

# Marginal Structural Models

James M. Robins, Harvard School of Public Health  
 677 Huntington Avenue, Boston, MA 02115, USA

*Key Words: causal inference; semiparametric models; longitudinal data.*

## 1. Introduction

Robins (1993, 1994, 1997, 1998) has developed a set of causal or counterfactual models, the structural nested models (SNMs). The purpose of this paper is to introduce a simpler 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 assumption. We then compare the strengths and weaknesses of MSMs versus SNMs for causal inference from complex longitudinal data. Our results provide an extension to continuous treatments of the semiparametric efficient propensity score estimator of an average treatment effect.

## 2. The Data

Consider a study where we observe  $n$  iid 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  $\bar{L}_{\bar{a}(u), \mathbf{0}}$  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:

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

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

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

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  $\bar{a} \in \bar{\mathcal{A}}$ ,

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

where for any variable  $\underline{Z}(t) = \{Z(u); u \geq t\}$ . 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

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

Here, the intensity process  $\lambda_A(t \mid \cdot)$  is  $\lim_{\delta t \rightarrow 0} pr[A(t + \delta t) \neq A(t^-) \mid A(t^-), \cdot] / \delta t$ .

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

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

We assume the  $\{L_{\bar{a}}(u); \bar{a} \in \bar{\mathcal{A}}^\dagger \subseteq \bar{\mathcal{A}}\}$  have a non-degenerate joint distribution whenever  $\bar{a}_1(u) \neq \bar{a}_2(u)$  for all  $\bar{a}_1, \bar{a}_2 \in \bar{\mathcal{A}}^\dagger$ . Then CAR implies sequential randomization (4) but the converse is not true. 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, if (7) is the sole restriction imposed, this 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.

## 2.1. 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.

**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:**  $E[Y_{\bar{a}} \mid V^\dagger] = g[\bar{a}(K), V^\dagger, \beta_0]$  where  $g(\cdot, \cdot, \cdot)$  is a known function.

**Model 1b – semiparametric regression:**  $\eta\{E[Y_{\bar{a}} \mid V^\dagger]\} = g[\bar{a}(K), V^\dagger, \beta_0] + g^\dagger(V^\dagger)$  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  $V^\dagger$ .”

**Model 1c – stratified transformation model:**  $pr [R(\bar{\alpha}, \beta_0) < t | V^\dagger] = F_0(t | V^\dagger)$ ,  $F_0(t | V^\dagger)$  an unknown distribution function,  $R(\bar{\alpha}, \beta) = r(Y_{\bar{\alpha}}, \bar{\alpha}, V^\dagger, \beta)$  is a known increasing function of  $Y_{\bar{\alpha}}$  satisfying  $r(y, \bar{\alpha}, V^\dagger, \beta) = y$  if  $\bar{\alpha} \equiv \mathbf{0}$  or  $\beta = 0$ .

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

**Model 1d – multivariate non-linear least squares:**  $E [Y_{\bar{\alpha}}(m) | V^\dagger] = g_m[\bar{\alpha}(m-1), V^\dagger, \beta_0]$ ,  $m = 1, \dots, K+1$  where the  $g_m(\cdot, \cdot, \cdot)$  are known.

**Model 2:**  $C = \infty$ ,  $Y_{\bar{\alpha}}$  is a failure time process, i.e.,  $Y_{\bar{\alpha}}$  jumps from 0 to 1 at some particular time and stays at 1. Then define the failure time  $T_{\bar{\alpha}}$  by the equation  $Y_{\bar{\alpha}}(T_{\bar{\alpha}}) = 1$  and  $Y_{\bar{\alpha}}(T_{\bar{\alpha}}^-) = 0$ . 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:**  $\lambda_{T_{\bar{\alpha}}}[t | V^\dagger] = \lambda_0(t) \exp[r(\bar{\alpha}(t^-), t, V^\dagger; \beta_0)]$  where  $r(\cdot)$  is a known function satisfying  $r(\mathbf{0}, t, 0; \beta) = 0$ .

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

**Model 2c – stratified time-dependent accelerated failure time model:**  $pr [r(T_{\bar{\alpha}}, \bar{\alpha}, V^\dagger, \beta_0) < t | V^\dagger] = F_0(t | V^\dagger)$  where  $r(u, \bar{\alpha}, V^\dagger, \beta) = r(u, \bar{\alpha}(u), V^\dagger, \beta)$  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

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

for  $\bar{\alpha} \in \bar{\mathcal{A}}$ , where  $R(\bar{\alpha}, \beta) = r(T_{\bar{\alpha}}, \bar{\alpha}, V^\dagger, \beta)$ .

### 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  $\bar{\mathcal{A}}$  is an ancillary (i.e., exogenous) covariate process, i.e.,

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

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

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

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  $\bar{Y}_{\bar{\alpha}}$  given  $V^\dagger$  specified by the MSM hold for the conditional distribution of the observable  $\bar{Y}(C)$  conditional on  $(\bar{\mathcal{A}}(C), V^\dagger)$ .

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  $\bar{\mathcal{A}}(C)$  as  $\bar{\mathcal{A}}$ . Thus, for our MSM models (1a) – (2c), (8) implies

**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:**  $R(\beta_0) \prod [^* \bar{A} | V^\dagger$  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

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

Semiparametric inference in models 1a - 2c without (9) and (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  $\{\widehat{\beta}^*(h, \phi)\}$  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  $h \in \mathcal{H}, \phi \in \Phi$ . 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. 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,  $\Lambda^\perp = \{U^*(\beta_0, h, \phi)\}$  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. 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  $0, 1, 2, \dots$ ,  $U_{tp}^*(\phi) = \sum_{k=0}^{int(C)} \phi(k, \bar{A}(k), \bar{L}(k^-)) - E^*[\phi(k, \bar{A}(k), \bar{L}(k^-)) | \bar{L}(k^-), \bar{A}(k^-)]$  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  $(\bar{A}(k), \bar{L}(k^-))$  with mean zero given  $(\bar{A}(k^-), \bar{L}(k^-))$ . If  $A(t)$  is a continuous time marked point process, then  $U_{tp}^*(\phi) = \int dM_A^*(u) \phi_1(u, \bar{A}(u^-), \bar{L}(u^-)) + \int dN_A(u) \{\phi_2(u, \bar{A}(u), \bar{L}(u^-)) - E^*[\phi_2(u, \bar{A}(u), \bar{L}(u^-)) | A(u) \neq A(u^-), \bar{L}(u^-), \bar{A}(u^-)]\}$  where  $dM_A^*(u) = dN_A(u) - \lambda_A^*[u | \bar{A}(u^-), \bar{L}(u^-)] du$  and  $dN_A(u) = I\{A(u) \neq A(u^-)\}$  counts jumps in the  $\bar{A}$  process.

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 models 1a - 2c with neither (9) nor (10) imposed.

**Model 1a:**  $\widehat{D}_{sm}^*(\beta, h) = U_{sm}^*(\beta, h) = h(\bar{A}, V^\dagger) \varepsilon(\beta)$  with  $\varepsilon(\beta) = Y - g(\bar{A}, V^\dagger, \beta)$  and  $h(\bar{A}, V^\dagger)$  is any  $\dim(\beta)$  vector function.

**Model 1b:**  $\eta(x) = x : \widehat{D}_{sm}^*(\beta, h) = U_{sm}^*(\beta, h) = \{\varepsilon(\beta) - h_1(\bar{A}, V^\dagger)\}$

$\{h_2(\bar{A}, V^\dagger) - E^*[h_2(\bar{A}, V^\dagger) | V^\dagger]\}$  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  $\widehat{D}_{sm}^*(\beta, h) \equiv U_{sm}^\dagger(h, \widehat{P}(\beta))$ , where  $P(\beta) = \text{expit}[g(\bar{A}, V^\dagger, \beta) + g^\dagger(V^\dagger)]$ ,  $\widehat{P}(\beta) = \text{expit}[g(\bar{A}, V^\dagger, \beta) + \widehat{g}^\dagger(V^\dagger)]$ ,  $\text{expit}(x) = e^x/(1+e^x)$ ,  $\widehat{g}^\dagger(V^\dagger)$  is a  $n^{\frac{1}{4}}$ -consistent estimate of  $g^\dagger(V^\dagger)$ , and  $U^\dagger(h, P(\beta)) \equiv \{Y - P(\beta)\} \{h(\bar{A}, V^\dagger) - E^*[h(\bar{A}, V^\dagger) P(\beta) \{1 - P(\beta)\} | V^\dagger] / E^*[P(\beta) \{1 - P(\beta)\} | V^\dagger]\}$ .

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

**Model 1d:** Let  $\varepsilon(\beta) = \{\varepsilon_1(\beta), \dots, \varepsilon_{K+1}(\beta)\}'$ ,  $\varepsilon_m(\beta) = Y(m) - g_m[\bar{A}(m-1), V^\dagger; \beta]$ . 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:**  $\widehat{D}_{sm}^*(\beta, h) = \int_0^\infty dN_T(u) \{h(u, \bar{A}(u), V^\dagger) - \tilde{\mathcal{L}}(h, u, \beta)\}$ , 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  $N_T(u) = I(T \leq u)$ .  $U_{sm}^*(\beta, h) = \int_0^\infty dM_T(u) \{h(u, \bar{A}(u), V^\dagger) - \mathcal{L}^*(h, u, V^\dagger, \beta)\}$  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  $\widehat{D}_{sm}^*(\beta, h)$  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  $\tilde{E}(\cdot | V^\dagger)$ .

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

**Remark:** Note that in model 2b and in model 1b with  $\eta(x) = \ln[x/(1-x)]$ , smooths are necessary 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.

### 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  $\widehat{D}_{sm}^*(\beta, h)$  by the inverse of a subject's probability of having had his observed treatment history. We continue to assume that

$$f[a(t) | \bar{L}(t^-), \bar{A}(t^-), V^\dagger]$$

is known for  $t \leq C$  (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,  $\bar{W}(t) = \exp\left[-\int_0^t \lambda_A[u | \bar{L}(u^-), \bar{A}(u^-)] du\right] \prod_{\{u: A(u) \neq A(u^-), u < t\}} f[A(u) | \bar{A}(u^-), \bar{L}(u^-), A(u) \neq A(u^-)]$  is known.

Now if  $A(t)$  jumps at non-random times, let  $\overset{\circ}{\mathcal{A}}(k, \bar{a}(k^-), v^\dagger) = \{a(k); f[a(k) | \bar{L}(k^-), \bar{A}(k^-) = \bar{a}(k^-), V^\dagger = v^\dagger] \neq 0 \text{ w.p.1}\}$  and set  $C^\dagger = \min\{k; A(k) \notin \overset{\circ}{\mathcal{A}}(k, \bar{A}(k^-), V^\dagger)\}$ .

If  $A(t)$  jumps in continuous time, let  $\overset{\circ}{\mathcal{A}}(t, \bar{a}(t^-), v^\dagger) = \{a(t); f[a(t) | \bar{L}(t^-), \bar{A}(t^-) = \bar{a}(t^-), A(t) \neq a(t^-), V^\dagger = v^\dagger] \neq 0 \text{ w.p.1} \text{ or } a(t) = a(t^-)\}$  and set  $C^\dagger = \inf\{t; A(t) \notin \overset{\circ}{\mathcal{A}}(t, \bar{A}(t^-), V^\dagger)\}$ . 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  $t \in \overset{\circ}{\mathcal{A}}(t, v^\dagger)$ .

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  $O^\dagger = (\bar{L}(C^\dagger), \bar{A}(C^\dagger))$ . 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.

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 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,  $\{\hat{\beta}(h, \phi)\}$  solving  $0 = \sum_i \hat{D}_i(\beta, h, \phi)$  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 $h$

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

**Theorem 3.2: a)** The semiparametric model (ii) characterized by (4), data  $O$ , and a MSM (with (11) not imposed), has the set of influence functions  $\{IF(h, \phi_{opt}(h))\}$ .

**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

$\{\kappa(h)^{-1} [U(h, \phi) - E[U(h, \phi) S'_\alpha] \{E[S_\alpha S'_\alpha]\}^{-1} S_\alpha]\}$  of influence functions of  $\{\hat{\beta}(h, \phi, \hat{\alpha})\}$  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  $\alpha$ ,  $S_\alpha$  is the subject-specific score for  $\alpha$ , and  $D(\beta, h, \phi, \hat{\alpha})$

is  $D(\beta, h, \phi)$  evaluated at  $\hat{\alpha}$ .

Theorem 3.2b 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).

## 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  $U(\beta_0, h, \phi)$

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

which is equal to

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

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  $V^\dagger \subset L(0) = V(0)$ . As the proof of the following theorem, available from the author, shows, we can let  $f^*(\bar{A} | V^\dagger) = 1$  w.p.1 without worrying that it is not a density, because it can be absorbed into  $h_{eff}(\bar{A}, V^\dagger)$ .

**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  $h_{eff}(\bar{A}, V^\dagger) [f \text{var}[\varepsilon | \bar{A}, L(0)] \{f(\bar{A} | L(0))\}^{-1} f(V^\bullet | V^\dagger) d\mu(V^\bullet)] + \int h_{eff}(\bar{a}, V^\dagger) \omega(\bar{a}, \bar{A}, V^\dagger) d\mu(\bar{a}) = \partial g(\bar{A}, V^\dagger, \beta_0) / \partial \beta$  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.

### 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  $\bar{a}_1, \bar{a}_2$

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

then

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

**Remark:** Since (14) places no restriction on the law of the observed data  $O$ , 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  $g[\mathbf{m}^{-1}(D_{eff})]$  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),

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



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,  $f(\bar{a} | \bar{L}_{\bar{a}}) = f(\bar{a} | L(0))$ , the LHS of (21) can be written

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

However, since by assumption (14),  $Y_{\bar{a}_j} \perp\!\!\!\perp Y_{\bar{a}_k} | \bar{L}(0)$  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,  $f[\varepsilon | \bar{A} = \bar{a}, L(0)] = f[\varepsilon(\bar{a}) | L(0)]$ ,  $w$ . 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. [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. When (8) is false, any variation in  $\mathcal{W}$  will be unavoidable variation. 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

### 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. 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  $\bar{L}(m) = \bar{\ell}(m), \bar{A}(m) = \bar{a}(m)$ . A structural nested distribution model (SNDM) 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  $\dot{R}(\beta) = \dot{r}(Y, \bar{L}(K), \bar{A}(K), \beta) \equiv \dot{R}_0(\beta)$ . 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  $\bar{\ell}(m)$ , 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.

The semiparametric information bound for  $\beta$  is greater if we impose a no-interaction SNDM than if we only imposed the corresponding STM. Theorem 5.1 indicates that a STM is the natural MSM analog of a SNDM. 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

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

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

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

which, for continuous  $Y$  and  $k = -1$ , equals

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

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. Robins (1997) proves the following.

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

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

## 5.2. Advantages of SNDMs with a continuous $Y$

1. Although (25) implies  $\beta_0 = 0$  for both a SNDM and a STM, only for a SNDM is (25) equivalent to  $\beta_0 = 0$ .

2. If the  $L(k)$  are discrete with only a moderate number of levels, then, even with  $f[a(k) \mid \bar{\ell}(k-1), \bar{a}(k-1)]$

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 lack of correlation of  $Y$  with  $A(k)$  within strata defined jointly by  $\bar{L}(k), \bar{A}(k-1)$  (Robins, 1997). Thus, even if  $A(k)$  is continuous, a test of lack of correlation of  $A(0)$  with  $Y$  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 Theorem 3.2a that  $\mathcal{W}$  can be consistently non-parametrically estimated which will not be possible due to the curse of dimensionality.

3. Henceforth, assume a correct model for  $f[a(t) | \bar{\ell}(t^-), \bar{a}(t^-)]$  is available. 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  $\dot{R}_i(\hat{\beta})$ , (iii) specify and estimate a parametric model for  $f[L(m) | \bar{L}(m-1), \bar{A}(m-1), \dot{R}(\hat{\beta})]$ , (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; 1998, Sec. 11), specify a parametric model for  $\nu^*(y, \bar{\ell}(m), \bar{a})$  where  $\nu^*(y, \bar{\ell}(m), \bar{a}) = \nu(y, \bar{\ell}(m), \bar{a}) - \nu(y, \{\bar{\ell}(m-1), \ell(m) = 0\}, \bar{a})$  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  $\bar{\ell}_{m-1}, \bar{a}_{m-1}$ . As discussed in Robins (1997, pg. 114-116; 1998, Sec. 11), 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), 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 much less straightforward and cannot be based on estimating equations similar to those in Theorem 3.1.

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).

### 5.3. Advantages of MSMs with continuous $Y$

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,  $F_{Y_{\bar{a}}}(y)$  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  $Y_{\bar{a}(m-1),0}$ . 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.

A hybrid approach is to impose a MSM model and then specify a model for  $\nu^*(y, \bar{\ell}(m), \bar{a})$  whose

parameter is the product of the parameter  $\beta$  of the MSM model with another parameter  $\psi$ , so that  $\psi$  is identified only if  $\beta \neq 0$ . Such a model can overcome objections 1 and 2 of Sec. 5.2 (but not objections 3-5) while retaining the advantages 1-2 of Sec. 5.3.

#### 5.4. Structural Nested Mean Models (SNMMs)

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  $\gamma^\dagger(\bar{\ell}(m), \bar{a}(m)) = \ln\{E[Y_{\bar{a}(m),0} | \bar{\ell}(m), \bar{a}(m)] / E[Y_{\bar{a}(m-1),0} | \bar{\ell}(m), \bar{a}(m)]\}$ . A 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) 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  $E[Y_g]$ . Advantage (5) is no longer so clear, since estimation of a SNMM when (22) is false is not as straightforward as for a SNDM. Advantages (1) - (2) in Sec. 5.3 of a MSM over a SNDM 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 a SNMM or multiplicative SNMM naturally imposes the fact that, for dichotomous  $Y$ ,  $E[Y_g] \in [0, 1]$  (and logistic SNMMs have not been developed). 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  $Y_{\bar{a}}(m)$ .

#### 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 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, 1998) have important limitations due to the curse of dimensionality. 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,  $E[Y_{\bar{a}_1, \bar{a}_2} - Y_{\bar{a}_1 \equiv 0, \bar{a}_2}] = g(\bar{a}, \beta_0)$ , it follows we are modelling the direct effect of treatment  $\bar{a}_1$ . Furthermore, the main effect of the second treatment  $g^\dagger(\bar{a}_2) = E[Y_{\bar{a}_1 \equiv 0, \bar{a}_2} - Y_{\bar{a}_1 \equiv 0, \bar{a}_2 \equiv 0}]$  is completely unrestricted. 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  $V^\dagger$  (under sequential randomization assumption (4)). 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  $\overline{A}_2(K)$  will be high-dimensional. Indeed, any choice of  $\eta(x)$  that guarantees that  $E[Y_{\overline{a}}] \in [0, 1]$  will fail to provide reasonable estimators of  $\beta_0$ , negating the advantage of this MSM for dichotomous  $Y$ .

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

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

with  $r(\overline{a}(t^-), t; \beta_0) = 0$  if  $\overline{a}_1(t^-) = \mathbf{0}$ . This is a model for the direct effect of treatment  $\overline{a}_1$  on the hazard of  $T$  with the main effect of  $\overline{a}_2$  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:** Consider the model

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

where  $R(\overline{a}, \beta) = r(T_{\overline{a}}, \overline{a}, \beta)$  satisfies  $r(t, \overline{a}, \beta) = t$  if  $\overline{a}_1 = \mathbf{0}$  or  $\beta = 0$ . The model for the observables induced by MSM 3c is isomorphic to that induced by MSM 2c with  $\overline{A}_2$  in the role of  $V^\dagger$ . Hence, model 3c can be used to estimate the direct effect of  $\overline{a}_1$  on  $T$  with the main effect of  $\overline{a}_2$  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.

## 6. Censoring

No new idea is required to account for and adjust for right censoring. Specifically, let  $Q$  be censoring time. 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  $a(u) = (a_1(u), a_2(u))$ . To want to adjust for censoring is only to say that interest is in the direct effect of  $\overline{a}_1$  when  $\overline{a}_2 \equiv 0$ , i.e., when censoring is abolished. As a concrete example, the Cox model MSM 2a in the presence of censoring would become

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

If  $\overline{A}$  is ancillary (now including the censoring process),  $\widehat{D}_{sm}^*(\beta, h)$  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  $\overline{W}(t)$  is now the probability that a subject would have his observed treatment and censoring history.

### Appendix 1:

By arguments as in Robins et al. (1994), Theorem 3.1 and 3.2b are easy corollaries of Theorem 3.2a. **Sketch of Proof of Theorem 3.2a:** 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

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

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  $\overset{\circ}{\Lambda} = \overset{\circ}{\Lambda}^*$  where  $\overset{\circ}{\Lambda} \equiv \Lambda \cap \{U_{tp}(\phi)\}^\perp \cap \{z(O^\dagger)\}$  and the fact that  $E[UB] = E^*[UWB]$  for any  $B \in \overset{\circ}{\Lambda}$  by Lemma 3.1.

## References

- Chamberlain, G.** (1987), Asymptotic Efficiency in Estimation with Conditional Moment Restrictions, *J. 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. (1996). Coarsening at random: characterizations, conjectures and counterexamples. **Proceedings of the First Seattle Symposium on Survival Analysis**, to appear.
- Heitjan, D.F.,** & Rubin, D.B. (1991), Ignorability and Coarse Data, **The Annals of Statistics**, 19, 2244-2253.
- Newey, W.K.** & McFadden, D. (1993). Estimation in large samples. **Handbook of Econometrics**, Vol. 4. Eds. McFadden, D., Engler, R. Amsterdam: North Holland.
- Robins, J.M.** (1993). Analytic methods for estimating HIV treatment and cofactor effects. **Methodological Issues of AIDS Mental Health Research**. Eds: Ostrow D.G., Kessler R. New York: Plenum Publishing. pp. 213-290.
- Robins, J.M.** (1994). Correcting for non-compliance in randomized trials using structural nested mean models. *Comm. Stat.*, 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, 69-117.
- Robins, J.M.,** Rotnitzky, A., Zhao LP. (1994). Estimation of regression coefficients when some regressors are not always observed. *JASA*, 89:846-866.
- Robins, J.M.** (1998). Testing and estimation of direct effects by reparameterizing directed acyclic graphs with structural nested models. **Causation and Computation**, Ed. G. Cooper and C. Glymour, MIT/AAAI Press (to appear).
- Sasieni, P.** (1992). Information bounds for the conditional hazard ratio in a nested family of regression models. *J. Royal Stat. Soc. B*, 54:617-635.
- van der Vaart, A.W.** (1991). On differentiable functionals. *Ann. Stat.*, 19:178-204.