-parametric adjoint operator  $\mathbf{s}^\dagger$  under CAR is  $\mathbf{s}^\dagger(Q) = E[Q | \bar{L}^F] = \int d\mu(\bar{a}) q(\bar{a}, \bar{L}_{\bar{a}}) f[\bar{a} | \bar{L}_{\bar{a}}]$ . Suppose for the remainder of this subsection that the  $\bar{A}$  jumps only at times  $0, 1, \dots, K$  and  $\bar{L}$  jumps at  $0^-, \dots, K + 1^-$  and  $C = K + 1$ . We then have by CAR

$$f[\bar{a}(k) | \bar{L}_{\bar{a}}] = \prod_{m=0}^k f[a(m) | \bar{a}(m-1), \bar{L}_{\bar{a}}(m)] \quad (13)$$

and  $f[\bar{a} | \bar{L}_{\bar{a}}] = f[\bar{a}(K) | \bar{L}_{\bar{a}}]$ . It is then easy to check the null space of  $\mathbf{s}^\dagger$ ,  $N(\mathbf{s}^\dagger) = \{U_{tp}(\phi)\}$ . Now define the non-parametric information operator,  $\mathbf{m} = \mathbf{s}^\dagger \mathbf{s} : \bar{L}^F \rightarrow R(\mathbf{s}^\dagger)$  where  $R(\mathbf{s}^\dagger)$  is the range of  $\mathbf{s}^\dagger$ . Note that  $R(\mathbf{s}^\dagger) = \{B = \int d\mu(\bar{a}) b(\bar{a}, \bar{L}_{\bar{a}})\}$ . Let  $\mathbf{m}^{-1} : R(\mathbf{s}^\dagger) \rightarrow R(\mathbf{s}^\dagger)$  be the inverse of  $\mathbf{m}$  on  $R(\mathbf{s}^\dagger)$ . Given  $\bar{a}_1, \bar{a}_2$ , let  $u_{12}$  be the smallest  $u$  with  $a_1(u) \neq a_2(u)$ . We then have by a direct calculation

**Theorem 4.2:** If for all  $\bar{a}_1, \bar{a}_2$

$$\bar{L}_{\bar{a}_1} \amalg \bar{L}_{\bar{a}_2} | \bar{L}_{\bar{a}_1}(u_{12}^-) \quad (14)$$

then

$$\mathbf{m}^{-1} \left[ \int d\mu(\bar{a}) b(\bar{a}, \bar{L}_{\bar{a}}) \right] = \int d\mu(\bar{a}) \left\{ \sum_{m=1}^{K+1} \{f[\bar{a}(m-1) | \bar{L}_{\bar{a}}]\}^{-1} \right. \\ \left. \{E[b(\bar{a}, \bar{L}_{\bar{a}}) | \bar{L}_{\bar{a}}(m)] - E[b(\bar{a}, \bar{L}_{\bar{a}}) | \bar{L}_{\bar{a}}(m-1)]\} + E[b(\bar{a}, \bar{L}_{\bar{a}}) | \bar{L}_{\bar{a}}(0)] \right\} . \quad (15)$$

**Remark:** Gill and Robins (1999) show that (14) places no restriction on the law of the observed data  $O$  even when sequential randomization and a MSM are imposed. We can and do always assume that (14) holds.

Now let  $S_{eff}^F$  and  $\Lambda^{F,\perp}$  be the efficient score and the orthogonal complement to the nuisance tangent space for the parameter  $\beta$  of our marginal structural model when we have data on  $\bar{L}^F$ . Then the efficient score  $S_{eff}$  based on data  $O$  under CAR is  $g[m^{-1}(D_{eff})]$  where  $D_{eff}$  is the unique member of

$\Lambda^{F,\perp} \cap R(\mathbf{s}^\dagger)$  satisfying

$$\Pi [\mathbf{m}^{-1} (D_{eff}) | \Lambda^{F,\perp}] = S_{eff}^F, \quad (16)$$

where  $\Pi$  is the Hilbert space projection operator. To show how to use this result to calculate  $S_{eff}$ , we revisit the example given in the last subsection.

**Example: Model 1a:** Consider MSM 1a as in Sec. 4.1. Then, by an extension of Theorem 8.3 of Robins et al. (1994)

$$\Lambda^{F,\perp} = \left\{ \int d\mu(\bar{a}) h(\bar{a}) \varepsilon(\bar{a}) \right\} \quad (17)$$

where (i)  $\varepsilon(\bar{a}) = \varepsilon(\bar{a}, \beta_0)$ , (ii)  $h(\bar{a})$  is a vector valued function of the dimension of  $\beta_0$ . Note that  $\Lambda^{F,\perp}$  is contained in  $R(\mathbf{s}^\dagger)$  as will be the case for MSMs with positive information.

**Remark:**

$$\text{If } \bar{\mathcal{A}} = \{\bar{a}_1, \dots, \bar{a}_S\} \text{ is finite,} \quad (18)$$

then  $\varepsilon(\bar{a})$  can be identified with the  $S$  vector that has components  $\varepsilon_s(\bar{a}_s) = Y_{\bar{a}_s} - g(\bar{a}_s, V^\dagger, \beta_0)$ . For arbitrary  $\bar{\mathcal{A}}$ ,  $\varepsilon(\bar{a})$  is a stochastic process with index set  $\bar{\mathcal{A}}$ . For  $\bar{a}, \bar{a}^* \in \bar{\mathcal{A}}$ , let  $\mathbf{cv}(\bar{a}, \bar{a}^*) = \text{cov}(\varepsilon(\bar{a}), \varepsilon(\bar{a}^*))$ . If  $\bar{\mathcal{A}}$  is given by (18),  $\mathbf{cv}(\bar{a}, \bar{a}^*)$  corresponds to the  $S \times S$  matrix with  $j, k$  entry  $\mathbf{cv}(\bar{a}_j, \bar{a}_k)$ . Let  $\mathbf{cv}^{-1}(\bar{a}^*, \bar{a}^*)$  be a (generalized) inverse of  $\mathbf{cv}(\bar{a}, \bar{a}^*)$ , i.e., by definition, for any function  $q(\bar{a}^*)$ ,  $\int [\int \mathbf{cv}^{-1}(\bar{a}^*, \bar{a}) \mathbf{cv}(\bar{a}, \bar{a}^*) d\mu(\bar{a})] q(\bar{a}^*) d\mu(\bar{a}^*) = q(\bar{a}^*)$ . In particular, if (18) holds,  $\mathbf{cv}^{-1}(\bar{a}^*, \bar{a})$  is just the inverse of the matrix identified with  $\mathbf{cv}(\bar{a}, \bar{a}^*)$ . Then, generalizing Chamberlain (1987),

$$S_{eff}^F = \int d\mu(\bar{a}) \left\{ \partial g(\bar{a}, V^\dagger; \beta_0) / \partial \beta \right\} \left[ \int \mathbf{cv}^{-1}(\bar{a}, \bar{a}^*) \varepsilon(\bar{a}^*) d\mu(\bar{a}^*) \right]. \quad (19)$$

If  $\bar{\mathcal{A}}$  is given by (18),  $\partial g(\bar{a}, V^\dagger; \beta_0) / \partial \beta$  can be identified with the  $\dim \beta \times S$  matrix with  $j, k$  entry  $\partial g(\bar{a}_k, V^\dagger; \beta_0) / \partial \beta_j$ . Again, generalizing Theorem 8.3 in Robins et al. (1994),

$$\begin{aligned} \Pi \left[ \int d\mu(\bar{a}) b(\bar{a}, \bar{\mathcal{L}}_{\bar{a}}) | \Lambda^{F,\perp} \right] &= \int d\mu(\bar{a}) E \left[ b(\bar{a}, \bar{\mathcal{L}}_{\bar{a}}) \varepsilon(\bar{a}) | V^\dagger \right] \\ &\left[ \int \mathbf{cv}^{-1}(\bar{a}, \bar{a}^*) \varepsilon(\bar{a}^*) d\mu(\bar{a}^*) \right]. \end{aligned} \quad (20)$$

Hence to solve (16), we need to find the solution  $h_{eff}(\bar{a}, V^\dagger)$  to the equation

$$E[\mathbf{m}^{-1} \left\{ \int d\mu(\bar{a}^*) h(\bar{a}^*, V^\dagger) \varepsilon(\bar{a}^*) \right\} \varepsilon(\bar{a}) | V^\dagger] = \partial g(\bar{a}, V^\dagger; \beta_0) / \partial \beta. \quad (21)$$

By (15) with  $K = 0$  and the fact that, by CAR,  $f(\bar{a} | \bar{\mathcal{L}}_{\bar{a}}) = f(\bar{a} | L(0))$ , the LHS of (21) can be written

$$\begin{aligned} E \left\{ \int d\mu(\bar{a}^*) h(\bar{a}^*, V^\dagger) \left\{ [\varepsilon(\bar{a}^*) - E[\varepsilon(\bar{a}^*) | L(0)]] \{f(\bar{a}^* | L(0))\}^{-1} + \right. \right. \\ \left. \left. E[\varepsilon(\bar{a}^*) | L(0)] \varepsilon(\bar{a}) | V^\dagger \right\} \right\} = E \left\{ \int d\mu(\bar{a}^*) h(\bar{a}^*, V^\dagger) \right. \\ \left. [cov[\varepsilon(\bar{a}^*), \varepsilon(\bar{a}) | L(0)] \{f(\bar{a}^* | L(0))\}^{-1} + E[\varepsilon(\bar{a}^*) | L(0)] E[\varepsilon(\bar{a}) | L(0)]] | V^\dagger \right\}. \end{aligned}$$

However, since by assumption (14),  $Y_{\bar{a}_j} \perp Y_{\bar{a}_k} | \bar{L}(0)$  for  $k \neq j$ , (21) reduces to

$$\begin{aligned} h(\bar{a}, V^\dagger) E \{ var[\varepsilon(\bar{a}) | L(0)] f(\bar{a} | L(0))^{-1} | V^\dagger \} + \\ \int d\mu(\bar{a}^*) h(\bar{a}^*, V^\dagger) E \{ E[\varepsilon(\bar{a}^*) | L(0)] E[\varepsilon(\bar{a}) | L(0)] | V^\dagger \} \\ = \partial g(\bar{a}, V^\dagger, \beta_0) / \partial \beta. \end{aligned}$$

Upon noting that, by CAR,  $f[\varepsilon | \bar{\mathcal{A}} = \bar{a}, L(0)] = f[\varepsilon(\bar{a}) | L(0)]$ ,  $w$ . We see that this is the same expression for  $h_{eff}(\bar{a}, V^\dagger)$  as obtained in Theorem 4.1.

### 4.3 A practical approach to obtaining reasonable efficiency

Estimation of  $h_{eff}$  is computationally difficult because of the need to solve integral equations without closed form solutions. A practical approach to choosing  $h$  and  $f^*(\bar{a} | V^\dagger)$  is important. Given a model for  $f[a(t) | \bar{a}(t^-), \bar{\ell}(t^-)]$  depending on parameter  $\alpha' = (\alpha'_1, \alpha'_2)$  such that  $\alpha_1 = 0 \Leftrightarrow f[a(t) | \bar{a}(t^-), \bar{\ell}(t^-)] = f[a(t) | \bar{a}(t^-), v^\dagger]$ , rather than choosing  $f^*[\bar{a} | V^\dagger]$ , we use  $f^*[\bar{a} | V^\dagger; \tilde{\alpha}_2]$  where  $\tilde{\alpha}_2$  is the MLE of  $\alpha_2$  with  $\alpha_1$  set to zero. This is exactly the approach we took in analyzing the 002 data in the Introduction. [The fact that  $f^*[\bar{a} | V^\dagger]$  is estimated does not influence the asymptotic distribution of  $\hat{\beta}(h, \phi)$ .] It follows that if (8) holds [i.e.,  $\bar{A}$  is an ancillary process],  $\mathcal{W}$  will converge to 1. Further, in each of the models 1a – 2c of Sec. 3.1, the efficient choice of  $h$ , say,  $h_{opt}$ , for solving  $\sum_i \hat{D}_{sm,i}^*(h) = 0$  when (9) is not imposed is well known. We suggest choosing  $h$  to be  $h_{opt}$  or an estimate  $\hat{h}_{opt}$  thereof, and choosing  $\phi$  to be an estimate of  $\phi_{opt}(\hat{h}_{opt})$ . Such a choice guarantees that if  $\bar{A}$  is an ancillary process, our estimate of  $\beta_0$  will be more efficient than the estimate based on solving  $0 = \sum_i \hat{D}_{sm,i}(h_{opt})$ . Specifically, in MSM 1a,  $h_{opt} = \{\partial \varepsilon(\beta_0) / \partial \beta\} \{var(\varepsilon(\beta_0) | \bar{A}, V^\dagger)\}^{-1}$ . For model 1b, Chamberlain (1988) gives  $h_{1,opt}$  and  $h_{2,opt}$ . In model 1c,  $h_{opt} = [\partial / \partial \beta] [\ln \{[\partial R(\beta_0) / \partial Y] f[R(\beta_0) | V^\dagger]\}]$ . In model 1d,  $h_{opt}$  is as in model 1a with  $\varepsilon(\beta_0)$  now a vector. In model 2a,  $h_{opt} = \partial \ln r[\bar{A}(u^-), u, V^\dagger, \beta_0] / \partial \beta$ . For model 2b,  $h_{opt}$  is given by Sasieni (1992). In model 2c,  $h_{opt,2} = h_{opt,1} \lambda_0(u | V^\dagger)$  and  $h_{opt,1} = \partial \ln \lambda_{R(\beta)}(u | V^\dagger) / \partial \beta|_{\beta=\beta_0}$ .

## 5 Comparison of MSMs and SNMs

We begin by recalling the definition of a structural nested distribution model.

### 5.1 Structural Nested Distribution Models

For concreteness, we consider the setting of the MSM model 1a - 1c with  $Y_{\bar{a}} = Y_{\bar{a}}(K+1)$  and the  $A$  process and  $L$  process jumping at non-random times  $0, \dots, K$  and  $0^-, \dots, K+1^-$  respectively. Henceforth, we take  $V^\dagger = \emptyset$ . Suppose  $Y$  is a continuous variable with a continuous distribution function  $F_Y(y) = pr[Y < y]$ . Let  $(\bar{a}(m), 0)$  denote the treatment history given by  $\bar{a}(m)$  through time  $m$  and zero at times  $m+1, \dots, K$ . Then let  $\gamma(y, \bar{\ell}(m), \bar{a}(m))$  be the unique function mapping quantiles of  $Y_{\bar{a}(m),0}$  into those of  $Y_{\bar{a}(m-1),0}$  conditional on  $\bar{L}(m) = \bar{\ell}(m), \bar{A}(m) = \bar{a}(m)$  so that  $\gamma(y, \bar{\ell}(m), \bar{a}(m))$  measures the magnitude of the effect of a final blip of treatment  $a(m)$  on quantiles of  $Y$  among subjects with observed history  $\{\bar{\ell}(m), \bar{a}(m)\}$ . A structural nested distribution model (SNDM) is a parametric model for this function. That is, it specifies that  $\gamma(y, \bar{\ell}(m), \bar{a}(m)) = \gamma(y, \bar{\ell}(m), \bar{a}(m), \beta_0)$  where  $\gamma(y, \bar{\ell}(m), \bar{a}(m), \beta)$  is a known increasing function of  $y$  satisfying  $\gamma(y, \bar{\ell}(m), \bar{a}(m), \beta) = y$  if  $a(m) = 0$  or  $\beta = 0$ . Recursively define random variables  $\dot{R}_K(\beta), \dots, \dot{R}_0(\beta)$  by  $\dot{R}_K(\beta) = \gamma(Y, \bar{L}(K), \bar{A}(K), \beta)$  and  $\dot{R}_m(\beta) = \dot{r}_m(Y, \bar{L}(K), \bar{A}(K), \beta) = \gamma(\dot{R}_{m+1}(\beta), \bar{L}(m), \bar{A}(m), \beta)$  and set  $\dot{R}(\beta) = \dot{r}(Y, \bar{L}(K), \bar{A}(K), \beta) \equiv \dot{R}_0(\beta)$ . [Heuristically,  $\dot{R}_m(\beta_0)$  is  $Y_{\bar{A}(m-1),0}$  and  $\bar{R}(\beta_0)$  is  $Y_0$ , where  $Y_0$  is the outcome when treatment is always withheld. This is heuristic because in fact it is only the conditional distributions through time  $m$  that are guaranteed to be the same.] Also let  $\dot{r}^{-1}(y, \bar{\ell}(K), \bar{a}(K), \beta)$  be the inverse of the function  $\dot{r}$  with respect to its first argument. If  $\gamma(y, \bar{\ell}(m), \bar{a}(m), \beta) = \gamma(y, \bar{a}(m), \beta)$  does not depend on  $\bar{\ell}(m)$  for each  $m$ , we say that the SNDM model has no interaction.

**Theorem 5.1:** Under (4), a no interaction SNDM model is a stratified transformation model (STM), i.e., MSM 1c, with  $R(\bar{a}, \beta) = \dot{r}(Y, \bar{a}, \beta)$  and  $R(\beta) = \dot{R}(\beta)$ . However, the converse is not true.

That is, a no-interaction SNDM is a MSM. The semiparametric information bound for  $\beta$  is greater if

we correctly impose a no-interaction SNDM than if we only imposed the corresponding STM. An SNDM will be a MSM only if (as a fact of nature) there is no interaction. Theorem 5.1 indicates that a STM is the natural MSM analog of a SNDM in this case. If  $\gamma(y, \bar{\ell}(m), \bar{a}(m))$  depends on  $\bar{\ell}(m)$ , we must choose between analyzing the data under a SNDM versus a STM. To understand the advantages and disadvantages of each, we need some additional background. Define a regime  $g = (g_0, \dots, g_K) \in \mathcal{G}$  to be a collection of functions  $g_m : \bar{\mathcal{L}}_m \rightarrow \mathcal{A}_m$ . Define  $g(\bar{\ell}(m)) = \{g_0(\bar{\ell}_0), \dots, g_m(\bar{\ell}(m))\}$ . Let  $Y_g$  be the counterfactual value of  $Y$  if regime  $g$  were followed. If  $g(\bar{\ell}_K) = \bar{a}_K \equiv \bar{a}$  does not depend on  $\bar{\ell}_K$ , then  $Y_g = Y_{\bar{a}}$  and we say  $g$  is non-dynamic; otherwise,  $g$  is dynamic. Let  $g[\bar{\ell}(k)]$  denote a realization of  $\bar{A}(k)$ . If we have sequential ignorability for regime  $g$ , i.e.,

$$Y_g \amalg A(t) \mid \bar{L}(t^-), \bar{A}(t^-) , \quad (22)$$

then, by Theorem 3.2 of Robins (1997), the law of  $Y_g$  is given by the G-computation algorithm formula

$$F_{Y_g}(y \mid \bar{\ell}(k), g[\bar{\ell}(k-1)]) = \iint F_Y(y \mid \bar{\ell}(K), g(\bar{\ell}(K))) \prod_{m=k+1}^K dF[\ell(m) \mid \bar{\ell}(m-1), g(\bar{\ell}(m-1))] . \quad (23)$$

We obtain  $F_{Y_g}(y)$  from (23) by substituting in  $k = -1$ . Using the fact that, for continuous  $Y$ ,  $Y$  is  $\dot{r}^{-1} \left( \dot{R}(\beta_0), \bar{L}(K), \bar{A}(K), \beta_0 \right)$ , it can be shown that (23) implies

$$F_{Y_g}(y) = \iint I[\dot{r}^{-1} \{u, \bar{\ell}(K), g(\bar{\ell}(K)), \beta_0\} > y] \prod_{m=0}^K dF[\ell(m) \mid \bar{\ell}(m-1), g(\bar{\ell}(m-1)), \dot{R}(\beta_0) = u] dF_{\dot{R}(\beta_0)}(u) . \quad (24)$$

In many settings, the g-null hypothesis that

$$F_{Y_{g1}}(y) = F_{Y_{g2}}(y) \text{ for all } g1, g2 \in \mathcal{G} \quad (25)$$

will be of interest. This is implied by the sharp null hypothesis of no treatment effect that  $Y_{g1} = Y_{g2}$  with probability 1, i.e. no subject's outcome is influenced by the treatment history they choose. Robins (1986, 1997) proves the following.

**Theorem 5.2:** Given (22) for  $g \in \mathcal{G}$ , (25) holds  $\Leftrightarrow$  (23) is the same for all  $g \Leftrightarrow \gamma(y, \bar{\ell}_m, \bar{a}_m) = y \Leftrightarrow$

$$Y \amalg A(k) \mid \bar{L}(k), \bar{A}(k-1), k = 0, \dots, K . \quad (26)$$

## 5.2 Advantages of SNDMs with a Continuous $Y$

We are now ready to compare the advantages and disadvantages of SNDMs and MSMs for continuous  $Y$ . We begin by reviewing the advantages of SNDMs.

1. Although (26) implies  $\beta_0 = 0$  for both a SNDM and a STM, only for a SNDM is (26) equivalent to  $\beta_0 = 0$ . What this means causally is the following. For a STM, the null hypothesis  $\beta_0 = 0$  is equivalent to the hypothesis that the distribution of  $Y_{\bar{a}}$  is the same for all non-dynamic regimes  $\bar{a}$ . This is a weaker hypothesis than the g-null hypothesis that says the distribution of  $Y_g$  is the same for all regimes, whether non-dynamic or dynamic. In most cases, it will be the latter null hypothesis (25) that will be of public health interest unless it were not possible to collect data on the covariates  $L_k$  which determine the treatment decisions for dynamic regimes.

2. If the  $L(k)$  are discrete with only a moderate number of levels, then, even with  $f[a(k) | \bar{\ell}(k-1), \bar{a}(k-1)]$  totally unrestricted, an asymptotically distribution-free g-null test of  $\beta_0 = 0$  (and thus of (25)) exists for a SNDM but, because of the curse of dimensionality, not for a STM. Specifically, a non-parametric g-null test is equivalent to a test of independence of  $Y$  and  $A(k)$  within strata defined jointly by  $\bar{L}(k), \bar{A}(k-1)$  (Robins, 1997). Thus, even if  $A(k)$  is continuous, a test of independence of  $A(0)$  and  $Y$  within levels of  $L(0)$  will be an asymptotic  $\alpha$ -level test under (25). In contrast, a test of  $\beta_0 = 0$  in a STM (without (25) additionally imposed) requires, by the Remark following Theorem 3.3a, that  $\mathcal{W}$  can be consistently non-parametrically estimated which will not be possible due to the curse of dimensionality. (Note that to estimate  $\mathcal{W}$ , we must be able to consistently estimate the density of  $A(k)$  given  $\bar{L}(k)$  and  $\bar{A}(k-1)$  for all  $k$ , which is not possible to do non-parametrically when the  $A(m)$  are continuous.) In other words, the stronger hypothesis (25) that  $\beta_0 = 0$  for a SNDM is easier to test non-parametrically than the weaker hypothesis that  $\beta_0 = 0$  for a STM.
3. Henceforth, assume a correct model for  $f[a(t) | \bar{\ell}(t^-), \bar{a}(t^-)]$  is available for all  $t$ . (We remind the reader that this would often be a false assumption.) Given a SNDM, with some difficulty the law of  $Y_g$  for dynamic  $g$  can be estimated using (24). In contrast, the law of  $Y_g$  for dynamic  $g$  is very hard to estimate given a STM. Specifically, given a SNDM, we estimate the law of  $Y_g$  as follows: (i) obtain an estimate  $\hat{\beta}$  by g-estimation (Robins, 1997), (ii) estimate  $F_{R(\beta_0)}^\bullet(u)$  by the empirical law of the  $\dot{R}_i(\hat{\beta})$  for  $i = 1, \dots, n$ , (iii) specify and estimate a parametric model for  $f\left[L(m) | \bar{L}(m-1), \bar{A}(m-1), \dot{R}(\hat{\beta})\right]$ , (iv) and then evaluate the estimated version of the integral (24) by Monte carlo.
- In contrast, given a STM, we must, as discussed in Robins (1997, pg. 114; 1998a, Sec. 11) and Robins et al. (1999), specify a parametric model for  $\nu^*(y, \bar{\ell}(m), \bar{a})$ , where one can choose to define  $\nu^*(y, \bar{\ell}(m), \bar{a})$  in either of two ways, leading to different parameterizations. Either  $\nu^*(y, \bar{\ell}(m), \bar{a}) \equiv \nu(y, \bar{\ell}(m), \bar{a}) - \nu(y, \{\bar{\ell}(m-1), \ell(m) = 0\}, \bar{a})$  and  $\nu(y, \bar{\ell}(m), \bar{a})$  maps quantiles of  $Y_{\bar{a}}$  given  $\bar{\ell}_m, \bar{a}_{m-1}$  into quantiles of  $Y_{\bar{a}}$  given  $\bar{\ell}_{m-1}, \bar{a}_{m-1}$ , or  $\nu^*(y, \bar{\ell}(m), \bar{a})$  is defined to be the ratio of the hazard evaluated at  $y$  of  $Y_{\bar{a}}$  given  $\bar{\ell}_m, \bar{a}_{m-1}$  to the hazard at  $y$  of  $Y_{\bar{a}}$  given  $\bar{\ell}_{m-1}, \bar{a}_{m-1}, \ell(m) = 0$ . Robins et al. (1999, Sec.8.7a) argue for the second option, since, in contrast to the first option, a parameterization in terms of hazard ratios is variation independent. As discussed in Robins (1997, pg. 114-116; 1998a, Sec. 11) and Robins et al. (1999), estimation of  $\nu^*(y, \bar{\ell}(m), \bar{a})$  is a computational nightmare; indeed, fully parametric Bayesian or likelihood-based inference for a MSM is computationally extremely burdensome.
4. As discussed in Robins (1997, Sec. 9; 1998bc) and Robins et al. (1999), for SNDMs, it is easy to perform a sensitivity analysis in which the fundamental assumption (22) of ignorable treatment assignment is no longer imposed. For a STM such a sensitivity analysis is somewhat less straightforward and sensitivity analysis methods for MSMs are described in Robins et al. (1999) and Appendix 3 below.
5. A parameter  $\beta_0$  of a SNDM, in contrast to that of a STM, can often still be consistently estimated if (22) is false but data are available on an instrumental variable. Specifically, suppose  $A(t) = (A_1(t), A_2(t))$  with  $A_1(t)$  recording a physician's prescribed treatment and  $A_2(t)$  recording treatment actually received. We might suppose (22) is false, but  $A_1(t) \perp\!\!\!\perp Y_g | \bar{L}(t^-), \bar{A}(t^-)$  is true if a predictor of  $Y_g$  and of  $A_2(t)$  was not recorded in  $\bar{L}(t^-)$ .  $A_1(t)$  is often then referred to as

an instrumental variable process, particularly when  $A_1(t)$  has no direct causal effect, i.e.,  $Y_{\bar{a}} = Y_{\bar{a}_2}$  w.p.1. In this setting, the parameter of a STM is not identified but the parameter of a SNDM can still in general be consistently estimated by g-estimation (Robins, 1993; 1998b).

4. MSMs, in contrast to SNMs, cannot be used if there exists a value of  $\ell_k$ , say  $\ell_k = 0$ , such that for all but one  $a_k \in \mathcal{A}_k$ ,  $f[a_k | \bar{\ell}_{k-1}, \ell_k = 0, \bar{a}_{k-1}] = 0$ , since then the artificial censoring time  $C^\dagger$  is zero with probability 1. An example would be a study of the effect of an occupational exposure on mortality with  $\ell_k = 0$  if a subject is off work at time  $k$ ,  $\ell_k = 1$  otherwise, and subjects off work can only receive exposure level  $a_k = 0$ .

### 5.3 Advantages of MSMs with Continuous $Y$ or with Failure Time Outcomes

1. Even in the presence of interaction [i.e.,  $\gamma(y, \bar{\ell}_m, \bar{a}_m)$  depends on  $\bar{\ell}_m$ ], given a STM, the distribution of a non-dynamic counterfactual outcome  $F_{Y_{\bar{a}}}(y)$  can be estimated by  $n^{-1} \sum_i I\left\{r^{-1} \left[R_i(\hat{\beta}), \bar{a}, \hat{\beta}\right] > y\right\}$  without requiring either integration or modelling of the conditional law of  $L(m)$ . In contrast, as described in point 3 of Sec. 5.2 above, for a SNDM, both integration and modelling are required.
2. Any MSM that can be easily estimated when (8) holds (i.e.,  $\bar{A}$  is an ancillary process) can be easily estimated when (8) is false. For example, we can use the Cox proportional hazards MSM 2a for a continuous failure time outcome  $T_{\bar{a}}$ . In contrast, a structural nested Cox model would model the ratio of the conditional hazard given  $\bar{\ell}_m, \bar{a}_m$  of  $Y_{\bar{a}(m),0}$  to that of  $Y_{\bar{a}(m-1),0}$  as a function of an unknown finite-dimensional parameter. Unfortunately, a structural nested Cox model does not admit any simple semiparametric estimators, and even complex estimators will fail due to the curse of dimensionality. Formally, the CODA information bound of Robins and Ritov (1997) for a structural nested Cox model is zero, even when  $f(a_k | \bar{\ell}_k, \bar{a}_{k-1})$  is completely known.

A possible hybrid approach is to impose a MSM model and then specify a model for  $\nu^*(y, \bar{\ell}(m), \bar{a})$  is as follows. The g-null hypothesis (26) is true if and only if the distribution of  $Y_{\bar{a}}$  given  $V^\dagger$  is the same for all  $\bar{a}$  (i.e., the parameter  $\beta_0$  of our MSM is zero) and  $\nu^*(y, \bar{\ell}(m), \bar{a})$  depends on  $\bar{a}$  only through  $\bar{a}_{m-1}$  (Robins, 1997, Appendix 3). Thus, we impose a MSM model depending on a parameter  $\beta$  and an additional model for  $\nu^*(y, \bar{\ell}(m), \bar{a})$  that depends on both the parameter  $\beta$  of the MSM model and another parameter  $\psi = (\psi_1, \psi_2)$  in such a way that  $\beta\psi_1 = 0$  if and only if  $\nu^*(y, \bar{\ell}(m), \bar{a})$  depends on  $\bar{a}$  only through  $\bar{a}_{m-1}$ . Specifically,  $\nu^*(y, \bar{\ell}(m), \bar{a}) / \nu^*(y, \bar{\ell}(m), (\bar{a}_{m-1}, a_m = 0, \dots, a_K = 0))$  depends on  $(\beta, \psi)$  only through the product  $\beta\psi_1$  and  $\nu^*(y, \bar{\ell}(m), (\bar{a}_{m-1}, a_m = 0, \dots, a_K = 0))$  depends only on  $\psi_2$ . Thus  $\psi_1$  is identified only if  $\beta \neq 0$ . Such a model can overcome objections 1 and 2 of Sec. 5.2 (but not objections 3-6) while retaining the advantages 1-2 of Sec. 5.3.

### 5.4 Structural Nested Mean Models (SNMMs)

We now turn to comparing MSMs and SNMs for discrete outcomes. Consider the set up of Sec. 5.1 but with  $Y$  discrete. For discrete outcomes, we define structural nested mean models. However, SNMMs are applicable to discrete and continuous outcomes.

Let  $\gamma(\bar{\ell}(m), \bar{a}(m)) = E[Y_{\bar{a}(m),0} - Y_{\bar{a}(m-1),0} | \bar{\ell}(m), \bar{a}(m)]$ . Let  $\gamma^\dagger(\bar{\ell}(m), \bar{a}(m)) = \ln\{E[Y_{\bar{a}(m),0} | \bar{\ell}(m), \bar{a}(m)] / E[Y_{\bar{a}(m-1),0} | \bar{\ell}(m), \bar{a}(m)]\}$ . An additive structural nested mean model

(SNMM) specifies  $\gamma(\bar{\ell}(m), \bar{a}(m)) = \gamma(\bar{\ell}(m), \bar{a}(m), \beta_0)$  with  $\gamma(\bar{\ell}_m, \bar{a}_m, \beta)$  a known function satisfying  $\gamma(\bar{\ell}_m, \bar{a}_m, \beta) = 0$  if  $a_m = 0$  or  $\beta = 0$ . A multiplicative SNMM specifies  $\gamma^\dagger(\bar{\ell}(m), \bar{a}(m)) = \gamma(\bar{\ell}(m), \bar{a}(m), \beta_0)$ . The g-null mean hypothesis is the hypothesis

$$E[Y_{g1}] = E[Y_{g2}], g_1, g_2 \in \mathcal{G}. \quad (27)$$

Robins (1997) proves the following.

**Theorem 5.2:** Given (22), (27) holds if and only if  $\gamma(\bar{\ell}(m), \bar{a}(m)) = 0 \Leftrightarrow \gamma^\dagger(\bar{\ell}(m), \bar{a}(m)) = 0 \Leftrightarrow E[Y | \bar{A}(k), \bar{L}(k)] = E[Y | \bar{A}(k-1), \bar{L}(k)], k = 0, \dots, K$ :

Advantages (1) - (4) and (6) of Sec. 5.2 of a SNM over a MSM for continuous  $Y$  also will hold (appropriately modified) for discrete  $Y$  when considering the g-null mean hypothesis or when estimating  $E[Y_g]$ . Advantages (1) - (2) in Sec. 5.3 of a MSM over a SNM for continuous outcome also hold in the discrete case.

An important advantage of MSMs over SNMs with  $Y$  dichotomous (or, more generally, when  $Y$  has finite support) is that neither an additive SNMM or multiplicative SNMM naturally imposes the fact that, for dichotomous  $Y$ ,  $E[Y_g] \in [0, 1]$ . In contrast, using the MSM model 1a with  $g(\bar{a}, \beta)$  a logistic function, the above restriction is naturally imposed. Analogously, in the setting of MSM model 1d, we can use standard marginal logistic models for the repeated measures outcomes  $Y_{\bar{a}}(m)$ . There exists logistic SNMMs that do impose that  $E[Y_g] \in [0, 1]$  (Robins et al., 1999). However, these logistic SNMMs are not very useful for semiparametric inference with high-dimensional data, since the CODA information bound for the parameter  $\psi$  of interest is zero even when  $f(a_k | \bar{\ell}_k, \bar{a}_{k-1})$  is known.

## 5.5 Direct Effect Models

In this section we show that some of the advantages of MSMs described in Secs. 5.3 and 5.4 are not retained in semiparametric models for the direct effect of a treatment  $a_1$  when a second treatment  $a_2$  is held fixed (set). In such a setting, both MSMs and direct effect SNMs (Robins, 1998a) have important limitations due to the curse of dimensionality if the functional form of the effect of the second treatment  $a_2$  on the outcome is left completely unrestricted. Let  $a(u) = (a_1(u), a_2(u))$  and, in a slight abuse of notation, set  $\bar{a}(u) = (\bar{a}_1(u), \bar{a}_2(u))$  and  $\bar{a} = (\bar{a}_1, \bar{a}_2)$ . Continue to assume  $V^\dagger = \emptyset$ . Consider the following.

**Model 3a – direct effect semiparametric regression:** Consider the set-up of MSM 1b, with

$$\eta\{E[Y_{\bar{a}}]\} = g[\bar{a}, \beta_0] + g^\dagger(\bar{a}_2)$$

where  $g[\bar{a}, \beta_0] = 0$  if  $\bar{a}_1 \equiv 0$  and  $g^\dagger(\cdot, \cdot)$  is unknown and unrestricted. Since, according to the model,  $\eta\{E[Y_{\bar{a}_1, \bar{a}_2}]\} - \eta\{E[Y_{\bar{a}_1=0, \bar{a}_2}]\} = g(\bar{a}, \beta_0)$ , it follows we are modelling the direct effect of treatment  $\bar{a}_1$ . Furthermore, the main effect of the second treatment  $g^\dagger(\bar{a}_2) = \eta\{E[Y_{\bar{a}_1=0, \bar{a}_2}]\} - \eta\{E[Y_{\bar{a}_1=0, \bar{a}_2=0}]\}$  is completely unrestricted. Under sequential randomization assumption (4), the model for the observables  $O$  induced by MSM 3a is isomorphic to that induced by MSM 1b with  $\bar{A}_2 \equiv \bar{A}_2(K)$  playing the role of  $V^\dagger$ . In particular, if  $\eta(x) = x$  [or  $\ln(x)$ ],  $\hat{\beta}(h, \phi)$  will perform well in moderate size samples, provided  $f[a(t) | \bar{L}(t^-), \bar{A}(t^-)]$  is known or can be parametrically modelled. However, as discussed in the final remark of Sec. 3.1, if  $\eta(x) = \ln[x/(1-x)]$ , reasonable estimators of  $\beta_0$  are unavailable because  $\bar{A}_2(K)$  will be high-dimensional. Indeed, any choice of  $\eta(x)$  that guarantees that  $E[Y_{\bar{a}}] \in [0, 1]$  will fail to provide reasonable estimators of  $\beta_0$ , negating the advantage of this MSM for dichotomous  $Y$ . This reflects the fact that the CODA information bound is zero, even when  $f(a_k | \bar{\ell}_k, \bar{a}_{k-1})$  is known.



**Model 3b – direct effect semiparametric Cox proportional hazards model:** Consider the set up of MSM 2b, with

$$\lambda_{T_{\bar{a}}}(t) = \lambda_{T_{\bar{a}_1=0, \bar{a}_2}}(t) \exp [r \{ \bar{a}(t^-), t; \beta_0 \}]$$

with  $r(\bar{a}(t^-), t; \beta_0) = 0$  if  $\bar{a}_1(t^-) = \mathbf{0}$ . This is a model for the direct effect of treatment  $\bar{a}_1$  on the hazard of  $T$  with the main effect of  $\bar{a}_2$  left unrestricted. Given (4), MSM 3b induces a model for the observables isomorphic to that induced by MSM 2b. This implies that, as discussed in the remark in Sec. 3.1, due to the curse of dimensionality, it will not be possible to obtain reasonable estimators of  $\beta_0$  negating advantage 2 of Sec. 5.3.

**Model 3c – direct effect semiparametric time-dependent accelerated failure time model:** Strikingly, the accelerated failure time MSM we now develop does not suffer from degradation due to the curse of dimensionality as did model 3b. Consider the model

$$\lambda_{R(\bar{a}, \beta_0)}(t) = \lambda_{T_{\bar{a}_1=0, \bar{a}_2}}(t)$$

where  $R(\bar{a}, \beta) = r(T_{\bar{a}}, \bar{a}, \beta)$  satisfies  $r(t, \bar{a}, \beta) = t$  if  $\bar{a}_1 = \mathbf{0}$  or  $\beta = 0$ . The model for the observables induced by MSM 3c is isomorphic to that induced by MSM 2c with  $\bar{A}_2$  in the role of  $V^\dagger$ . Hence, the association model 3c can be used to estimate the direct effect of  $\bar{a}_1$  on  $T$  with the main effect of  $\bar{a}_2$  unrestricted. MSM 3c is the natural MSM associated with a structural nested failure time model (SNFTM) (Robins, 1993, App. 1; 1998) since a direct-effect SNFTM without interaction is a MSM 3c. In the presence of interaction, the MSM 3c retains advantage 1 of Sec. 5.3.

## Appendix 1:

By arguments as in Robins et al. (1994), Theorem 3.1 and 3.3b are easy corollaries of Theorem 3.3a.

**Sketch of Proof of Theorem 3.3a:** For convenience, assume the  $L$  process and  $A$  process jump at times  $0^-, 1^-, \dots$  and  $0, 1, \dots$  respectively. Then by Theorem 3.2 of Robins (1997), Eq. (4) implies the G-computation algorithm formula

$$f_{\bar{Y}_{\bar{a}}(k)} [\bar{y}(k) | v^\dagger] = \int \prod_{m=0}^k f [y(m), v(m) | \bar{y}(m-1), \bar{v}(m-1), \bar{a}(m-1), v^\dagger] \prod_{m=1}^k d\mu [v(m)] d\mu (\dot{v}), \quad (\text{A1})$$

with  $\dot{v} \equiv v(0) \setminus v^\dagger$ . Thus, if for some  $j < k$ , the proposition  $f [a(j) | \bar{a}(j-1), \bar{L}(j), v^\dagger] \neq 0$  w.p.1 given  $V^\dagger$  is false, then (A1) is not identified. Hence, a MSM model places no (local) restrictions on  $f [L(m) | \bar{L}(m-1), \bar{A}(m-1)]$  for  $m > C^\dagger$ . Hence, in semiparametric model (ii), every function of  $O$  with mean zero given  $O^\dagger$  is in the nuisance tangent space for the model. It follows that all members of  $\Lambda^\perp$  in model (ii) depend on the data only through  $O^\dagger$ .

Because of our assumed knowledge of  $\Lambda^{\perp,*}$  (the orthogonal complement to the nuisance tangent space under  $F^*$ ), it is sufficient to show that  $U \in \Lambda^\perp \Leftrightarrow UW \in \Lambda^{\perp,*}$  when  $F^*$  is chosen such that  $C^{\dagger,*}$  is equal to  $C^\dagger$ . This follows from the fact that  $\mathring{\Lambda} = \mathring{\Lambda}^*$  where  $\mathring{\Lambda} \equiv \Lambda \cap \{U_{tp}(\phi)\}^\perp \cap \{z(O^\dagger)\}$  and the fact that  $E[UB] = E^*[UWB]$  for any  $B \in \mathring{\Lambda}$  by Lemma 3.1.

## Appendix 2

**Proof:** In our model, Eq. (12) states that

$$E[h(\bar{A}, V^\dagger) \varepsilon \mathcal{W}^{-1} \{h_{eff}(\bar{A}, V^\dagger) \varepsilon \mathcal{W}^{-1} - h_{eff}(\bar{A}, V^\dagger) \mathcal{W}^{-1} E[\varepsilon | \bar{A}, L(0)] +$$

$E [h_{eff}(\bar{A}, V^\dagger) \mathcal{W}^{-1} E [\varepsilon | \bar{A}, L(0)] | L(0)] = \kappa(h) \equiv E [h(\bar{A}, V^\dagger) \mathcal{W}^{-1} \partial g(\bar{A}, V^\dagger, \beta_0) / \partial \beta]$ . This can be rewritten as

$$\begin{aligned} E[h(\bar{A}, V^\dagger) \{h_{eff}(\bar{A}, V^\dagger) \mathcal{W}^{-2} var [\varepsilon | \bar{A}, L(0)] + \mathcal{B}(h_{eff})\}] = \\ E [h(\bar{A}, V^\dagger) \mathcal{W}^{-1} \partial g(\bar{A}, V^\dagger, \beta_0) / \partial \beta] \end{aligned} \quad (\text{A2.1})$$

where  $\mathcal{B}(h_{eff}) = \varepsilon \mathcal{W}^{-1} E \{h_{eff}(\bar{A}, V^\dagger) \mathcal{W}^{-1} E [\varepsilon | \bar{A}, L(0)] | L(0)\}$ . Now (A2.1) is true for all  $h(\bar{A}, V^\dagger)$  if and only if

$$\begin{aligned} h_{eff}(\bar{A}, V^\dagger) E [\mathcal{W}^{-2} var [\varepsilon | \bar{A}, L(0)] | \bar{A}, V^\dagger] + \\ E [\mathcal{B}(h_{eff}) | \bar{A}, V^\dagger] = E [\mathcal{W}^{-1} \partial g(\bar{A}, V^\dagger, \beta_0) / \partial \beta | \bar{A}, V^\dagger] . \end{aligned} \quad (\text{A2.2})$$

To simplify (A2.2), note for any  $q[\bar{A}, L(0)]$ ,  $E [q(\bar{A}, L(0)) \mathcal{W}^{-1} | \bar{A}, V^\dagger] = f^*(\bar{A} | V^\dagger) \int q(\bar{A}, L(0)) \{f(\bar{A} | L(0))\}^{-1} \{f(\bar{A} | L(0), V^\dagger) f[L(0) | V^\dagger] / f(\bar{A} | V^\dagger)\} d\mu(V^\bullet) = f^*(\bar{A} | V^\dagger) \{f(\bar{A} | V^\dagger)\}^{-1} \int q(\bar{A}, L(0)) f(V^\bullet | V^\dagger) d\mu(V^\bullet)$ . Thus, we have

$$E [\mathcal{W}^{-1} | \bar{A}, V^\dagger] = f^*(\bar{A} | V^\dagger) / f(\bar{A} | V^\dagger) . \quad (*)$$

$$\begin{aligned} E [\mathcal{W}^{-2} var [\varepsilon | \bar{A}, L(0)] | \bar{A}, V^\dagger] = \\ f^*(\bar{A} | V^\dagger) \{f(\bar{A} | V^\dagger)\}^{-1} \int \mathcal{W}^{-1} var [\varepsilon | \bar{A}, L(0)] f(V^\bullet | V^\dagger) d\mu(V^\bullet) \end{aligned} \quad (**)$$

and

$$\begin{aligned} E [\mathcal{B}(h_{eff}) | \bar{A}, V^\dagger] = \\ f^*(\bar{A} | V^\dagger) \{f(\bar{A} | V^\dagger)\}^{-1} \int E [\varepsilon | \bar{A}, L(0)] f(V^\bullet | V^\dagger) d\mu(V^\bullet) \\ \left\{ \int E [\varepsilon | \bar{a}, L(0)] f^*(\bar{a} | V^\dagger) h_{eff}(\bar{a}, V^\dagger) d\mu(\bar{a}) \right\} = \\ \int h_{eff}(\bar{a}, V^\dagger) f^*(\bar{a} | V^\dagger) d\mu(\bar{a}) \omega(\bar{a}, \bar{A}, V^\dagger) . \end{aligned} \quad (***)$$

Substituting \*, \*\*, \*\*\* into (A2.2) proves the theorem.

### Appendix 3: Sensitivity analysis for continuous and failure-time outcomes

Suppose rather than making the assumption (4) of sequential randomization, we instead assume for a model with a continuous outcome  $Y$  measured at end of follow-up at time  $K + 1^-$  (e.g., model 1c of Sec. 2.2) or a continuous failure-time outcome  $T$  (models 2a - 2c of Sec. 2.2) the existence of a known function  $q_m(y, \bar{\ell}_m, \bar{a}, a_m^*)$  such that for  $Y_{\bar{a}}$  the continuous counterfactual variable measured at end of follow-up

$$pr [Y_{\bar{a}} < y | \bar{\ell}(m), \bar{a}(m-1), a^*(m)] = pr [q(Y_{\bar{a}}, \bar{\ell}(m), \bar{a}, a^*(m)) < y | \bar{\ell}(m), \bar{a}(m)] \quad (\text{A3.1})$$

and, for  $T_{\bar{a}}$  the counterfactual failure time variable,

$$pr [T_{\bar{a}} < y | \bar{\ell}(m), \bar{a}(m-1), a^*(m), T \geq m] = pr [q(T_{\bar{a}}, \bar{\ell}(m), \bar{a}, a^*(m)) < y | \bar{\ell}(m), \bar{a}(m), T \geq m] . \quad (\text{A3.2})$$

The chosen conditional quantile - quantile function  $q_m(y, \bar{\ell}_m, \bar{a}, a_m^*)$  must satisfy

$$q_m(y, \bar{\ell}(m), \bar{a}, a^*(m)) = y \text{ if } a^*(m) = a(m) \quad (\text{A3.3})$$

$$q_m(y, \bar{\ell}(m), \bar{a}, a^*(m)) \text{ is increasing in } y \quad (\text{A3.4})$$

and, in the failure time case (A3.2),

$$q_m(y, \bar{\ell}(m), \bar{a}, a^*(m)) > m \quad (\text{A3.5a})$$

and writing  $\mu = q_m(y, \bar{\ell}_m, \bar{a}, a_m^*)$ , then

$$q_m(y, \bar{\ell}_m, \bar{a}, a_m^*) \text{ is a function of } \bar{a} \text{ only through } \bar{a}[\max(y, u)] . \quad (\text{A3.5b})$$

This last restriction follows by consistency assumption (2). Note that the sequential randomization assumption (4) implies that  $q_m(y, \bar{\ell}_m, \bar{a}, a_m^*) \equiv y$ .

We now sketch how to construct a regular asymptotically linear (RAL) estimator of the parameter  $\beta_0$  of a MSM such as model (1c) or model (2a) - (2c). Consider first the continuous outcome  $Y$  measured at end of follow-up. We shall replace each subject's observed outcome  $Y$  with  $J$  pseudo- $Y$ 's obtained by use of the following algorithm.

- Step 1: Do for  $j = 1, \dots, J$ .
- Step 2: Do for  $m = 0, \dots, K$ .
  - Draw  $a_j^*(m)$  from  $f[a(m) | \bar{L}(m), \bar{A}(m-1)]$ .
- Step 3: Set  $Y_{(K+1)j} = Y$ .
- Step 4: Do for  $m = K, \dots, 0$   $Y_{mj} = q_m(Y_{(m+1)j}, \bar{L}(m), \bar{A}, a_j^*(m))$ .
- Create a new data set with  $n \times J$  observations,  $O_{ij} = (\bar{A}_i(K), \bar{L}_i(K), Y_{i0j})$ ,  $i = 1, \dots, n, j = 1, \dots, J$ . Now for each observation  $O_{ij}$ , calculate  $\hat{D}_{sm}(\beta, a)$  and  $D_{tp}(\phi)$  as described in the paragraph following Lemma 3.1, with  $Y_{i0j}$  in place of the actual data  $Y_i$ . Then let  $\hat{\beta}(h, \phi)$  solve  $0 = \sum_{i=1}^n \sum_{j=1}^J \hat{D}_{ij}(\beta, h, \phi)$  where, for each observation  $O_{ij}$ ,  $\hat{D}(\beta, h, \phi) = \hat{D}_{sm}(\beta, h) / W + D_{tp}(\phi)$ .

Then it can be shown that, subject to regularity conditions, under the model characterized by (A3.1), an appropriate MSM, such as model 1c, and (11),  $\hat{\beta}(h, \phi)$  will be a RAL estimator of  $\beta_0$ .

The above algorithm can be modified so as to apply to a study with failure time outcomes under the assumption that the treatment process only jumps at times  $0, 1, 2, 3, \dots$  as follows.

To describe our algorithm for failure-time outcomes, we first discuss how to obtain a function  $q_m(y, \bar{\ell}_m, \bar{a}, a_m^*)$  guaranteed to satisfy (A3.5a) and (A3.5b).

Define, for  $m \leq [y-1]$  where  $[x]$  is the greatest integer less than or equal to  $x$ , the function  $q^*(x, \bar{\ell}(m), \bar{a}[y], a^*(m))$  by

$$\begin{aligned} & pr [T_{(\bar{a}[y], 0)} > u | \bar{\ell}(m), \bar{a}(m-1), a^*(m), t \geq m, T_{(\bar{a}[y-1], 0)} > [y]] = \\ & pr [q^*(T_{(\bar{a}[y-1], 0)}, \bar{\ell}(m), \bar{a}[y], a^*(m)) > u | \bar{\ell}(m), \bar{a}(m-1), a^*(m), T_{(\bar{a}[y-1], 0)} > [y]] . \end{aligned} \quad (\text{A3.6})$$

Then  $q_m(y, \bar{\ell}(m), \bar{a}, a^*(m))$  is determined by  $q_m^*(x, \bar{\ell}(m), \bar{a}[y], a^*(m))$  and  $q_m(y, \bar{\ell}(m), (\bar{a}[y], 0), a_m^*)$  through the following algorithm.

- Step 1:  $p \leftarrow q_m(y, \bar{\ell}(m), (\bar{a}[y], 0), a_m^*)$ .

- Step 2: Do for  $k = 1, 2, \dots$   
if  $p < [y + k]$ , stop and declare  $q_m(y, \bar{\ell}(m), \bar{a}, a^*(m)) = p$ .  
Otherwise,  $p \leftarrow q_m^*(p, \bar{\ell}(m), \bar{a}[y + k], a^*(m))$ .

In conducting a sensitivity analysis, we choose  $q_m(y, \bar{\ell}(m), (\bar{a}[y], 0), a^*(m))$  and  $q_m^*(x, \bar{\ell}(m), \bar{a}[y], a^*(m))$  restricted only by the fact that the first function  $q_m$  is increasing in  $y$ , exceeds  $m$  and equals  $y$  if  $a^*(m) = a(m)$  and that the second function  $q_m^*$  is increasing in  $y$  and exceeds  $[y]$ . Then we use the above algorithm to compute  $q_m(y, \bar{\ell}(m), \bar{a}, a_m^*)$ , which will then be guaranteed to satisfy (A3.5a) and (A3.5b).

We now describe how to construct a RAL estimator for the parameter of  $\beta_0$  of a MSM failure time model such as model (2a). We assume we have on each of  $n$  subjects the data  $\Delta = I(T = C), X = \min(T, C), \bar{A}(X), \bar{L}(X)$  where  $T$  is the failure time variable and  $C$  is the censoring variable. The algorithm goes as follows.

- Step 1: Do for  $j = 1, \dots, J$
- Step 2: Set  $K = [X]$
- Step 3: For  $s = K + 1, K + 2, \dots$ 
  - Draw  $a_j(s)$  from a chosen density  $f^*[a(s) | \bar{A}_K, a_j(K + 1), \dots, a_j(s - 1)]$
- Step 4: Set  $X_{(K+1)j} = X$
- Step 5: Do for  $m = K, K - 1, \dots, 0$ 
  - Draw  $a_j^*(m)$  from  $f[a(m) | \bar{L}(m), \bar{A}(m - 1), T > m]$
  - Define  $\bar{A}_j = (\bar{A}(K), a_j(K + 1), a_j(K + 2), \dots)$  and set  $X_{mj} = q(X_{(m+1)j}, \bar{L}(m), \bar{A}_j, a_j^*(m))$
- Step 6: Create a new data set with  $n \times J$  observations

$$O_{ij} = (\bar{A}_{ij}(X_{ioj}), X_{ioj}, \Delta_i) .$$

We then fit the Cox model (2a) as we described previously in the paper but based on the  $n \times j$  observations  $O_{ij}$ , with each observation on subject  $i$  associated with the same weight  $\mathcal{W}_i$ , where in calculating the numerator of  $\mathcal{W}_i \equiv \mathcal{W}(X_i)$ , we must use density  $f^*$  that was used in step 3 above. The resulting estimator will be consistent under the assumption that the hazard of censoring at time  $t$  given all the data only depends on the observed past.

Robins et al. (1999, Sec. 8.7b) discuss some potential problems with the sensitivity analysis methods discussed in this section due to the lack of a variation independent parameterization.

## BIBLIOGRAPHY:

- Chamberlain, G. 1987. Asymptotic Efficiency in Estimation with Conditional Moment Restrictions, *Journal of Econometrics*, 34: 305-324.
- Chamberlain, G. 1988. Efficiency bounds for semiparametric regression. *Technical Report*, Department of Statistics, University of Wisconsin.

- Gill, R.D., van der Laan, M.J., & Robins, J.M. 1997. Coarsening at random: characterizations, conjectures and counterexamples. *Proceedings of the First Seattle Symposium on Survival Analysis*, pp. 255-294.
- Gill, R.D. and Robins, J.M. 1999. Causal inference from complex longitudinal data: The continuous case. Unpublished manuscript.
- Heitjan, D.F., and Rubin, D.B. 1991. Ignorability and Coarse Data. *The Annals of Statistics*, 19: 2244-2253.
- Holland, P. 1986. Statistics and Causal Inference. *Journal of the American Statistical Association*, 81: 945-961.
- Lewis, D. 1973. Causation. *Journal of Philosophy*, 70: 556-567.
- Liang, K-Y., and Zeger, S.L. 1986. Longitudinal Data Analysis Using Generalized Linear Model. *Biometrika*, 73: 13-22.
- Lin, D.Y., Wei, L-J. 1989. The robust inference for the Cox proportional hazards model. *Journal of the American Statistical Association*, 84: 1074-1078.
- Newey, W.K. and McFadden, D. 1993. Estimation in large samples. **Handbook of Econometrics**, Vol. 4. Eds. McFadden, D., Engler, R. Amsterdam: North Holland.
- Pearl J. 1995. Causal Diagrams for Empirical Research. *Biometrika*, 82: 669-688.
- Robins J.M.1986. A new approach to causal inference in mortality studies with a sustained exposure period - application to control of the healthy worker survivor effect. *Mathematical Modeling*, 7: 1393-1512.
- Robins J.M. 1987. Addendum to "A new approach to causal inference in mortality studies with sustained exposure periods-Application to control of the healthy worker survivor effect." *Computers and Mathematics with Applications*, 14: 923-945.
- Robins J.M. 1993. Analytic methods for HIV treatment and cofactor effects. **AIDS Epidemiology - Methodological Issues**. Eds. Ostrow DG; Kessler R. Plenum Publishing, New York. pp. 213-290.
- Robins, J.M. 1994. Correcting for non-compliance in randomized trials using structural nested mean models. *Communications in Statistics*, 23: 2379-2412.
- Robins, J.M. 1997. Causal inference from complex longitudinal data. In: **Latent Variable Modeling and Applications to Causality. Lecture Notes in Statistics (120)**. M. Berkane, Editor. NY: Springer Verlag. pp. 69-117.
- Robins, J.M. 1998a. Testing and estimation of direct effects by reparameterizing directed acyclic graphs with structural nested models. **Computation, Causation, and Discovery**. Eds. C. Glymour and G. Cooper. Cambridge, MA: The MIT Press. Forthcoming.
- Robins, J.M. 1998b. Structural nested failure time models. **Survival Analysis**, P.K. Andersen and N. Keiding, Section Editors. **The Encyclopedia of Biostatistics**, P. Armitage and T. Colton, Editors. Chichester, UK: John Wiley & Sons. pp. 4372-4389.
- Robins, J.M. 1998c. Correction for non-compliance in equivalence trials. *Statistics in Medicine*, 17: 269-302.
- Robins, J.M. 1998d. Marginal structural models. *1997 Proceedings of the American Statistical Association, Section on Bayesian Statistical Science*, pp. 1-10.
- Robins, J.M., and Greenland S. 1989. The probability of causation under a stochastic model for individual risk. *Biometrics*, 45: 1125-1138.
- Robins, J.M., and Greenland S. 1994. Adjusting for differential rates of PCP prophylaxis in high-versus low-dose AZT treatment arms in an AIDS randomized trial. *Journal of the American Statistical Association*, 89: 737-749.
- Robins, J.M., Rotnitzky, A., Zhao LP. 1994. Estimation of regression coefficients when some regressors are not always observed. *Journal of the American Statistical Association*, 89: 846-866.
- Robins, J.M. and Ritov, Y. 1997. A curse of dimensionality appropriate (CODA) asymptotic theory for semiparametric models. *Statistics in Medicine*, 16: 285-319.
- Robins, J.M., Rotnitzky, A., and Scharfstein, D.O. 1999. Sensitivity analysis for selection bias and unmeasured confounding in missing data and causal inference models. In: **Statistical Models in Epidemiology**. Halloran E., Editor. Springer-Verlag. Forthcoming.

Robins, J.M., Rotnitzky, A., and Zhao LP. 1994. Estimation of regression coefficients when some regressors are not always observed. *Journal of the American Statistical Association*, 89: 846-866.

Robins, J.M., and Tsiatis A.A. 1992. Semiparametric estimation of an accelerated failure time model with time-dependent covariates. *Biometrika*, 79: 311-319.

Robins, J.M. and Wasserman L. 1997. Estimation of Effects of Sequential Treatments by Reparameterizing Directed Acyclic Graphs. *Proceedings of the Thirteenth Conference on Uncertainty in Artificial Intelligence, Providence Rhode Island, August 1-3, 1997*. Dan Geiger and Prakash Shenoy (Eds.), Morgan Kaufmann, San Francisco. pp. 409-420.

Rosenbaum, P.R. and Rubin, D.B. 1983. The Central Role of the Propensity Score in Observational Studies for Causal Effects. *Biometrika*, 70: 41-55.

Rubin, D.B. 1976. Inference and Missing Data. *Biometrika*, 63: 581-592.

Rubin, D.B. 1978. Bayesian Inference for Causal Effects: The Role of Randomization. *The Annals of Statistics*, 6: 34-58.

Sasieni, P. 1992. Information bounds for the conditional hazard ratio in a nested family of regression models. *Journal of the Royal Statistical Society, Series B*, 54: 617-635.

P. Spirtes, C. Glymour, R. Scheines. 1993. **Causation, Prediction, and Search. Lecture Notes in Statistics 81**. New York: Springer-Verlag.

van der Vaart, A.W. 1991. On differentiable functionals. *Annals of Statistics*, 19: 178-204.