

Outline

Chapter 5

April 8, 2022

- 1 Semiparametric models
- 2 Causal effects of a point exposure

1

2

Outline

- 1 Semiparametric models
- 2 Causal effects of a point exposure

- In Chapter 3, we developed theoretical results for influence functions of parameters in finite-dimensional parametric models $\{f(X; \theta) : \theta \in \mathbb{R}^p\}$, p finite and where θ can be partitioned as

$$\theta = (\beta^T, \eta^T)^T : \beta \in \mathbb{R}^s, \eta \in \mathbb{R}^r, p = r + s$$

with β the parameter of interest and η the nuisance parameter.

- In chapter 4 we considered influence functions of functionals $\psi(\theta) = \psi(F_\theta)$, where θ is an infinite dimensional parameter ranging in some functional space Θ , i.e. $\{f(X; \theta) : \theta \in \Theta\}$ is a nonparametric model
- In this chapter, we extend this theory to semiparametric models where the parameter $\theta = (\beta, \eta)$ can be partitioned into a finite parameter of interest β and an infinite dimensional nuisance parameter η .

3

4

- In this chapter we will consider semiparametric models that can be represented using the class of densities $\mathcal{M}_{F_{\theta^*}} = \{f(X; \theta) : \beta \in \mathbb{R}^s, \eta \in \Delta\}$, where β the parameter of interest and the infinite dimensional parameter η are variationally independent.
- That is, any choice of β and η in the neighborhood about the true β^* and η^* would result in a density $f(X; \theta)$ in the semiparametric model.
- This will allow us to define $\partial f(X; \theta) / \partial \beta|_{\theta^*}$.
- However keep in mind that some pbs lend themselves more naturally to models represented by the class of densities $f(X; \theta)$, where θ is infinite dimensional and $\psi(\theta)$ is an explicit functional of θ . We will make the distinction explicit when necessary.

5

- Example 1: Restricted moment model (See Chapter 1)

$$\begin{aligned} Y_i &= \mu(X_i, \beta) + \varepsilon_i, \\ \mathbb{E}(\varepsilon_i | X_i) &= 0 \end{aligned}$$

$\mu(X_i, \beta)$ is known up to β , density of ε_i , η_1 has conditional mean zero given X_i and otherwise unrestricted and density η_2 of X_i is unrestricted. Thus $\theta = (\beta, \eta_1, \eta_2) \in \mathbb{R}^s \times \Gamma_1 \times \Gamma_2$.

- Example 2: Cox proportional hazards model (see Chapter 1)

$$\lambda(t|Z; \beta) = \lambda_0(t) \exp(\beta' Z)$$

where $\lambda(t|Z; \beta)$ is the conditional hazard function of T at t given Z , let η_1 index the density of Z which is unrestricted, then:

$$\theta = (\beta, \lambda_0, \eta_1) \in \mathbb{R}^p \times \Gamma \times \Gamma_1,$$

β is the regression parameter of interest (log hazards ratio), while $\lambda_0(t)$ is an unrestricted function of t .

6

Example 3: Semilinear or semiparametric regression model (see Chapter 1)

$$\mathbb{E}(Y|X, Z) = \beta^T X + \eta(Z),$$

β is the finite dimensional parameter of interest and the model is otherwise unrestricted, i.e. $\eta(Z)$ is unrestricted and joint density of X and Z is unrestricted.

- To compute influence functions in a semiparametric model, we define a regular parametric submodel

$$\mathcal{M}_{\text{sub}, F_{\theta^*}} = \left\{ f(X; \theta = (\beta^T, \eta^T)^T) : \beta \in \mathbb{R}^s, \eta \in \mathbb{R}^r \right\}$$

such that

- $\mathcal{M}_{\text{sub}, F_{\theta^*}} \subset \mathcal{M}_{F_{\theta^*}} = \{f(X; \theta = (\beta, \eta)) : \beta \in \mathbb{R}^s, \eta \in \Delta\}$ and
- F_{θ^*} matches $f(X; \bar{\theta} = (\bar{\beta}, \bar{\eta}))$ for some $(\bar{\beta}, \bar{\eta})$ within $\mathbb{R}^s \times \mathbb{R}^r$
- As in Chapter 4, throughout we will assume the same regularity or smoothness conditions as where needed there for pathwise differentiation.

7

8

- One clarification: the term parametric *submodel* and parametric *model* can be confusing.
- A parametric model is a model whose probability densities are characterized through a finite number of parameters that the analyst believes will suffice to identify the unknown data generating mechanism
- In contrast, a parametric submodel is a conceptual tool used to develop theory for semiparametric models. Such parametric models cannot be used for data analysis as they require knowing what the true data generating mechanism is.

9

- Consider Example 2, a parametric submodel of the Cox proportional model may be taken to be: let $h_1(t) \dots h_r(t)$ denote r different functions of t specified by the analyst

$$\mathcal{M}_{\text{sub}, F_{\theta^*}} = \left\{ \begin{array}{l} \text{class of densities with hazard function:} \\ \lambda_0^*(t) \exp(\eta_1 h_1(t) + \dots + \eta_r h_r(t)) \exp(\beta' Z) : \\ \beta \in \mathbb{R}^s, \eta \in \mathbb{R}^r \end{array} \right\}$$

- Because $(\beta^T, \eta^T)^T$ are unspecified, this model is indeed finite dimensional
- For any choice of (β^T, η^T) , the resulting density follows a proportional hazards model and is therefore contained in the semiparametric model
- the truth is obtained at $(\beta^T, \eta^T) = (\beta^{*T}, 0^T)$
- The parametric submodel is "anchored" at the true model which is of course unknown and therefore such a parametric model is not useful for data analysis.

10

- In contrast, suppose one is willing to consider the parametric model

$$\lambda(t|Z; \beta) = \lambda_0 \exp(\beta' Z), \lambda_0 \text{ and } \beta \text{ unknown}$$

that is the baseline hazard is constant over time conditional on Z , i.e. the failure time is a conditional exponential model. This is a model which satisfies the proportional hazards assumption and may be used for data analysis, however if the model is incorrect, that is the true data generating model is not exponential, the resulting inferences may be completely incorrect and possibly meaningless.

11

Influence Functions for Semiparametric Estimators

- The nuisance tangent space for a semiparametric model denoted by Λ_{nuis} is defined as the mean-square closure of parametric submodel nuisance tangent spaces, where a parametric submodel nuisance tangent space is the set of elements

$$\Lambda_{\text{sub}, \text{nuis}} = \{A^{s \times r} S_{\eta}(X; \theta^*)\}$$

$S_{\eta}(X; \theta^*) = \partial \log f(X; \beta^*, \eta) / \partial \eta|_{\eta^*}$ is the score vector for the nuisance parameter η for some parametric submodel assuming β^* is known, and $A^{s \times r}$ is a conformable matrix with s rows.

- Specifically the mean-square closure of the space above is defined as the space of functions $\Lambda_{\text{nuis}} = \{h(X) \text{ such that } \mathbb{E}(h^T(X)h(X)) < \infty\}$ and there exists a sequence of parametric submodels indexed by j such that

$$\|h(X) - A_j^{s \times r} S_{\eta, j}(X; \theta^*)\|^2 \xrightarrow{j \rightarrow \infty} 0.$$

12

- Recall that the efficient influence function of β for a parametric sub-model $\mathcal{M}_{\text{sub}, F_{\theta^*}}$ with finite dimensional nuisance parameter η is given by

$$\varphi_{\beta}^{\text{eff}}(X) = \{\mathbb{E} S_{\beta}^{\text{eff}} S_{\beta}^{\text{eff}T}\}^{-1} S_{\beta}^{\text{eff}}(X; \beta^*, \eta^*),$$

where

$$S_{\beta}^{\text{eff}}(X; \beta^*, \eta^*) = S_{\beta}(X; \beta^*, \eta^*) - \Pi(S_{\beta} | \Lambda_{\text{sub}, \text{nuis}}),$$

the CR efficiency bound for estimating β in the submodel is

$$\{\mathbb{E} S_{\beta}^{\text{eff}} S_{\beta}^{\text{eff}T}\}^{-1}.$$

13

- The variance of any RAL semiparametric influence function must be greater or equal to

$$\{\mathbb{E} S_{\beta, \eta}^{\text{eff}} S_{\beta, \eta}^{\text{eff}T}\}^{-1}$$

for all parametric submodels.

- Hence the variance of any semiparametric influence function must be greater than or equal to

$$\mathcal{V} = \sup_{\mathcal{M}_{\text{sub}, F_{\theta^*}}} \{\mathbb{E} S_{\beta, \eta}^{\text{eff}} S_{\beta, \eta}^{\text{eff}T}\}^{-1}.$$

- The supremum \mathcal{V} is defined to be semiparametric efficiency bound. Any semiparametric estimator $\hat{\beta}_n$ with asymptotic variance achieving this bound for $f^*(X) = f(X; \beta^*, \eta^*)$ is said to be semiparametric locally efficient at f^* . If the same estimator $\hat{\beta}_n$ is semiparametric efficient regardless of f^* , then we say such estimator is globally semiparametric efficient.

14

Theorem 5.1. The Semiparametric efficiency bound \mathcal{V} is equal to

$$\mathcal{V} = \{\mathbb{E} S_{\text{eff}} S_{\text{eff}}^T\}^{-1}$$

where

$$S_{\text{eff}}(X; \beta^*, \eta^*) = S_{\beta}(X; \beta^*, \eta^*) - \Pi(S_{\beta} | \Lambda_{\text{nuis}})$$

is defined as the semiparametric efficient score for β .

Proof: For simplicity, take β to be a scalar parameter although this can easily be extended to $s > 1$. Since $\mathcal{M}_{\text{sub}, F_{\theta^*}} \subset \mathcal{M}_{F_{\theta^*}}$, this implies that

$$\begin{aligned} & \|S_{\beta}(X; \beta^*, \eta^*) - \Pi(S_{\beta} | \Lambda_{\text{nuis}})\| \\ & \leq \|S_{\beta}(X; \beta^*, \eta^*) - \Pi(S_{\beta} | \Lambda_{\text{sub}, \text{nuis}})\| \end{aligned}$$

or

$$\|S_{\text{eff}}(X; \beta^*, \eta^*)\| \leq \|S_{\beta, \eta}^{\text{eff}}\|$$

for all $\mathcal{M}_{\text{nuis}, F_{\theta^*}}$. Hence

$$\|S_{\text{eff}}(X; \beta^*, \eta^*)\|^{-2} \geq \sup_{\mathcal{M}_{\text{nuis}, F_{\theta^*}}} \|S_{\beta, \eta}^{\text{eff}}\|^{-2} = \mathcal{V}$$

To complete the proof requires also showing that

$$\|S_{\text{eff}}(X; \beta^*, \eta^*)\|^{-2} \leq \mathcal{V}$$

15

16

Because $\Pi(S_\beta(X)|\Lambda_{\text{nuis}}) \in \Lambda_{\text{nuis}}$, this means that there exists a sequence of parametric submodels with nuisance scores S_{η_j} and vectors A_j such that

$$\|\Pi(S_\beta(X)|\Lambda_{\text{nuis}}) - A_j^T S_{\eta,j}(X; \theta^*)\|^2 \xrightarrow{j \rightarrow \infty} 0.$$

Therefore

$$\begin{aligned} \mathcal{V}^{-1} &\leq \|S_{\beta, \eta_j}^{\text{eff}}\|^2 \\ &= \|S_\beta(X; \beta^*, \eta^*) - \Pi(S_\beta| \Lambda_{\text{nuis}, \text{sub}, j})\|^2 \\ &\leq \|S_\beta(X; \beta^*, \eta^*) - A_j^T S_{\eta,j}(X; \theta^*)\|^2 \\ &= \|S_\beta(X; \beta^*, \eta^*) - \Pi(S_\beta(X)|\Lambda_{\text{nuis}})\|^2 \\ &\quad + \|\Pi(S_\beta(X)|\Lambda_{\text{nuis}}) - A_j^T S_{\eta,j}(X; \theta^*)\|^2 \end{aligned}$$

taking $j \rightarrow \infty$ implies

$$\begin{aligned} \mathcal{V}^{-1} &\leq \|S_\beta(X; \beta^*, \eta^*) - \Pi(S_\beta(X)|\Lambda_{\text{nuis}})\|^2 \\ &= \|S_{\text{eff}}(X; \beta^*, \eta^*)\|^2 \end{aligned}$$

proving the result.

Definition 5.1. The efficient influence function is defined as the influence function of a semiparametric RAL estimator, if it exists, that achieves the semiparametric efficiency bound.

Theorem 5.2. Any RAL estimator for β must have an influence function $\varphi(X)$ that satisfies

$$(i) \mathbb{E}\{\varphi(X) S_\beta^T(X)\} = \mathbb{E}\{\varphi(X) S_{\text{eff}}(X; \beta^*, \eta^*)^T\} = I$$

and (ii) $\Pi(\varphi(X)|\Lambda_{\text{nuis}}) = 0$, i.e. $\varphi(X)$ is orthogonal to the nuisance tangent space. The efficient influence function is equal to

$$\begin{aligned} \varphi_{\text{eff}}(X) &= \mathbb{E}\{S_{\text{eff}} S_{\text{eff}}^T\}^{-1} S_{\text{eff}}(X; \beta^*, \eta^*) \\ S_{\text{eff}}(X; \beta^*, \eta^*) &= S_\beta(X; \beta^*, \eta^*) - \Pi(S_\beta| \Lambda_{\text{nuis}}) \\ &= \Pi(S_\beta| \Lambda_{\text{nuis}}^\perp) \end{aligned}$$

Proof of (ii) Note that from previous chapter under any parametric submodel

$$\frac{\partial \beta(F_t)}{\partial t} = \mathbb{E}\{\varphi(X) S_t(X)^T\}$$

taking t in the direction of η implies that

$$0 = \mathbb{E}\{\varphi(X) S_{\eta_t}(X)^T\}$$

for any nuisance submodel, which implies that for any sequence of submodels $S_{\eta,j}(X; \theta^*)$ and matrices A_j ,

$$\|h(X) - A_j^T S_{\eta,j}(X; \theta^*)\|^2 \xrightarrow{j \rightarrow \infty} 0 \text{ for } h \in \Lambda_{\text{nuis}}$$

$$0 = \mathbb{E}\{\varphi(X) [h(X) - A_j^T S_{\eta_t}(X)^T]\} + \mathbb{E}\{\varphi(X) h^T(X)\}$$

$$|\mathbb{E}\{\varphi(X) [h(X) - A_j^T S_{\eta_t}(X)^T]\}| \leq \|\varphi\| \|h - A_j^T S_{\eta_t}\|$$

the result follows by letting $j \rightarrow \infty$.

Part (i) of the theorem follows from the fact that φ is the IF of a RAL estimator and (ii).

Theorem 5.3. If a semiparametric RAL estimator for β exists, then the influence function of this estimator must belong to the space of influence functions, the linear variety $\{\varphi(X) + \Lambda^\perp\}$ where $\varphi(X)$ is the influence function of any semiparametric RAL estimator of β and Λ is the semiparametric tangent space, i.e. the closed linear span of all scores for the model, including score of β . Furthermore, if a RAL estimator for β exists that achieves the semiparametric efficiency bound, then the IF of this estimator must be the unique element $\varphi_{\text{eff}}(X) = \varphi(X) - \Pi(\varphi(X)|\Lambda^\perp) = \Pi(\varphi(X)|\Lambda)$ and the space of IF may be written $\varphi_{\text{eff}}(X) \oplus \Lambda^\perp \subseteq \Lambda_{\text{nuis}}^\perp$.

In the rest of the course, we will use the space of influence functions or even the space orthogonal to the nuisance tangent space for semiparametric models to construct "good" estimating equations of functionals. "good" in the sense that they have small bias and in many instances will give RAL estimators that are locally or globally efficient in the semiparametric model. We illustrate this fact with the familiar restricted mean model next, linking it to the class of so called generalized estimating equations (GEE).

- The semiparametric restricted mean model:

$$\begin{aligned} Y &= m(Z; \beta) + \varepsilon \\ \mathbb{E}\{\varepsilon|Z\} &= 0. \end{aligned}$$

- Note that the likelihood for one observation can be written

$$\begin{aligned} f_{Y,Z}(y, z; \beta, \eta) &= f_{\varepsilon,Z}(y - m(Z; \beta), z; \beta, \eta) \\ &= \eta_1(\varepsilon, z) \eta_2(z) \end{aligned}$$

where $\eta_1(\varepsilon, z)$ is the conditional density of ε given z is only restricted by conditions

$$\begin{aligned} \int \eta_1(\varepsilon, z) d\mu(\varepsilon) &= 1 \\ \int \varepsilon \eta_1(\varepsilon, z) d\mu(\varepsilon) &= 0 \\ \int \eta_2(z) d\mu(z) &= 1 \end{aligned}$$

21

22

- Result 5.1: As we will show below, the nuisance tangent space of RMM is given by:

$$\begin{aligned} \Lambda_{\text{nuis}} &= \Lambda_{\text{nuis},1} \oplus \Lambda_{\text{nuis},2} \subset L_2^0 \\ \Lambda_{\text{nuis},1} &= \left\{ a_1(\varepsilon, Z) : \mathbb{E}\{a_1(\varepsilon, Z)|Z\} = 0 \right. \\ &\quad \left. \mathbb{E}\{\varepsilon a_1^T(\varepsilon, Z)|Z\} = 0 \right\} \cap L_2^0 \\ \Lambda_{\text{nuis},2} &= \{a_2(Z) : \mathbb{E}\{a_2(Z)\} = 0\} \cap L_2^0 \end{aligned}$$

23

- Result 5.2: The ortho-complement to the nuisance tangent space

$$\begin{aligned} \Lambda_{\text{nuis}}^\perp &= \Lambda_{\text{nuis},1}^\perp \cap \Lambda_{\text{nuis},2}^\perp \cap L_2^0 \\ &= \{h(Z) \varepsilon(\beta^*) : h\} \cap L_2^0 \end{aligned}$$

- Result 5.3: The set of influence functions is therefore given by

$$\left\{ \mathbb{E}\{h(Z) \varepsilon(\beta^*) S_\beta^T\}^{-1} h(Z) \varepsilon(\beta^*) : h \right\} \cap L_2^0$$

where S_β^T is the score of β under the RMM.

24

► Result 5.4: The efficient score for β is given by

$$\begin{aligned} S_{\text{eff}}(Z; \beta^*, \eta^*) &= \Pi(S_\beta | \Lambda_{\text{nuis}}^\perp) \\ &= h^{\text{opt}}(Z) \varepsilon(\beta) \\ h^{\text{opt}}(Z) &= \mathbb{E}\{S_\beta \varepsilon^T(\beta) | Z\} \mathbb{E}\{\varepsilon(\beta) \varepsilon^T(\beta) | Z\}^{-1} \end{aligned}$$

► Proof Result 5.1. Consider a regular parametric submodel $\eta_t(\varepsilon, z) = \eta_{1,t}(\varepsilon, z) \eta_{2,t}(z)$ with unknown parameter $t \in (-\epsilon, \epsilon)$ such that

$$\eta_{t=0}(\varepsilon, z) = f(\varepsilon, z).$$

It is straightforward to verify that corresponding score equations may be written

$$\frac{\partial \log \eta_t(\varepsilon, z)}{\partial t} = a_1(\varepsilon, z) + a_2(z)$$

where $\mathbb{E}\{a_1(\varepsilon, Z) | Z\} = \mathbb{E}\{a_2(Z)\} = 0$ and thus, $\mathbb{E}\{a_1(\varepsilon, Z) a_2(Z)\} = 0$. The first term a_1 is a score for the conditional density of $\varepsilon | Z$, while a_2 is a score for the density of Z .

25

26

Furthermore, we also have that for all $t \in (-\epsilon, \epsilon)$:

$$\int \eta_{1,t}(\varepsilon, z) \varepsilon^T d\mu(\varepsilon) = 0,$$

therefore

$$\int \frac{\partial \eta_{1,t}(\varepsilon, z)}{\partial t} \Big|_{t=0} \varepsilon^T d\mu(\varepsilon) = 0,$$

which in turn implies that $\mathbb{E}\{a_1(\varepsilon, Z) \varepsilon^T | Z\} = 0$ must hold for all scores $a_1(\varepsilon, Z)$ in the semiparametric model.

► From these results, we have that:

$$\begin{aligned} \Lambda_{\text{nuis}} &= \Lambda_{\text{nuis},1} \oplus \Lambda_{\text{nuis},2} \subset L_2^0 \\ \Lambda_{\text{nuis},1} &= \left\{ a_1(\varepsilon, Z) : \mathbb{E}\{a_1(\varepsilon, Z) | Z\} = 0 \right. \\ &\quad \left. \mathbb{E}\{\varepsilon a_1^T(\varepsilon, Z) | Z\} = 0 \right\} \cap L_2^0 \\ \Lambda_{\text{nuis},2} &= \{a_2(Z) : \mathbb{E}\{a_2(Z)\} = 0\} \cap L_2^0 \end{aligned}$$

► Technically, for each element of Λ_{nuis} we must also exhibit parametric submodels with corresponding score equal to the given element. I will leave it as hwk to work this out. See Tsiatis Section 4.5.

27

28

- Proof of Result 5.2. The ortho-complement to the nuisance tangent space

$$\begin{aligned}\Lambda_{\text{nuis}}^\perp &= (\Lambda_{\text{nuis},1} \oplus \Lambda_{\text{nuis},2})^\perp \\ &= \Lambda_{\text{nuis},1}^\perp \cap \Lambda_{\text{nuis},2}^\perp \cap L_2^0 \text{ (why?)}\end{aligned}$$

Note that any function $b(\varepsilon, Z) \in \Lambda_{\text{nuis},2}^\perp \cap L_2^0$ must satisfy

$$\mathbb{E}(b(Z, \varepsilon) a_2(Z)) = 0$$

for all $a_2(Z) \in \Lambda_{\text{nuis},2} \cap L_2^0$, that is $\mathbb{E}(b(Z, \varepsilon) | Z) = 0$. (why?)

- Therefore, to characterize $\Lambda_{\text{nuis}}^\perp$ entails finding all $b(\varepsilon, Z) \in \Lambda_{\text{nuis},2}^\perp$ such that $\mathbb{E}(b(X, \varepsilon) a_1(\varepsilon, Z)) = 0$ for all $a_1(\varepsilon, Z) \in \Lambda_{\text{nuis},1}^\perp$.
- Now recall that $a_1(\varepsilon, Z) \in \Lambda_{\text{nuis},1}^\perp$ implies $\mathbb{E}\{a_1(\varepsilon, Z) \varepsilon | Z = 0\}$; therefore, a natural candidate for $\Lambda_{\text{nuis},2}^\perp$ is given by

$$\{b(Z, \varepsilon) = h(Z) \varepsilon, h(Z) \text{ unrestricted}\} \cap L_2^0.$$

29

- To confirm that this is in fact correct, note that the subspace of $\Lambda_{\text{nuis},2}^\perp$ orthogonal to ε given Z is given by

$$\{C(\varepsilon, Z) - \mathbb{E}(C\varepsilon^T | Z) \mathbb{E}(\varepsilon\varepsilon^T | Z) \varepsilon : \mathbb{E}(C | Z) = 0\} \cap L_2^0.$$

Each element of this set is the residual of a projection of an element of $\{C(\varepsilon, Z) : \mathbb{E}(C | Z) = 0\} \cap L_2^0$ on the space of functions given by $\{b(X, \varepsilon) = h(Z) \varepsilon, h(Z) \text{ unrestricted}\} \cap L_2^0$ proving the result. In other words,

$$\Lambda_{\text{nuis},1} = \{C(\varepsilon, Z) - \mathbb{E}(C\varepsilon^T | Z) \mathbb{E}(\varepsilon\varepsilon^T | Z) \varepsilon : \mathbb{E}(C | Z) = 0\} \cap L_2^0$$

is an alternative characterization of $\Lambda_{\text{nuis},1}$.

- Proof of Result 5.4: Follows directly from definition of orthogonal projection. i.e.

$$\begin{aligned}S_{\text{eff}}(\varepsilon, Z; \beta^*, \eta^*) &= \Pi(S_\beta | \Lambda_{\text{nuis}}^\perp) \\ &= \Pi(S_\beta | \{b(Z, \varepsilon) = h(Z) \varepsilon, h(Z) \text{ unrestricted}\} \cap L_2^0) \\ &= \mathbb{E}\{S_\beta \varepsilon^T(\beta^*) | Z\} \mathbb{E}\{\varepsilon(\beta^*) \varepsilon^T(\beta^*) | Z\}^{-1} \varepsilon(\beta^*).\end{aligned}$$

30

Note also that for all submodels

$$\int \eta_1(\varepsilon(\beta), z) \varepsilon(\beta)^T d\mu(\varepsilon) = 0 \text{ for all } \beta$$

therefore

$$\begin{aligned}0 &= \int \frac{\partial \eta_1(\varepsilon(\beta), z)}{\partial \beta} \Big|_{\beta^*} \varepsilon(\beta^*)^T d\mu(\varepsilon) \\ &\quad - \int \eta_1(\varepsilon(\beta^*), z) \frac{\partial m(z; \beta)^T}{\partial \beta} \Big|_{\beta^*} d\mu(\varepsilon)\end{aligned}$$

which implies that

$$\mathbb{E}\{S_\beta \varepsilon^T(\beta^*) | Z\} = \frac{\partial m(Z; \beta)^T}{\partial \beta} \Big|_{\beta^*}.$$

31

We conclude that

$$\begin{aligned}S_{\text{eff}}(X; \beta^*, \eta^*) &= \mathbb{E}\{S_\beta \varepsilon^T(\beta^*) | Z\} \mathbb{E}\{\varepsilon(\beta^*) \varepsilon^T(\beta^*) | Z\}^{-1} \varepsilon(\beta^*) \\ &= D^T(Z) V(Z)^{-1} \varepsilon(\beta^*)\end{aligned}$$

where

$$\begin{aligned}D(Z) &= \frac{\partial m(Z; \beta)}{\partial \beta^T} \Big|_{\beta^*} \\ V(Z) &= \mathbb{E}\{\varepsilon(\beta) \varepsilon^T(\beta) | Z\}^{-1}\end{aligned}$$

32

Note also that by Result 5.3: The set of influence functions is given by

$$\left\{ \mathbb{E} \left\{ h(Z) \varepsilon(\beta^*) S_{\beta}^T \right\}^{-1} h(Z) \varepsilon(\beta^*) : h \right\} \cap L_2^0$$

Using a similar argument as above one can show that

$$\mathbb{E} \left\{ h(Z) \varepsilon(\beta) S_{\beta}^T \right\} = \mathbb{E} \left\{ h(Z) D(Z) \right\}$$

and therefore the set of influence functions is given by

$$\left\{ \mathbb{E} \left\{ h(Z) D(Z) \right\}^{-1} h(Z) \varepsilon : h \right\} \cap L_2^0$$

and the efficient influence function is given by

$$\mathbb{E} \left\{ D^T(Z) V(Z)^{-1} D(Z) \right\}^{-1} D^T(Z) V(Z)^{-1} \varepsilon$$

and the semiparametric efficiency bound is given by

$$\mathbb{E} \left\{ D^T(Z) V(Z)^{-1} D(Z) \right\}^{-1}$$

33

Adaptive Semiparametric Estimators for the Restricted Mean Model

- The class of influence functions is particularly useful as it provides very good candidate estimating equations.
- Specifically, let $u_h(X; \beta) = U_h(\beta) = h(Z)(Y - m(Z; \beta)) = h(Z)\varepsilon(\beta)$ with h of dimensions $\dim(\beta) \times \dim(\varepsilon)$. Note that $\{U_h(\beta); h\} \cap L_2^0 = \Lambda_{\text{nuis}}^{\perp}$. As we have previously shown, under fairly weak regularity conditions, the solution to equation

$$\sum_i U_{h,i}(\hat{\beta}_h) = 0 \quad (1)$$

for user specified h is a RAL estimator w influence function given by $\mathbb{E} \left\{ h(Z) D(Z) \right\}^{-1} h(Z) \varepsilon(\beta)$. The estimating equations (1) is an example of so-called generalized estimating equations. Using semiparametric theory, we have shown that any RAL estimator must have influence function (asymptotically) equivalent to one in our class!!!

34

Adaptive Semiparametric Estimators for the Restricted Mean Model

- It is reasonable to restrict attention to the semiparametric efficient estimator which is given by $\hat{\beta}_{h_{opt}}$ where

$$h_{opt}(Z) = D^T(Z) V(Z)^{-1} \varepsilon(\beta)$$

$$D(Z) = \frac{\partial m(Z; \beta)}{\partial \beta^T} \Big|_{\beta^*}$$

$$V(Z) = \mathbb{E} \left\{ \varepsilon(\beta^*) \varepsilon^T(\beta^*) | Z \right\}^{-1}$$

Thus this estimator is not feasible as it requires $D(Z)$ which depends on β^* and on $V(Z) = \text{var}(Y|Z)$.

35

Adaptive Semiparametric Estimators for the Restricted Mean Model

- The dependence on β^* is not much of an issue, as one can substitute in an (inefficient) preliminary estimator of β say by solving the GEE with user specified choice of $h_I(Z)$. For example when $m(Z; \beta) = g^{-1}((1, Z')\beta)$ where g is a link function, then it is standard to take $h_I(Z) = (1, Z)$. We denote such an initial estimator $\hat{\beta}_I$.
- Dependence of the efficient estimator on $\text{var}(Y|Z) = \text{var}(\varepsilon|Z)$ is potentially more challenging, since our semiparetric model allows the distribution of $\varepsilon|Z$ to be unrestricted. One may wish to proceed by estimating $\text{var}(Y|Z)$ nonparametrically using say kernel smoothing or similar techniques, however in finite sample the resulting estimator is likely to be very unstable particularly if Z is multi dimensional.

36

Adaptive Semiparametric Estimators for the Restricted Mean Model

- An alternative approach is to specify a working model for $V(Z)$, say $V(Z; \beta, \gamma)$ with finite dimensional parameter γ . which in the case of scalar Y may be estimated by solving the estimating equation:

$$0 = \sum_i T(Z_i, \hat{\beta}_I, \gamma) \left\{ (Y_i - m(Z; \hat{\beta}_I))^2 - V(Z; \hat{\beta}_I, \gamma) \right\}$$

where $T(Z_i, \hat{\beta}_I)$ is user-specified, e.g. $\partial V(Z; \hat{\beta}_I, \gamma) / \partial \gamma$, and of the same dimension as γ . Denote the resulting estimator $\hat{\gamma}$.

37

Adaptive Semiparametric Estimators for the Restricted Mean Model

- Replacing $V(Z; \hat{\beta}_I, \hat{\gamma})$ for $V(Z)$ into $h_{opt}(Z)$ will result under mild regularity conditions, in a RAL estimator for β with influence function

$$\mathbb{E} \left\{ D^T(Z) V(Z; \beta^*, \gamma^*)^{-1} D(Z) \right\}^{-1} D^T(Z) V(Z; \beta^*, \gamma^*)^{-1} \varepsilon$$

where γ^* is the probability limit of $\hat{\gamma}$.

- Consequently, solutions to such estimating equations where one uses a working variance model is what is called a locally efficient estimator because the resulting estimator is RAL under the semiparametric model irrespective of whether the variance model is correct or not, but it only achieves the semiparametric efficiency bound if the working variance model is correctly specified.

38

Adaptive Semiparametric Estimators for the Restricted Mean Model

- One must take care when estimating the asymptotic variance of locally efficient estimator, note that its asymptotic variance is given by the variance of its influence function

$$\begin{aligned} & \mathbb{E} \left\{ D^T(Z) V(Z; \beta^*, \gamma^*)^{-1} D(Z) \right\}^{-1} \\ & \mathbb{E} \left\{ D^T(Z) V(Z; \beta^*, \gamma^*)^{-1} V(Z) V(Z; \beta^*, \gamma^*)^{-1} D(Z) \right\} \\ & \mathbb{E} \left\{ D^T(Z) V(Z; \beta^*, \gamma^*)^{-1} D(Z) \right\}^{-1} \end{aligned}$$

the so-called sandwich variance estimator.

- This only simplifies to the SEB if the variance model is correct.
- Alternatively, one may use the nonparametric bootstrap.

39

Semiparametric location-shift model

- Suppose we observe iid data (Y, Z) Y is continuous that follows a semiparametric location-shift model

$$Y = m(Z; \beta^*) + \varepsilon,$$

where $m(Z, \cdot)$ is known, $\beta^* \in \mathbb{R}^s$ and $\varepsilon_i \perp Z_i$,

$$\eta(\varepsilon, Z) = \eta_1(\varepsilon) \eta_2(Z),$$

η_1 and η_2 unrestricted.

- We aim to make inferences about the finite dimensional parameter β^* .
- Note that for identification, we assume $m(0; \beta^*) = m(Z; 0) = 0$

40

Semiparametric location-shift model

► Result 5.5. As we will show below, the nuisance tangent space of LSM is given by:

$$\begin{aligned}\Lambda_{\text{nuis}} &= \Lambda_{\text{nuis},1} \oplus \Lambda_{\text{nuis},2} \subset L_2^0 \\ \Lambda_{\text{nuis},1} &= \{a_1(\varepsilon) : \mathbb{E}\{a_1(\varepsilon)\} = 0\} \cap L_2^0 \\ \Lambda_{\text{nuis},2} &= \{a_2(Z) : \mathbb{E}\{a_2(Z)\} = 0\} \cap L_2^0\end{aligned}$$

41

► Result 5.6: The ortho-complement to the nuisance tangent space

$$\begin{aligned}\Lambda_{\text{nuis}}^\perp &= \Lambda_{\text{nuis},1}^\perp \cap \Lambda_{\text{nuis},2}^\perp \cap L_2^0 \\ &= \left\{ \begin{pmatrix} h(Z, \varepsilon) - \mathbb{E}\{h(Z, \varepsilon) | Z\} \\ -\mathbb{E}\{h(Z, \varepsilon) | \varepsilon\} + \mathbb{E}\{h(Z, \varepsilon)\} \end{pmatrix} : h \right\} \cap L_2^0\end{aligned}$$

► Result 5.7: The set of influence functions is therefore given by

$$\left\{ \begin{pmatrix} \mathbb{E}\{\tilde{h}(Z, \varepsilon) S_\beta^T\}^{-1} \tilde{h}(Z, \varepsilon) \\ \tilde{h}(Z, \varepsilon) = \begin{pmatrix} h(Z, \varepsilon) - \mathbb{E}\{h(Z, \varepsilon) | Z\} \\ -\mathbb{E}\{h(Z, \varepsilon) | \varepsilon\} + \mathbb{E}\{h(Z, \varepsilon)\} \end{pmatrix} \end{pmatrix} : h \right\} \cap L_2^0$$

where S_β^T is the score of β under the LSM

42

► Result 5.8: The efficient score for β is given by

$$\begin{aligned}S_{\text{eff}}(X; \beta^*, \eta^*) &= \Pi(S_\beta | \Lambda_{\text{nuis}}^\perp) \\ &= \tilde{h}_{\text{opt}}(Z, \varepsilon) \\ h_{\text{opt}}(Z, \varepsilon) &= \dot{m}(Z; \beta^*) \kappa(\varepsilon) \\ \dot{m}(Z; \beta^*) &= \frac{\partial m(Z; \beta^*)}{\partial \beta} \Big|_{\beta^*} \\ \kappa(\varepsilon) &= -\frac{\partial \log \eta_1(\varepsilon)}{\partial \varepsilon}\end{aligned}$$

Because $\mathbb{E}\kappa(\varepsilon) = 0$ (show this), we have that

$$S_{\text{eff}}(X; \beta^*, \eta^*) = \left\{ \dot{m}(Z; \beta^*) - \mathbb{E}(\dot{m}(Z; \beta^*)) \right\} \kappa(\varepsilon)$$

43

► Proof Result 5.5. Consider a regular parametric submodel $\eta_t(\varepsilon, z) = \eta_{1,t}(\varepsilon) \eta_{2,t}(z)$ with unknown parameter $t \in (-\epsilon, \epsilon)$ such that $\eta_{t=0}(\varepsilon, z) = f(\varepsilon, z)$. It is straightforward to verify that corresponding score equations may be written

$$\frac{\partial \log \eta_t(\varepsilon, z)}{\partial t} = a_1(\varepsilon) + a_2(z)$$

where $\mathbb{E}\{a_1(\varepsilon)\} = \mathbb{E}\{a_2(Z)\} = 0$ and $\mathbb{E}\{a_1(\varepsilon) a_2(Z)\} = 0$ by the independence assumption. The first term a_1 is a score for the conditional density of ε , while a_2 is a score for the density of Z .

44

- From these results, we have that:

$$\begin{aligned}\Lambda_{\text{nuis}} &= \Lambda_{\text{nuis},1} \oplus \Lambda_{\text{nuis},2} \subset L_2^0 \\ \Lambda_{\text{nuis},1} &= \{a_1(\varepsilon) : \mathbb{E}\{a_1(\varepsilon)\} = 0\} \cap L_2^0 \\ \Lambda_{\text{nuis},2} &= \{a_2(Z) : \mathbb{E}\{a_2(Z)\} = 0\} \cap L_2^0\end{aligned}$$

- Technically, for each element of Λ_{nuis} we must also exhibit parametric submodels with corresponding score equal to the given element. I will leave it as hwk to work this out.

45

- Proof of Result 5.6. The ortho-complement to the nuisance tangent space

$$\begin{aligned}\Lambda_{\text{nuis}}^\perp &= (\Lambda_{\text{nuis},1} \oplus \Lambda_{\text{nuis},2})^\perp \\ &= \Lambda_{\text{nuis},1}^\perp \cap \Lambda_{\text{nuis},2}^\perp \cap L_2^0 \text{ (why?)}\end{aligned}$$

Note that any function $b(\varepsilon, Z) \in \Lambda_{\text{nuis},2}^\perp \cap L_2^0$ must satisfy

$$\mathbb{E}(b(Z, \varepsilon) a_2(Z)) = 0$$

for all $a_2(Z) \in \Lambda_{\text{nuis},2} \cap L_2^0$, that is $\mathbb{E}(b(Z, \varepsilon) | Z) = 0$. (why?)

- Therefore, to characterize $\Lambda_{\text{nuis}}^\perp$ entails finding all $b(\varepsilon, Z) \in \Lambda_{\text{nuis},2}^\perp$ such that $\mathbb{E}(b(X, \varepsilon) a_1(\varepsilon)) = 0$ for all $a_1(\varepsilon) \in \Lambda_{\text{nuis},1}^\perp$; that is $\mathbb{E}(b(X, \varepsilon) | \varepsilon) = 0$.
- Therefore, $\Lambda_{\text{nuis}}^\perp = \{b(\varepsilon, Z) : \mathbb{E}(b(X, \varepsilon) | \varepsilon) = \mathbb{E}(b(Z, \varepsilon) | Z) = 0\} \cap L_2^0$.

46

- We conjecture that

$$\Lambda_{\text{nuis}}^\perp = \left\{ \begin{pmatrix} h(Z, \varepsilon) - \mathbb{E}\{h(Z, \varepsilon) | Z\} \\ -\mathbb{E}\{h(Z, \varepsilon) | \varepsilon\} + \mathbb{E}\{h(Z, \varepsilon)\} \end{pmatrix} \right\} \cap L_2^0$$

- Proving the result entails showing that

$$\supseteq \left\{ \begin{pmatrix} h(Z, \varepsilon) - \mathbb{E}\{h(Z, \varepsilon) | Z\} \\ -\mathbb{E}\{h(Z, \varepsilon) | \varepsilon\} + \mathbb{E}\{h(Z, \varepsilon)\} \end{pmatrix} \right\} \cap L_2^0$$

and

$$\subseteq \left\{ \begin{pmatrix} h(Z, \varepsilon) - \mathbb{E}\{h(Z, \varepsilon) | Z\} \\ -\mathbb{E}\{h(Z, \varepsilon) | \varepsilon\} + \mathbb{E}\{h(Z, \varepsilon)\} \end{pmatrix} \right\} \cap L_2^0$$

- The first inclusion is easy to verify while the second follows from the fact that any $h(Z, \varepsilon)$ in $\Lambda_{\text{nuis}}^\perp$ must satisfy $\mathbb{E}\{h(Z, \varepsilon) | Z\} = \mathbb{E}\{h(Z, \varepsilon) | \varepsilon\} = 0$.

47

- Proof of Result 5.8: Follows directly from previous result on how to project onto $\Lambda_{\text{nuis}}^\perp$

$$\begin{aligned}S_{\text{eff}}(X; \beta^*, \eta^*) &= \Pi(S_\beta | \Lambda_{\text{nuis}}^\perp) \\ &= S_\beta - \mathbb{E}\{S_\beta | Z\} - \mathbb{E}\{S_\beta | \varepsilon\} + \mathbb{E}\{S_\beta\}\end{aligned}$$

- To show this, it suffices to note that $S_\beta - \mathbb{E}\{S_\beta | Z\} - \mathbb{E}\{S_\beta | \varepsilon\} + \mathbb{E}\{S_\beta\} \in \Lambda_{\text{nuis}}^\perp$ and

$$\begin{aligned}S_\beta - \Pi(S_\beta | \Lambda_{\text{nuis}}^\perp) &= \mathbb{E}\{S_\beta | Z\} + \mathbb{E}\{S_\beta | \varepsilon\} - \mathbb{E}\{S_\beta\}\end{aligned}$$

is orthogonal to $\Lambda_{\text{nuis}}^\perp$.

48

► Next, note that the score for β is given by

$$\begin{aligned} S_\beta &= \frac{\partial \log \eta_1(\varepsilon(\beta))}{\partial \beta} \\ &= -\frac{\partial \log \eta_1(\varepsilon(\beta))}{\partial \beta} \dot{m}(Z; \beta^*) \\ &= \dot{m}(Z; \beta^*) \kappa(\varepsilon) \end{aligned}$$

We conclude that

$$\begin{aligned} S_{\text{eff}}(X; \beta^*, \eta^*) &= \dot{m}(Z; \beta^*) \kappa(\varepsilon) - \mathbb{E} \left\{ \dot{m}(Z; \beta^*) \kappa(\varepsilon) | Z \right\} \\ &\quad - \mathbb{E} \left\{ \dot{m}(Z; \beta^*) \kappa(\varepsilon) | \varepsilon \right\} + \mathbb{E} \left\{ \dot{m}(Z; \beta^*) \kappa(\varepsilon) \right\} \\ &= \left\{ \dot{m}(Z; \beta^*) - \mathbb{E} \left\{ \dot{m}(Z; \beta^*) \right\} \right\} \kappa(\varepsilon) \end{aligned}$$

by independence and the fact that $\kappa(\varepsilon)$ has mean zero (why is that?)

49

► Note also that the set of influence functions is given by

$$\left\{ \begin{array}{c} \mathbb{E} \left\{ \tilde{h}(Z, \varepsilon) S_\beta^T \right\}^{-1} \tilde{h}(Z, \varepsilon) \\ \tilde{h}(Z, \varepsilon) = \begin{pmatrix} h(Z, \varepsilon) - \mathbb{E} \{ h(Z, \varepsilon) | Z \} \\ -\mathbb{E} \{ h(Z, \varepsilon) | \varepsilon \} + \mathbb{E} \{ h(Z, \varepsilon) \} \end{pmatrix} : h \end{array} \right\} \cap L_2^0$$

Note that because $\mathbb{E}_\beta \left\{ \tilde{h}(Z, \varepsilon(\beta)) \right\} = 0$ for all β , then

$$\mathbb{E} \left\{ \tilde{h}(Z, \varepsilon) S_\beta^T \right\} = -\mathbb{E} \left\{ \frac{\partial \tilde{h}(Z, \varepsilon)}{\partial \beta^T} \right\}$$

and therefore influence functions are of the form

$$-\mathbb{E} \left\{ \frac{\partial \tilde{h}(Z, \varepsilon)}{\partial \beta^T} \right\} \tilde{h}(Z, \varepsilon)$$

50

Adaptive Semiparametric Estimators for the Location Shift Model

► Therefore the efficient influence function is given by

$$\begin{aligned} &\mathbb{E} \left\{ \left\{ \dot{m}(Z; \beta^*) - \mathbb{E} \left\{ \dot{m}(Z; \beta^*) \right\} \right\}^{\otimes 2} \kappa(\varepsilon)^2 \right\}^{-1} \\ &\times \left\{ \dot{m}(Z; \beta^*) - \mathbb{E} \left\{ \dot{m}(Z; \beta^*) \right\} \right\} \kappa(\varepsilon) \end{aligned}$$

and the semiparametric efficiency bound is given by

$$\mathbb{E} \left\{ \kappa(\varepsilon)^2 \right\}^{-1} \left[\mathbb{E} \left\{ \left\{ \dot{m}(Z; \beta^*) - \mathbb{E} \left\{ \dot{m}(Z; \beta^*) \right\} \right\}^{\otimes 2} \right\} \right]^{-1}$$

51

► Note that as in the case of the restricted mean model, constructing a semiparametric estimator based on the efficient score requires knowing the density of ε or more precisely knowing $\kappa(\varepsilon)$. Consider a feasible semiparametric efficient estimator which is given by $\hat{\beta}_{\hat{\kappa}_{opt}}$ which solves

$$\mathbb{P}_n \left\{ \dot{m}(Z; \beta) - \mathbb{P}_n \left\{ \dot{m}(Z; \beta) \right\} \right\} \hat{\kappa} \left(Y - m(Z; \hat{\beta}_{\hat{\kappa}_{opt}}) \right) = 0$$

where we have replace $\mathbb{E} \left\{ \dot{m}(Z; \beta) \right\}$ with its sample version $\mathbb{P}_n \left\{ \dot{m}(Z; \beta) \right\}$ and κ with an estimator $\hat{\kappa}$.

52

Adaptive Semiparametric Estimators for the Location Shift Model

- Ideally, one may wish to proceed by estimating $\kappa(\varepsilon) = -\partial \log \eta_1(\varepsilon) / \partial \varepsilon$ in order to remain in the semiparametric model which does not restrict the latter. however in finite sample the resulting estimator is likely to be very unstable unless sample size is very large, particularly because it's the ratio of a derivative of density divided by a density , i.e. density appears in the denominator.

53

Adaptive Semiparametric Estimators for the for the Location Shift Model

- An alternative approach is to specify a working model for $\eta_1(\varepsilon)$, say $\eta_1(\varepsilon; \gamma)$ with finite dimensional parameter γ . which may be estimated by standard maximum likelihood, i.e. by solving :

$$\hat{\gamma} = \arg \max_{\gamma} \mathbb{P}_n \log \eta_1 \left(\varepsilon \left(\hat{\beta}_I \right); \gamma \right)$$

where $\hat{\beta}_I$ is an inefficient initial estimator of β solving

$$\mathbb{P}_n \left[\tilde{h} \left(\varepsilon \left(\hat{\beta}_I \right), Z \right) \right] = 0$$

for user-specified h , e.g. $h \left(\varepsilon \left(\hat{\beta}_I \right), Z \right) = \varepsilon \left(\hat{\beta}_I \right) \dot{m} (Z; \beta)$

54

Adaptive Semiparametric Estimators for the for the Location Shift Model

- Replacing $\eta_1(\varepsilon; \hat{\gamma})$ for $\eta_1(\varepsilon; \hat{\gamma})$ into \tilde{h}_{opt} will result under mild regularity conditions, in a RAL estimator for β with influence function

$$\mathbb{E} \left\{ \left\{ \dot{m} (Z; \beta^*) - \mathbb{E} \left\{ \dot{m} (Z; \beta^*) \right\} \right\} \frac{\partial \kappa (\varepsilon; \gamma^*)}{\partial \beta^T} \Big|_{\beta^*} \right\}^{-1} \\ \times \left\{ \dot{m} (Z; \beta^*) - \mathbb{E} \left\{ \dot{m} (Z; \beta^*) \right\} \right\} \left\{ \kappa (\varepsilon; \gamma^*) - \mathbb{E} \left\{ \kappa (\varepsilon; \gamma^*) \right\} \right\}$$

where γ^* is the probability limit of $\hat{\gamma}$, You will show this as homework.

55

Adaptive Semiparametric Estimators for the for the Location Shift Model

- Consequently, solutions to such estimating equations where one uses a model is a locally efficient estimator because the resulting estimator is RAL under the semiparametric model irrespective of whether the working model fpr the density of the uncentered residual is correct or not, but it only achieves the semiparametric efficiency bound if the working variance model is correctly specified.
- Note that resulting influence function has (up to a scale) an extra term $\left\{ \dot{m} (Z; \beta^*) - \mathbb{E} \left\{ \dot{m} (Z; \beta^*) \right\} \right\} \mathbb{E} \left\{ \kappa (\varepsilon; \gamma^*) \right\}$, This term is technically the price one pays for having to estimate $\mathbb{E} \left\{ \dot{m} (Z; \beta^*) \right\}$, however this term is exactly zero at the submodel where $\kappa (\varepsilon; \gamma^*)$ and therefore $\mathbb{E} \left\{ \kappa (\varepsilon; \gamma^*) \right\} = \mathbb{E} \left\{ \kappa (\varepsilon) \right\} = 0$ but not otherwise.

56

Cox Proportional Hazards Model with Censored Data

- Suppose we wish to make inferences about the log hazards ratio parameter β in a Cox PH model

$$\begin{aligned}\lambda_{T|Z}(t|z) &= \lim_{h \rightarrow 0} \left\{ \frac{\Pr(t \leq T \leq t+h | T \geq t, Z=z)}{h} \right\} \\ &= \lambda_0(t) \exp\{\beta^T z\}\end{aligned}$$

where $\lambda_0(t)$ is an unrestricted function of t .

57

Cox Proportional Hazards Model with Censored Data

- However, we observe data (V, Δ, Z) where

$$V = \min(T, C)$$

with T the underlying failure time and C time to censoring, or time on study, and

$$\Delta = I(T \leq C)$$

is the failure indicator, and throughout we assume conditional independent censoring

$$T \perp C | Z$$

which is needed for identification.

58

Cox Proportional Hazards Model with Censored Data

- The goal is to find semiparametric consistent and asymptotically normal and efficient estimators for β using data (V_i, Z_i, Δ_i) without making any additional assumptions on the underlying baseline hazard function λ_0 , on the conditional density of C given Z or on the distribution of Z .
- We will use the semiparametric theory we have developed to find the space of all influence functions of RAL estimators for β in this model which will motivate a class of RAL estimators among which we will derive the semiparametric efficient estimator.
- This will be accomplished by deriving the semiparametric nuisance tangent space and its orthogonal complement when influence functions lie and the efficient score.

59

Cox Proportional Hazards Model with Censored Data

- The density for a single observation is given by (check this)

$$\begin{aligned}p_{V, \Delta, Z}(V, \delta, Z) &= \{\lambda_0(v) \exp(\beta^T z)\}^\delta \exp\{-\Lambda_0(v) \exp(\beta^T z)\} \\ &\quad \times \{p_{C|Z}(v|z)\}^{1-\delta} \left\{ \int_v^\infty p_{C|Z}(u|z) du \right\}^\delta p_Z(z)\end{aligned}$$

where

$$\Lambda_0(v) = \int_0^v \lambda_0(u) du$$

is the cumulative baseline hazard function; likewise, let $\lambda_{C|Z}(v|z)$ and $\Lambda_{C|Z}(v|z)$ denote the conditional hazard and cumulative hazard functions of C given Z .

60

Cox Proportional Hazards Model with Censored Data

- The density may then be written

$$p_{V,\Delta,Z}(V,\delta,Z) = \{\lambda_0(v) \exp(\beta^T z)\}^\delta \exp\{-\Lambda_0(v) \exp(\beta^T z)\} \\ \times \{\lambda_{C|Z}(v|z)\}^{1-\delta} \{\exp\{-\Lambda_{C|Z}(v|z)\}\} p_Z(z)$$

- The semiparametric model is characterized by an r -dimensional parameter β of primary interest, and a infinite dimensional nuisance parameter $\eta = (\lambda_0(v), \lambda_{C|Z}(v|z), p_Z(z))$ where $\lambda_0(v)$ is an unrestricted positive function of v , $\lambda_{C|Z}(v|z)$ is an unrestricted positive function of v and z , and $p_Z(z)$ is an unrestricted density of z such that $\int p_Z(z) dz = 1$.
- The goal is to characterize the orthogonal complement of the nuisance tangent space which contains all influence functions of regular and asymptotically linear estimators of β , and to characterize the efficient score for β .

61

Cox Proportional Hazards Model with Censored Data

- Let $N(u) = I(T \leq u, \Delta = 1)$ denote the indicator of whether the person is observed to have an event prior to time u and $Y(u) = I(V \geq u)$ denote the indicator of being at risk at time u . The corresponding martingale increment is

$$dM(u) = dN(u) - \lambda_0(u) \exp(\beta^{*T} z) Y(u) du.$$

- Likewise, let $N_C(u) = I(C \leq u, \Delta = 0)$ denote the indicator of whether the person is observed to be censored prior to time u and $Y(u) = I(V \geq u)$ denote the indicator of being at risk at time u . The corresponding martingale increment is

$$dM_C(u) = dN_C(u) - \lambda_{C|Z}(u|z) Y(u) du.$$

62

Cox Proportional Hazards Model with Censored Data

- Result 5.9. The nuisance tangent space

$$\Lambda_{\text{nuis}} = \Lambda_1 \oplus \Lambda_2 \oplus \Lambda_3$$

where

$$\Lambda_1 = \left\{ \int a_1(u) dM(u) : \text{for all } r\text{-dimensional functions } a_1(u) \right\} \cap L_2^0 \\ \Lambda_2 = \left\{ \int a_2(u, Z) dM_C(u) : \text{for all } r\text{-dimensional functions } a_2(u, z) \right\} \cap L_2^0 \\ \Lambda_3 = \left\{ \begin{matrix} a_3(Z) \\ \text{for all } r\text{-dimensional functions } a_3(z) \text{ with} \\ \mathbb{E}(a_3(Z)) = 0 \end{matrix} \right\} \cap L_2^0$$

63

Cox Proportional Hazards Model with Censored Data

- Result 5.10. The orthogonal complement to the nuisance tangent space is given by

$$\Lambda_{\text{nuis}}^\perp = \Lambda_1^\perp \cap \Lambda_2^\perp \cap \Lambda_3^\perp \\ = \left\{ \int [a(u, Z) - a^\dagger(u)] dM(u) : \text{for all } r\text{-dimensional functions } a(u, Z) \right\} \cap L_2^0 \\ \text{where } a^\dagger(u) = \frac{\mathbb{E}\{a(u, Z) \exp(\beta^{*T} Z) Y(u)\}}{\mathbb{E}\{\exp(\beta^{*T} Z) Y(u)\}}$$

The efficient score of β is obtained for the choice $a(u, Z) = Z$, i.e.

$$S_\beta^{\text{eff}} = \Pi(S_\beta | \Lambda_{\text{nuis}}^\perp) \\ = \int \left[Z - \frac{\mathbb{E}\{Z \exp(\beta^{*T} Z) Y(u)\}}{\mathbb{E}\{\exp(\beta^{*T} Z) Y(u)\}} \right] dM(u)$$

64

Cox Proportional Hazards Model with Censored Data

- Proof of Result 5.9. Recall that the likelihood for one observation is given by

$$p_{V,\Delta,Z}(V, \delta, Z) = \left\{ \lambda_0(v) \exp(\beta^T z) \right\}^\delta \exp \left\{ -\Lambda_0(v) \exp(\beta^T z) \right\} \\ \times \left\{ \lambda_{C|Z}(v|z) \right\}^{1-\delta} \left\{ \exp \left\{ -\Lambda_{C|Z}(v|z) \right\} \right\} p_Z(z)$$

- Consider a parametric submodel in the direction of $\lambda_{C|Z}(v|z)$ assuming the rest of the likelihood is known

$$\left\{ \lambda_{C|Z,\theta}(v|z) \right\}^{1-\delta} \left\{ \exp \left\{ -\Lambda_{C|Z,\theta}(v|z) \right\} \right\},$$

such that $\lambda_{C|Z,0}(v|z) = \lambda_{C|Z}(v|z)$. Let $\nabla_\theta(\cdot) = \partial(\cdot)/\partial\theta|_{\theta=0}$.

65

Cox Proportional Hazards Model with Censored Data

- The corresponding score equation is given by

$$(1 - \Delta) \frac{\nabla_\theta \lambda_{C|Z,\theta}(V|Z)}{\lambda_{C|Z,\theta}(V|Z)} - \int_0^V \frac{\nabla_\theta \lambda_{C|Z,\theta}(u|Z)}{\lambda_{C|Z,\theta}(u|Z)} \lambda_{C|Z,\theta}(u|z) du \\ = (1 - \Delta) \frac{\nabla_\theta \lambda_{C|Z,\theta}(V|Z)}{\lambda_{C|Z,\theta}(V|Z)} \\ - \int \frac{\nabla_\theta \lambda_{C|Z,\theta}(u|Z)}{\lambda_{C|Z,\theta}(u|Z)} \lambda_{C|Z,\theta}(u|z) Y(u) du \\ = \int \frac{\nabla_\theta \lambda_{C|Z,\theta}(u|Z)}{\lambda_{C|Z,\theta}(u|Z)} dM_C(u)$$

66

Cox Proportional Hazards Model with Censored Data

- Since the Cox semiparametric model for $T|Z$ places no restriction on $\nabla_\theta \lambda_{C|Z,\theta}(u|z)$, the function $\nabla_\theta \lambda_{C|Z,\theta}(u|z) / \lambda_{C|Z,\theta}(u|z)$ can be any function of (v, z) and therefore the tangent space for the censoring mechanism is given by

$$\Lambda_2 = \left\{ \int a_2(u, Z) dM_C(u) : \text{for all } r\text{-dimensional functions } a_2(u, z) \right\} \cap L_2^0$$

- Next consider submodels $\lambda_{0,\theta}(v)$ in the direction of $\lambda_0(v)$, the corresponding likelihood contribution is given by

$$\left\{ \lambda_{0,\theta}(v) \exp(\beta^T z) \right\}^\delta \exp \left\{ -\Lambda_{0,\theta}(v) \exp(\beta^T z) \right\}$$

with score equation given by

$$\Delta \frac{\nabla_\theta \lambda_{0,\theta}(V)}{\lambda_{0,\theta}(V)} - \int_0^V \frac{\nabla_\theta \lambda_{0,\theta}(u)}{\lambda_{0,\theta}(u)} \exp(\beta^T Z) \lambda_{0,\theta}(u) du \\ = \int \frac{\nabla_\theta \lambda_{0,\theta}(u)}{\lambda_{0,\theta}(u)} dM(u)$$

67

Cox Proportional Hazards Model with Censored Data

- Since the Cox semiparametric model for $T|Z$ places no restriction on the baseline hazard, the function $\nabla_\theta \lambda_{0,\theta}(u) / \lambda_{0,\theta}(u)$ can be any function of v therefore the tangent space for the baseline hazard function is given by

$$\Lambda_1 = \left\{ \int a_1(u) dM(u, Z) : \text{for all } r\text{-dimensional functions } a_1(u) \right\} \cap L_2^0$$

- It is straightforward to show that scores for the density of Z span the space

$$\Lambda_3 = \{ \text{All } r\text{-dimensional functions } a_3(Z) \} \cap L_2^0$$

68

Cox Proportional Hazards Model with Censored Data

- Therefore the nuisance tangent space is given by

$$\Lambda_{\text{nuis}} = \Lambda_1 + \Lambda_2 + \Lambda_3$$

- Note however that because $M_C(u)$ and $M(u)$ are martingales, $0 = \mathbb{E}(M_C(u)|Z) = \mathbb{E}(M(u)|Z)$.
- Furthermore, because T and Z are conditionally independent, it is straightforward to show that

$$\mathbb{E} \left\{ \int a_1(u) dM(u, Z) \int a_2(u, Z) dM_C(u) \right\} = 0$$

for all $a_1(u)$ and $a_2(u, Z)$. Together, these properties imply that

$$\Lambda_{\text{nuis}} = \Lambda_1 \oplus \Lambda_2 \oplus \Lambda_3$$

proving the result.

69

Cox Proportional Hazards Model with Censored Data

- Proof of Result 5.10. The orthogonal complement to the nuisance tangent space is given by

$$\Lambda_{\text{nuis}}^\perp = \Lambda_1^\perp \cap \Lambda_2^\perp \cap \Lambda_3^\perp$$

- We first show that

$$\Lambda_2^\perp \cap \Lambda_3^\perp = \left\{ \int [a(u, Z)] dM(u) : \text{for all } r\text{-dimensional functions } a(u, Z) \right\}$$

- To see this, consider the nonparametric model where $\lambda_{T|Z}(t|z)$ is allowed to remain unrestricted. Then it is possible to show that the model places no restriction on the joint density of the observed data $p(v, \delta, z)$.

70

Cox Proportional Hazards Model with Censored Data

- One can use the previous proof to show that in this model, the *maximal* tangent space is given by

$$\Lambda^* = \Lambda_1^* \oplus \Lambda_2 \oplus \Lambda_3 = L_2^0$$

where

$$\Lambda_1^* = \left\{ \int a(u, Z) dM(u) : \text{for all } r\text{-dimensional functions } a(u, Z) \right\} \cap L_2^0$$

Therefore $\Lambda_1^* = (\Lambda_2 \oplus \Lambda_3)^\perp = \Lambda_2^\perp \cap \Lambda_3^\perp$.

- Finally, we have that $\Lambda_{\text{nuis}}^\perp$ is the subspace of Λ_1^* also orthogonal to Λ_1 ,

$$\left\{ \int a^*(u, Z) dM(u) : \text{for all } a^*(u, Z) \text{ st } \mathbb{E} \left\{ \int a^*(u, Z) dM(u) \int a^*(u) dM(u) \right\} = 0 \text{ for all } a^*(u) \right\} \cap L_2^0$$

71

Cox Proportional Hazards Model with Censored Data

- Note that all such functions can be written $\int \{a(u, Z) - a^\dagger(u)\} dM(u)$ for arbitrary $a(u, Z)$ where $\int \{a^\dagger(u)\} dM(u)$ is the projection of $\int a(u, Z) dM(u)$ onto the Hilbert subspace $\{\int a(u) dM(u) : a(u) \in L_2^0\}$.

- That is for any choice of $a(u, Z)$,

$$\mathbb{E} \left\{ \int \{a(u, Z) - a^\dagger(u)\} dM(u) \int a(u) dM(u) \right\} = 0, \text{ for all } a$$

72

Cox Proportional Hazards Model with Censored Data

- Using the fact that the covariance of martingale stochastic integrals is given by the expectation of the predictable covariance process, one can show that this implies that for all a

$$\mathbb{E} \left\{ \int \{a(u, Z) - a^\dagger(u)\} a(u) Y(u) \lambda_0(u) \exp(\beta^{*T} Z) du \right\} = 0, \text{ that is}$$

$$\mathbb{E} \{ \{a(u, Z) - a^\dagger(u)\} Y(u) \exp(\beta^{*T} Z) \} = 0$$

proving the result

$$a^\dagger(u) = \frac{\mathbb{E} \{ a(u, Z) Y(u) \exp(\beta^{*T} Z) \}}{\mathbb{E} \{ Y(u) \exp(\beta^{*T} Z) \}}$$

and therefore

$$\Lambda_{\text{nuis}}^\perp = \Lambda_1^\perp \cap \Lambda_2^\perp \cap \Lambda_3^\perp$$

$$= \left\{ \int [a(u, Z) - a^\dagger(u)] dM(u) : \text{for all } r\text{-dimensional functions } a(u, Z) \right\} \cap L_2^0$$

73

Cox Proportional Hazards Model with Censored Data

- A globally semiparametric efficient estimator is obtained by deriving an empirical efficient estimating function

$$\int \left[Z - \frac{\mathbb{P}_n \{ Z \exp(\beta^T Z) Y(u) \}}{\mathbb{P}_n \{ \exp(\beta^T Z) Y(u) \}} \right] dM(u)$$

where recall $dM(u) = \{dN(u) - \lambda_0(u) \exp(\beta^T z) Y(u) du\}$. The empirical estimating equation simplifies

$$0 = \mathbb{P}_n \int \left[Z - \frac{\mathbb{P}_n \{ Z \exp(\beta^T Z) Y(u) \}}{\mathbb{P}_n \{ \exp(\beta^T Z) Y(u) \}} \right] dN(u)$$

so that it does depend on the baseline hazard function, which is remarkable. why?

- The solution to this equation is therefore globally semiparametric efficient. Why?

74

Efficiency in randomized trials w baseline covariates

- Consider a double-blind placebo controlled randomized trial to evaluate the average causal effect of an intervention A which is binary on an outcome Y .
- By randomization the average causal effect on the additive scale is given by

$$\beta = \mathbb{E}(Y|A=1) - \mathbb{E}(Y|A=0)$$

the difference in average outcome between the two treatment arms.

- Throughout we will assume that $\Pr(A=1) = 1/2$, although other randomization schemes are easily handled.
- Suppose that in addition to A and Y , as routinely done in randomized trials, we also observe a large collection of baseline characteristics L .

75

Efficiency in randomized trials w baseline covariates

- By randomization, we have that treatment assignment is independent of all pretreatment characteristics, i.e. $A \perp L$.
- Our semiparametric model \mathcal{M} is therefore the set of all regular law $f_{Y,A,L}(y, a, l)$, such that

$$f_{Y,A,L}(y, a, l) = f_{Y|A,L}(y|a, l) f_{A|L}(a) f_L(l)$$

- Our goal is to characterize all influence functions and therefore all RAL estimators of $\beta = \mathbb{E}(Y|A=1) - \mathbb{E}(Y|A=0)$ in this model.
- Additionally, we wish to characterize the efficient influence function in this model

76

Efficiency in randomized trials w baseline covariates

- Tangent space of the model:
- Result 5.11. The tangent space of this semiparametric model is given by

$$\Lambda = \Lambda_1 \oplus \Lambda_2 \oplus \Lambda_3$$

where

$$\begin{aligned}\Lambda_1 &= \{b_1(Y, A, L) : \mathbb{E}(B_1|A, L) = 0\} \cap L_2^0 \\ \Lambda_2 &= \{b_2(A) : \mathbb{E}(B_2|L) = \mathbb{E}(B_2) = 0\} \cap L_2^0 \\ \Lambda_3 &= \{b_3(L) : \mathbb{E}(B_3) = 0\} \cap L_2^0\end{aligned}$$

- Result 5.12. The orthocomplement to the tangent space is given by

$$\begin{aligned}\Lambda^\perp &= \Lambda_1^\perp \cap \Lambda_2^\perp \cap \Lambda_3^\perp \\ &= \left\{ \left(A - \frac{1}{2} \right) (b(L) - \mathbb{E}(b(L))) : b(L) \right\} \cap L_2^0\end{aligned}$$

77

Efficiency in randomized trials w baseline covariates

- Influence functions:
- Consider the nonparametric model \mathcal{M}_{np} given by the set of all regular law $f_{Y,A,L}(y, a, l)$ where A and L are no longer assumed to be independent. As this is a nonparametric model the following result provides the unique influence function of $\beta = \mathbb{E}(Y|A=1) - \mathbb{E}(Y|A=0)$
- Result 5.13. The efficient IF of β in the nonparametric model \mathcal{M}_{np} is given by

$$IF_{\beta, np} = 2I(A=1)(Y - \mu_1) - 2I(A=0)(Y - \mu_0)$$

- Result 5.14. The set of IFs of β in the semiparametric model \mathcal{M} is given by

$$IF_{\beta, np} \oplus \Lambda^\perp = \left\{ \begin{aligned} &2I(A=1)(Y - \mu_1) - 2I(A=0)(Y - \mu_0) \\ &+ \left(A - \frac{1}{2} \right) (b(L) - \mathbb{E}(b(L))) \end{aligned} : b(L) \right\} \cap L_2^0$$

where $\mu_0 = \mathbb{E}(Y|A=0)$, $\mu_1 = \mathbb{E}(Y|A=1)$.

78

Efficiency in randomized trials w baseline covariates

- Efficient IF:
- Result 5.15. The efficient IF of β in the semiparametric model \mathcal{M} is given by

$$\begin{aligned}IF_\beta^{\text{eff}} &= \Pi(IF_{\beta, np}|\Lambda) \\ &= \Pi(IF_{\beta, np}|\Lambda_1 \oplus \Lambda_2 \oplus \Lambda_3) \\ &= \Pi(IF_{\beta, np}|\Lambda_1) \oplus \Pi(IF_{\beta, np}|\Lambda_2) \oplus \Pi(IF_{\beta, np}|\Lambda_3) \\ &= 2I(A=1)(Y - E(Y|A=1, L)) \\ &\quad + E(Y|A=1, L) - \mu_1 \\ &\quad - 2I(A=0)(Y - E(Y|A=0, L)) \\ &\quad - E(Y|A=0, L) + \mu_0\end{aligned}$$

79

Efficiency in randomized trials w baseline covariates

- Efficient IF:
- The efficient IF of β in the semiparametric model \mathcal{M} may also be written

$$\begin{aligned}IF_\beta^{\text{eff}} &= 2I(A=1)(Y - \mu_1) - (2I(A=1) - 1)(E(Y|A=1, L) - \mu_1) \\ &\quad - 2I(A=0)(Y - \mu_0) + (2I(A=0) - 1)(E(Y|A=0, L) - \mu_0) \\ &= 2I(A=1)(Y - \mu_1) - 2I(A=0)(Y - \mu_0) \\ &\quad - (2I(A=1) - 1)(E(Y|A=1, L) - E(Y|A=0, L) - \mu_0 + \mu_1)\end{aligned}$$

which is in the set of influence functions

$$IF_{\beta, np} \oplus \Lambda^\perp = \left\{ \begin{aligned} &2I(A=1)(Y - \mu_1) - 2I(A=0)(Y - \mu_0) \\ &+ \left(A - \frac{1}{2} \right) (b(L) - \mathbb{E}(b(L))) \end{aligned} : b(L) \right\} \cap L_2^0$$

with $b(L) = b^{\text{opt}}(L) = -2(E(Y|A=1, L) - E(Y|A=0, L))$.

80

Efficiency in randomized trials w baseline covariates

- Efficient IF:
- Now we note that $IF_{\beta,np} = IF_{\beta,np} - \Pi(IF_{\beta,np}|\Lambda^\perp) + \Pi(IF_{\beta,np}|\Lambda^\perp)$
- Because $IF_{\beta,np} - \Pi(IF_{\beta,np}|\Lambda^\perp) \in \Lambda$ and $\Pi(IF_{\beta,np}|\Lambda^\perp) \in \Pi(IF_{\beta,np}|\Lambda^\perp)$, we conclude that

$$IF_{\beta}^{\text{eff}} = IF_{\beta,np} - \Pi(IF_{\beta,np}|\Lambda^\perp)$$

which implies that

$$\begin{aligned} & \Pi(IF_{\beta,np}|\Lambda^\perp) \\ = & \left(I(A=1) - \frac{1}{2} \right) \\ & \times 2(E(Y|A=1, L) - E(Y|A=0, L) - \beta) \end{aligned}$$

- This further implies that

$$\begin{aligned} \text{var}(IF_{\beta}^{\text{eff}}) &= \text{var}(IF_{\beta,np} - \Pi(IF_{\beta,np}|\Lambda^\perp)) \\ &\leq \text{var}(IF_{\beta,np}) \end{aligned}$$

81

Efficiency in randomized trials w baseline covariates

- Efficient IF:
- It is straightforward to establish that $IF_{\beta,np}$ is the influence function of the nonparametric estimator of β given by

$$\hat{\beta} = \frac{\mathbb{P}_n(YA)}{\mathbb{P}_n(A)} - \frac{\mathbb{P}_n(Y(1-A))}{\mathbb{P}_n((1-A))}$$

Check this.

- As a consequence, we have established that the difference in sample means which is the standard to analyze randomized trial is inefficient when as usually the case baseline characteristics are available in a randomized trial.

82

Efficiency in randomized trials w baseline covariates

- Efficient IF:
- A semiparametric locally efficient estimator can be obtained by using the EIF as an estimating equation,
- First note that

$$\begin{aligned} IF_{\beta}^{\text{eff}} &= 2I(A=1)(Y - E(Y|A=1, L)) \\ &\quad - 2I(A=0)(Y - E(Y|A=0, L)) \\ &\quad + E(Y|A=1, L) - E(Y|A=0, L) - \beta \end{aligned}$$

depends on the unknown conditional mean function $E(Y|A, L)$ which we must estimate in order to obtain a feasible estimating function. Suppose we posit a simple parametric model $b(A, L; \eta)$, say

$$b(A, L; \eta) = (1, A, AL', L')\eta$$

and suppose that η is estimated by OLS.

83

Efficiency in randomized trials w baseline covariates

- Efficient IF:
- $\tilde{\beta}$ can be shown to be consistent and asymptotically normal in model \mathcal{M} and to achieve the semiparametric efficiency bound if $b(1, L; \eta)$ is correctly specified, so that $b(a, l; \hat{\eta})$ converges to $\mathbb{E}(Y|a, l)$ for all (a, l) in probability.
- That is $\tilde{\beta}$ is completely robust to model misspecification of $b(1, L; \eta)$ and achieve the efficiency bound if $b(1, L; \eta)$ is correct model. In other words $\tilde{\beta}$ is a semiparametric locally efficient estimator.
- One may be tempted to use a nonparametric model for $b(1, L; \eta)$, this will typically only perform well if L is very low dimensional, e.g. continuous w one or at most 2 components, or discrete w few levels.

84

Outline

- 1 Semiparametric models
- 2 Causal effects of a point exposure
 - Observational study: The case of point exposure
 - G-computation
 - G-computation: Mathematical Interlude
 - Summary
 - G-computation: Efficiency considerations
 - IPTW estimator in a randomized trial: an efficiency paradox

Observational study: The case of point exposure

- Suppose that randomization no longer holds, because the observed data comes from a point exposure/cross-sectional observational study, with observed data $\{L, A, Y\}$.
- L is a rich vector of covariates that satisfies :
(NUCA) No unmeasured confounding assumption holds: $\{Y_0, Y_1\} \perp A | L$ where Y_a is the potential outcome had exposure been a
- Then we say that there are no unmeasured confounders for the effect of A on Y .

Observational study: The case of point exposure

- The intuition behind (NUCA) is similar to that of RA. Mainly, that we have measured enough covariates L , so that within levels of L , the data mimicks a randomized trial with the randomization probabilities now allowed to depend on L .
- Conceptually, this can be achieved only if we are able to measure all common causes of A and Y (that is all risk factors for Y that also determine A).

Observational study: The case of point exposure

- Next we show that the no unmeasured confounding assumption is sufficient to again identify $\{E(Y_a) : a\}$ and thus $\psi = E(Y_1) - E(Y_0)$
Without loss of generality, suppose L is categorical; then
$$E(Y_a) = E(E(Y_a|L)) = \sum_l E(Y_a|L=l) f_L(l)$$
$$\stackrel{NUCA}{=} \sum_l E(Y_a|A=a, L=l) f_L(l)$$
$$\stackrel{CA}{=} \sum_l E(Y|A=a, L=l) f_L(l)$$
$$\equiv g(a)$$

Observational study: The case of point exposure

Observational study: The case of point exposure

- ▶ $g(a)$ is known as the *direct standardization* of $E(Y|A=a, L)$. It is a special case of Robins' *G-formula* (which we will discuss in the longitudinal case).
- ▶ Thus
$$\psi = g(1) - g(0) = \sum_l \{E(Y|A=1, L=l) - E(Y|A=0, L=l)\} f_L(l)$$
is the *standardized risk difference*.
- ▶ Under NUCA, we see that crude association \neq causation, as
$$\sum_l E(Y_a|A=a, L=l) f_L(l) = E(Y_a) \neq E(Y|A=a) = \sum_l E(Y|A=a, L=l) f_L(l|A=a).$$

- ▶ So that the crude risk difference does not have a causal interpretation. However, if NUCA holds, and either of the following conditions holds:

$$Y \perp\!\!\!\perp L|A \text{ or } A \perp\!\!\!\perp L \tag{2}$$

then $E(Y_a) = E(Y|A=a)$ and L is a non-confounder, so that this implies that RA actually holds.

Observational study: The case of point exposure

Observational study: The case of point exposure

Proof:

- ▶ In the first case,
$$E(Y_a) = \sum_l E(Y|A=a, L=l) f_L(l) = \sum_l E(Y|A=a) f_L(l) = E(Y|A=a);$$
- ▶ In the second case,
$$E(Y_a) = \sum_l E(Y|A=a, L=l) f_L(l) = \sum_l E(Y|A=a, L=l) f_L(l|A=a) = E(Y|A=a)$$

- ▶ In general, the point exposure G-formula is written

$$E(Y_a) = \int E(Y|A=a, L=l) dF(l)$$

- ▶ The left-hand side is the mean of a counterfactual (latent variable), the right-hand side is a functional of the observed data, which is always well defined but only has a causal interpretation under the no unmeasured assumption.

Observational study: The case of point exposure

- ▶ This functional is not a conditional expectation of the observed data, therefore, cannot be estimated directly such as the crude. However, the plug-in principle may be used as discussed below.
- ▶ The G-formula is not restricted to the mean, other versions:
 - ▶ $E(Y_a|V=v) = \int E(Y|A=a, L=l) dF(l|V=v)$, where V is contained in L
 - ▶ $f(Y_a|V=v) = \int f(Y|A=a, L=l) dF(l|V=v)$.

93

G-computation

- ▶ Given the observed data $O_i = (Y_i, A_i, L_i)$, G-computation generally refers to nonparametric inference on the G-formula

$$g(a) = \sum_l E(Y|A=a, L=l) f_L(l).$$

- ▶ A natural nonparametric estimator of $g(a)$ is given by the nonparametric plug-in estimator, which requires nonparametric estimates of $E(Y|A=a, L=l) = b(a, l)$ written $\hat{b}(a, l)$ and of $F_L(l)$ written $\hat{F}_L(l)$.

94

G-computation

- ▶ Until otherwise stated, assume both A and L are categorical variables with low to moderate number of levels, so that $\hat{b}(a, l)$ is given by the stratified sample average:

$$\hat{b}(a, l) = \sum_{i=1}^n I(A_i = a, L_i = l) Y_i / \sum_{i=1}^n I(A_i = a, L_i = l)$$

$$\text{and } \hat{f}_L(l) = n^{-1} \sum_{i=1}^n I(L_i = l)$$

- ▶ The nonparametric estimator of the G-formula is given by:

$$\begin{aligned} \hat{g}(a) &= \sum_l \hat{b}(a, l) \hat{f}_L(l) = \sum_l \hat{b}(a, l) n^{-1} \sum_{i=1}^n I(L_i = l) \\ &= n^{-1} \sum_{i=1}^n \sum_l \hat{b}(a, l) I(L_i = l) = n^{-1} \sum_{i=1}^n \hat{b}(a, L_i) \end{aligned}$$

95

G-computation: Mathematical Interlude

- ▶ How to use $\hat{g}(a)$ to make inference on $\psi = g(1) - g(0)$?
- ▶ e.g. Want 95%CI of the average causal effect of A on Y .
- ▶ Wald type CI:
 - ▶ $\hat{\psi} \pm 1.96 \sqrt{\widehat{\text{var}}(\hat{\psi})}$, where $\hat{\psi} = \hat{g}(1) - \hat{g}(0)$ and $\widehat{\text{var}}(\hat{\psi})$ is a consistent estimator of the variance of $\hat{\psi}$.
 - ▶ Coverage of this CI depends on the asymptotic distribution of $\hat{\psi}$ and on identifying a reasonable estimator of $\widehat{\text{var}}(\hat{\psi})$.
 - ▶ Difficulty, $\hat{g}(a) = n^{-1} \sum_{i=1}^n \hat{b}(a, L_i)$ is a sample average, however of terms that are estimated using $\{Y_i, A_i, L_i : i = 1, \dots, n\}$ so that though $b(a, L_i)$ and $b(a, L_j), i \neq j$ are independent, $\hat{b}(a, L_i)$ and $\hat{b}(a, L_j), i \neq j$ are highly dependent and standard *i.i.d* Central Limit Theorem doesn't apply.

96

► However it is still possible to show asymptotic normality of $\hat{\psi}$; in fact, this is an instructive exercise, write:

$$\begin{aligned} n^{1/2}\hat{g}(a) &= n^{-1/2} \sum_l \sum_{i=1}^n \hat{b}(a, l) I(L_i = l) \\ &= n^{-1/2} \sum_l \sum_{i=1}^n I(L_i = l) \left\{ \frac{\sum_{s=1}^n I(A_s = a, L_s = l) Y_s}{\sum_{j=1}^n I(A_j = a, L_j = l)} \right\} \\ &= n^{-1/2} \sum_l \sum_{s=1}^n I(A_s = a, L_s = l) Y_s \left\{ \frac{\sum_{i=1}^n I(L_i = l)}{\sum_{j=1}^n I(A_j = a, L_j = l)} \right\} \end{aligned}$$

$$\begin{aligned} &= n^{-1/2} \sum_l \sum_{s=1}^n \left\{ I(A_s = a, L_s = l) (Y_s - b(a, l)) \right. \\ &\quad \left. \times \left\{ \frac{n^{-1} \sum_{i=1}^n I(L_i = l)}{n^{-1} \sum_{j=1}^n I(A_j = a, L_j = l)} \right\} \right\} \\ &+ n^{-1/2} \sum_l b(a, l) \left\{ \sum_{i=1}^n I(L_i = l) \right\} \\ &= n^{-1/2} \sum_l \sum_{s=1}^n I(A_s = a, L_s = l) (Y_s - b(a, l)) f_{A|L}^{-1}(a|l) \\ &+ n^{-1/2} \sum_l b(a, l) \left\{ \sum_{i=1}^n I(L_i = l) \right\} + o_p(1) \\ &= n^{-1/2} \sum_{i=1}^n \frac{I(A_i = a)}{f_{A|L}(A_i|L_i)} (Y_i - b(A_i, L_i)) + b(a, L_i) + o_p(1) \end{aligned}$$

97

98

so that we finally obtain

$$\begin{aligned} &n^{1/2}(\hat{g}(a) - g(a)) \\ &= n^{-1/2} \sum_{i=1}^n \left\{ \frac{I(A_i = a)}{f_{A|L}(A_i|L_i)} (Y_i - b(A_i, L_i)) + b(a, L_i) - g(a) \right\} + o_p(1) \\ &= n^{-1/2} \sum_{i=1}^n IF_i(a) + o_p(1) \\ &\sim N\left(0, E\left(IF_i(a)^2\right)\right) \end{aligned}$$

A wald type 95%CI for $\psi = g(1) - g(0)$ is given by:

$$\hat{g}(1) - \hat{g}(0) \pm 1.96 \sqrt{n^{-1} \sum_i \left(\widehat{IF}_i(0) - \widehat{IF}_i(1) \right)^2}$$

$$\text{where } \widehat{IF}_i(a) = \left\{ \frac{I(A_i = a)}{\hat{f}_{A|L}(a|L_i)} \left(Y_i - \hat{b}(A_i, L_i) \right) + \hat{b}(a, L_i) - \hat{g}(a) \right\},$$

and $\hat{f}_{A|L}(a|L_i = l) = \frac{\sum_i I(A_i = a, L_i = l)}{\sum_i I(L_i = l)}$ is the nonparametric estimator of $f_{A|L}(a|L_i = l)$, the probability of receiving treatment $A = a$ given $L = l$.

99

100

Summary

- ▶ The term $b(a, L_i) - g(a)$ reflects the variability due to the estimation of $F_L(l)$, whereas $\frac{I(A_i=a)}{f_{A|L}(a|L_i)} (Y_i - b(a, L_i))$ captures the variability due to the estimation of $b(a, L_i)$.
- ▶ It is interesting to note that the influence function and thus the variance of $\hat{\psi} = \hat{g}(1) - \hat{g}(0)$ depends on the treatment process $\{f_{A|L}(a|L_i) : a\}$, even though the pluggin estimator described above appears not to.
- ▶ In fact, we give a different representation of the nonparametric estimator of the G-formula which makes explicit its dependence on the estimated treatment process.

101

- ▶ Goal: Make inference on the causal mean difference: $\psi = E(Y_1 - Y_0) = E(Y_1) - E(Y_0)$
 - ▶ where $Y_a, a = 0, 1$ is the counterfactual/potential outcome that a patient would experience if possibly countering to fact he/she receives treatment $a = 0, 1$.
 - ▶ $E(Y_a), a = 0, 1$ has the interpretation of the population outcome mean that we would observe if everyone in the population was forced to take treatment $a = 0, 1$.
 - ▶ This is to be contrasted with $E(Y|A = a)$ which is the outcome mean among those who did take treatment $A = a$, which may or may not equal $E(Y_a)$.

102

- ▶ Identifying assumptions:
 - ▶ Randomization: $\{Y_0, Y_1\} \perp\!\!\!\perp A \Rightarrow E(Y_a) = E(Y|A = a) \Leftrightarrow$ **crude association is causal**.
 - ▶ No unmeasured confounding: $\{Y_0, Y_1\} \perp\!\!\!\perp A | L \Rightarrow E(Y_a) = g(a) = \sum_l E(Y|A = a, L = l) f_L(l) \Rightarrow$ **crude association is not causal**. To get causal effect ψ , we must compute the g-formula $g(a), a = 0, 1$; This naturally adjust for the confounders L
 - ▶ No unmeasured confounding essentially means that we have collected all common causes of A and Y in L . This is a strong assumption which requires subject matter expertise.

103

- ▶ Inference on the causal effect ψ :
 - ▶ If randomization holds standard crude analysis still applies, such as two sample t-tests ...etc .
 - ▶ under no unmeasured confounding, we must obtain point and interval estimates of ψ .
 - ▶ To get point estimates nonparametrically, we use the pluggin principle: $\hat{g}(a) = \sum_l \hat{b}(a, l) \hat{f}_L(l)$, where $\hat{b}(a, l)$ and $\hat{f}_L(l)$ are nonparametric estimates of **nuisance parameters** not of primary interest but required to define our causal contrast.
 - ▶ Nonparametric method here is appealing because we do not want to specify models for $b(a, l)$ and $f_L(l)$ if not necessary as neither is of primary interest.

104

- We also proved that:

$$\sqrt{n}(\hat{\psi} - \psi) = \sqrt{n}(\hat{g}(1) - \hat{g}(0) - \psi) \approx N(0, E[(IF(0) - IF(1))]^2)$$

$$\text{where } IF_i(a) = \left\{ \frac{I(A_i=a)}{\hat{f}_{A|L}(a|L_i)} (Y_i - b(A_i, L_i)) + b(a, L_i) - g(a) \right\}$$

- which we used to construct a Wald-type 95%CI

$$\hat{\psi} \pm n^{-1/2} 1.96 \sqrt{n^{-1} \sum_i \left(\widehat{IF}_i(0) - \widehat{IF}_i(1) \right)^2}$$

$$\text{where } \widehat{IF}_i(a) = \left\{ \frac{I(A_i=a)}{\hat{f}_{A|L}(a|L_i)} (Y_i - \hat{b}(A_i, L_i)) + \hat{b}(a, L_i) - \hat{g}(a) \right\}.$$

- Finally we found that: $\hat{g}(a) = \sum_l \hat{b}(a, l) \hat{f}_L(l) = \frac{\sum_{s=1}^n I(A_s=a) Y_s \hat{f}_{A|L}^{-1}(A_s|L_s)}{\sum_{s=1}^n I(A_s=a) \hat{f}_{A|L}^{-1}(A_s|L_s)}$ has an exact dual representation as an inverse-probability of treatment weighted (iptw) estimator.

105

G-computation: Efficiency considerations

The observed data likelihood is given by

$$\begin{aligned} & \prod_i f_{Y_i, L_i, A_i}(Y_i, L_i, A_i) \\ &= \prod_i \int f_{Y_{0i}, Y_{1i}, L_i, A_i}(Y_{0i}, Y_{1i}, L_i, A_i) dY_{1-A_i} \\ &= \prod_i \int f_{Y_{0i}, Y_{1i}, L_i}(Y_{0i}, Y_{1i}, L_i) f_{A|Y_{0i}, Y_{1i}, L_i}(A_i|Y_{0i}, Y_{1i}, L_i) dY_{1-A_i} \\ &\stackrel{NUCA}{=} \prod_i f_{A|Y_{0i}, Y_{1i}, L_i}(A_i|L_i) \int f_{Y_{0i}, Y_{1i}, L_i}(Y_{0i}, Y_{1i}, L_i) dY_{1-A_i} \\ &= \left\{ \prod_i f_{A|Y_{0i}, Y_{1i}, L_i}(A_i|L_i) \right\} \times \prod_i f_{Y_{0i}|L_i}(Y_{0i}|L_i)^{(1-A_i)} f_{Y_{1i}|L_i}(Y_{1i}|L_i)^{A_i} f_L(L_i) \end{aligned}$$

106

Thus the observed data likelihood factorizes into a part that depends on the treatment process (first factor) and a part that depends on the full data process (second factor). For simplicity, assume that Y is binary, so that the likelihood becomes

$$\begin{aligned} & \prod_i f_{A|Y_{0i}, Y_{1i}, L_i}(A_i|L_i) \\ & \times \prod_i \left\{ \begin{array}{l} \Pr(Y_{0i}=1|L_i)^{Y_{0i}(1-A_i)} \Pr(Y_{0i}=0|L_i)^{(1-Y_{0i})(1-A_i)} \\ \times \Pr(Y_{1i}=1|L_i)^{Y_{1i}A_i} \Pr(Y_{1i}=0|L_i)^{(1-Y_{1i})A_i} f_L(L_i) \end{array} \right\} \\ &= \prod_i f_{A|Y_{0i}, Y_{1i}, L_i}(A_i|L_i) \times \\ & f_L(L_i) \times \prod_i \left\{ \begin{array}{l} b(0, L_i)^{Y_{0i}(1-A_i)} (1-b(0, L_i))^{(1-Y_{0i})(1-A_i)} \\ \times b(1, L_i)^{Y_{1i}A_i} (1-b(1, L_i))^{(1-Y_{1i})A_i} \end{array} \right\} \end{aligned}$$

107

- The factorization is similar to that encountered in missing data problems, under 'missingness at random'.
- The treatment process is said to be ancillary to inferences on $\{g(0), g(1)\}$ and thus on $\psi = g(1) - g(0)$
- It is straightforward to show that the MLE of $g(a)$ is $\hat{g}(a) = \sum_l \hat{b}(a, l) \hat{f}_L(l)$, $0, 1$. Thus the pluggin estimator attains the Cramer-Rao efficiency bound.
- Moreover, because the likelihood factorizes, the Cramer-Rao efficiency bound is the same whether or not the treatment process is known. Though the bound itself (the asymptotic variance of $\hat{g}(1) - \hat{g}(0)$) depends on the treatment process!!!

108

- We now show that the nonparametric G-computation estimator $\hat{g}(1) - \hat{g}(0)$ has a dual representation as an inverse-probability of treatment weighted (iptw) estimator:

$$\begin{aligned}\hat{g}(a) &= \sum_l \hat{b}(a, l) \hat{f}_L(l) = n^{-1} \sum_l \sum_{i=1}^n \hat{b}(a, l) I(L_i = l) \\ &= n^{-1} \sum_l \sum_{i=1}^n I(L_i = l) \left\{ \frac{\sum_{s=1}^n I(A_s = a, L_s = l) Y_s}{\sum_{j=1}^n I(A_j = a, L_j = l)} \right\} \\ &= n^{-1} \sum_l \sum_{s=1}^n I(A_s = a, L_s = l) Y_s \left\{ \frac{\sum_{i=1}^n I(L_i = l)}{\sum_{j=1}^n I(A_j = a, L_j = l)} \right\} \\ &= n^{-1} \sum_{s=1}^n I(A_s = a) Y_s \hat{f}_{A|L}^{-1}(A_s | L_s) = \frac{\sum_{s=1}^n I(A_s = a) Y_s \hat{f}_{A|L}^{-1}(A_s | L_s)}{\sum_{s=1}^n I(A_s = a) \hat{f}_{A|L}^{-1}(A_s | L_s)}\end{aligned}$$

109

- $\hat{f}_{A|L}(a|L=l) = \frac{\sum_i I(A_i=a, L_i=l)}{\sum_i I(L_i=l)}$ is the nonparametric estimator of the treatment conditional probability mass function $f_{A|L}(a|L=l)$.
- So that the nonparametric G-formula estimator has a (exact) dual representation as an inverse-probability-of-treatment weighted (iptw) estimator, and this also shows that the latter is asymptotically efficient.

110

- One can easily show (try it) that the iptw estimator $\hat{g}(a)$ is the solution to the estimating equation: $\sum_{i=1}^n \frac{I(A_i=a)}{\hat{f}_{A|L}(a|L_i=l)} (Y_i - \hat{g}(a)) = 0$.
- In contrast to the estimating equation for the conditional mean $\mu(a) = E(Y|A=a)$ given by: $\sum_{i=1}^n I(A_i=a) (Y_i - \hat{\mu}(a)) = 0$.
- We again see that if $f_{A|L}(a|L_i=l) = f_A(a)$, then $E(Y|A=a) = E(Y_a)$ and association is causation.
- If $f_{A|L}(a|L_i=l) \neq f_A(a)$ and L is a risk factor of Y , then to obtain a counterfactual mean $E(Y_a)$, the recipe is to weight the estimating function for the conditional mean by $\hat{f}_{A|L}^{-1}(a|L_i=l)$ to adjust for confounding by L .

111

- Why does weighing adjust for confounding?
 - Essentially, creates a 'pseudo-population' in which A is no longer associated with Y ; so that the crude mean in this population has a counterfactual interpretation.
 - For example, suppose L is binary and we observe :
- | N | A | L | $E(Y A=a, L=l)$ |
|------|-----|-----|-----------------|
| 4000 | 1 | 0 | 24 |
| 3000 | 1 | 1 | 36 |
| 8000 | 0 | 0 | 10 |
| 9000 | 0 | 1 | 22 |
- So that $f_{A|L}(A=1|L=1) = 1/4$, $f_{A|L}(A=1|L=0) = 1/3$ and thus L predicts A .
 - Moreover $E(Y|A=1, L=l) = 24 + 12l$ so that L predict Y given A .

112

IPTW estimator in a randomized trial: an efficiency paradox

- The crude mean $E(Y|A=1) = \sum_l E(Y|A=1, L=l) f(L=l|A=1) = 24 \times 4/7 + 36 \times 3/7 = 204/7$

- Whereas $E(Y_1) = \sum_l E(Y|A=1, L=l) f(L=l) = 24 \times 1/2 + 36 \times 1/2 = 210/7$

- Create a pseudo population by reweighting the numbers in each row by $f_{A|L}^{-1}$

N	$f(A L)$	$Pseudo - N$	A	L	$E(Y A=a, L=l)$
4000	1/3	$4000 \times 3 = 12,000$	1	0	24
3000	1/4	$3000 \times 4 = 12,000$	1	1	36
8000	2/3	$8000 \times 3/2 = 12,000$	0	0	10
9000	3/4	$9000 \times 4/3 = 12,000$	0	1	22

- So that in the Pseudo population, the crude analysis gives: $E^*(Y|A=1) = 24 \times 1/2 + 36 \times 1/2 = 210/7 = E(Y_1)$

- Do a similar calculation for $E(Y_0)$

- Lets revisit the familiar setting where A is randomized. Denote by L a pretreatment categorical covariate that is a strong predictor of the outcome Y .

- By randomization, we have: $\{Y_0, Y_1, L\} \perp\!\!\!\perp A$

- And as before, the crude estimator

$$\tilde{\psi} = \frac{\sum_{s=1}^n I(A_s=1) Y_s}{\sum_{s=1}^n I(A_s=1)} - \frac{\sum_{s=1}^n I(A_s=0) Y_s}{\sum_{s=1}^n I(A_s=0)}$$

is a valid estimator of the average causal effect ψ_0 .

113

114

- Note that this estimator may also be written as an iptw estimator with known weights $f_{A|L}^{-1}(A_s|L_s) = (1/2)^{-1}$

$$\tilde{\psi} = \frac{\sum_{s=1}^n I(A_s=1) Y_s (1/2)^{-1}}{\sum_{s=1}^n I(A_s=1) (1/2)^{-1}} - \frac{\sum_{s=1}^n I(A_s=0) Y_s (1/2)^{-1}}{\sum_{s=1}^n I(A_s=0) (1/2)^{-1}}$$

- It is easy to show that the asymptotic variance of this estimator is given by

$$var(\sqrt{n}\tilde{\psi}) = E((Y - g(1))^2 | A=1) + E((Y - g(0))^2 | A=0)$$

- Compare this estimator to the one used by a statistician who decides to ignore randomization but rather assumes that $\{Y_0, Y_1\} \perp\!\!\!\perp A|L$, so that she may use G-computation

$$\begin{aligned} \hat{\psi} &= \hat{g}(1) - \hat{g}(0) = \sum_l (\hat{b}(1, l) - \hat{b}(0, l)) \hat{f}_L(l) \\ &= \frac{\sum_{s=1}^n I(A_s=1) Y_s \hat{f}_{A|L}^{-1}(A_s|L_s)}{\sum_{s=1}^n I(A_s=1) \hat{f}_{A|L}^{-1}(A_s|L_s)} - \frac{\sum_{s=1}^n I(A_s=0) Y_s \hat{f}_{A|L}^{-1}(A_s|L_s)}{\sum_{s=1}^n I(A_s=0) \hat{f}_{A|L}^{-1}(A_s|L_s)} \end{aligned}$$

with asymptotic variance $var(\hat{\psi}) \neq var(\tilde{\psi})$

- In fact, it can be shown (check this) that $\hat{\psi}$ is still asymptotically efficient and therefore $var(\hat{\psi}) \leq var(\tilde{\psi})$

115

116

- ▶ But this is paradoxal as we find that the estimator $\tilde{\psi}$ that uses available information on the treatment mechanism, is less efficient than $\hat{\psi}$ which ignores this information
- ▶ However, as argued with the likelihood on slide (97), the treatment process is ancillary, that is the efficiency bound for estimating ψ is the same whether or not one knows the treatment mechanism.
- ▶ The paradox is resolved by realizing that $\tilde{\psi}$ uses the available information, but in a very inefficient way as it ignores the fact that L is a strong correlate of Y .

117

- ▶ In contrast, the G-formula and the iptw with nonparametrically estimated weights correctly incorporate this information.
- ▶ Take home messages:
 - ▶ More information is generally a good thing, but not helpful if used inefficiently.
 - ▶ Standard methods for analyzing clinical trials such as simple two sample z-test, may be highly inefficient as they ignore potentially highly informative baseline covariates.

118