Zhiwei Zhang*, Chunling Liu, Shujie Ma, and Min Zhang

# Estimating Mann–Whitney-Type Causal Effects for Right-Censored Survival Outcomes

**Abstract:** Mann–Whitney-type causal effects are clinically relevant, easy to interpret, and readily applicable to a wide range of study settings. This article considers estimation of such effects when the outcome variable is a survival time subject to right censoring. We derive and discuss several methods: an outcome regression method based on a regression model for the survival outcome, an inverse probability weighting method based on models for treatment assignment and censoring, and two doubly robust methods that involve both types of models and that remain valid under correct specification of the outcome model or the other two models. The methods are compared in a simulation study and applied to an observational study of hospitalized pneumonia.

## 1 Introduction

Comparative effectiveness research frequently involves comparing one treatment ($a = 1$) with another ($a = 0$) with respect to a clinical outcome of interest. For a generic subject in the target population, let $T(a)$, $a \in \{0, 1\}$, denote the potential outcome under treatment $a$ [1]. The causal effect of $a = 1$ versus $a = 0$ is commonly assessed by comparing the means of $T(1)$ and $T(0)$, if they are defined. Recently, there has been growing interest in a class of Mann–Whitney-type effect measures defined as $\theta = E\{h(T_i(1), T_j(0))\}$, where $h(\cdot, \cdot)$ is a specified function and the subscripts $i$ and $j$ denote two independent subjects. A simple example of $h$ is given by

$$h(t_1, t_0) = I(t_1 > t_0) - I(t_1 < t_0), \tag{1}$$

where $I(\cdot)$ is the indicator function; see, for example, Agresti [2, page 58]. Related definitions include $h(t_1, t_0) = I(t_1 > t_0) + 0.5I(t_1 = t_0)$, which forms the basis for the rank sum test [3] and the Mann–Whitney statistic [4]. For an ordered categorical outcome, which does not have a mean, the Mann–Whitney-type effect based on (1) provides a natural overall effect measure. Such effect measures have been recommended for clinical trials with arbitrary (ordinal or higher level) outcomes because of their clinical relevance and interpretability [e. g., 5, 6, 7]. They are particularly useful in analyzing an ordinal composite outcome that combines a quantitative outcome with death and possibly treatment discontinuation due to adverse events or lack of efficacy [8, 9]. As a possible drawback of such effect measures, Lumley [10] points out a lack of transitivity when $\theta$ based on (1) is used to order multiple treatments. On the other hand, for comparing two treatments, $\theta$ does provide unique insights that are not available from a comparison of means (if feasible at all).

In this article, we consider estimation of $\theta$ when the outcome of interest is a survival time subject to right censoring, which has not been considered in previous studies of Mann–Whitney-type causal effects. If $T$ is

*Corresponding author: Zhiwei Zhang,** University of California, Riverside, Department of Statistics, 900 University Ave, Riverside, United States, e-mail: zhiwei.zhang@ucr.edu
**Chunling Liu,** Hong Kong Polytechnic University, Department of Applied Mathematics, Hong Kong, China, e-mail: catherine.chunling.liu@polyu.edu.hk
**Shujie Ma,** University of California, Riverside, Department of Statistics, 900 University Ave, Riverside, United States, e-mail: shujie.ma@ucr.edu
**Min Zhang,** University of Michigan, Department of Biostatistics, Ann Arbor, United States, e-mail: mzhangst@umich.edu

a survival time, $\theta$ remains well defined and is just as meaningful as it is for another continuous outcome. In fact, the class of Mann–Whitney-type causal effects includes some well known and commonly used effect measures for survival outcomes. For example, if $h(t_1, t_0) = I(t_1 \leq t) - I(t_0 \leq t)$, then $\theta = F_1(t) - F_0(t)$, where $F_a(t) = P\{T(a) \leq t\}$, $a = 0, 1$. As another example, if $h(t_1, t_0) = t_1 \wedge \tau - t_0 \wedge \tau$, where $\wedge$ denotes minimum and $\tau$ is a positive constant chosen to ensure identifiability, then $\theta$ is the difference in mean restricted survival time. A new effect measure for survival outcomes is obtained by taking $h$ to be a truncated version of (1):

$$h(t_1, t_0) = I(t_1 \wedge \tau > t_0 \wedge \tau) - I(t_1 \wedge \tau < t_0 \wedge \tau), \tag{2}$$

for which $\theta = P\{T_i(1) \wedge \tau > T_j(0) \wedge \tau\} - P\{T_i(1) \wedge \tau < T_j(0) \wedge \tau\}$ can be interpreted as a win-lose probability difference in comparing the restricted survival times of two randomly chosen subjects who are randomly assigned to treatments 1 and 0.

Estimation of $\theta$ for a fully observed outcome has been considered by Chen et al. [11], Vermeulen et al. [12] and Zhang et al. [13], with focus on adjusting for confounders in observational studies and using auxiliary information to improve efficiency in randomized clinical trials. These methods cannot be used to estimate $\theta$ for a survival outcome subject to right censoring, which represents an additional challenge. We are not aware of any existing method for estimating $\theta$ for a right-censored survival outcome, except in the aforementioned special cases where $\theta$ is a difference in distribution function [14, 15, 16] or a difference in mean restricted survival time [17, 18, 19]. In both of these special cases, $h(t_1, t_0)$ is additive in the sense that it can be written as a function of $t_1$ minus a function of $t_0$. Here, we consider estimation of $\theta$ for a general function $h$, such as (2), that is not assumed to be additive. We derive and compare several methods: an outcome regression (OR) method based on a regression model for the survival outcome, an inverse probability weighting (IPW) method based on models for treatment assignment and censoring, and two doubly robust (DR) methods that involve both types of models and that remain valid under correct specification of the outcome model or the other two models. One of the DR methods is a straightforward extension of Zhang and Schaubel [19], and the other one is a new method based directly on the efficient influence function [20, 21] for estimating $\theta$.

In the next section, we describe these methods and discuss their asymptotic behavior. In Section 3, we report a simulation study and present a real application. Concluding remarks are given in Section 4. Technical details are provided in appendices.

# 2 Methodology

## 2.1 Notation and assumptions

Let $A$ denote the actual treatment received by an individual subject in a study, which may be a randomized clinical trial or an observational study. Assuming consistency or stable unit treatment value, the actual survival time is then $T = T(A) = AT(1) + (1 - A)T(0)$. Let $\boldsymbol{W}$ be a vector of relevant covariates measured at baseline (before treatment). We assume that treatment assignment is ignorable [22] given $\boldsymbol{W}$ in the sense that

$$P\{A = 1 | \boldsymbol{W}, T(0), T(1)\} = P(A = 1 | \boldsymbol{W}) =: p(\boldsymbol{W}). \tag{3}$$

The ignorability assumption is trivially true in a randomized clinical trial. In an observational study, the assumption requires that $\boldsymbol{W}$ contain enough information to explain any association between $A$ and the potential outcomes. We also assume positivity:

$$0 < p(\boldsymbol{W}) < 1 \quad \text{with probability 1}, \tag{4}$$

which is trivially true in a clinical trial but not trivial in an observational study.

The actual survival time $T$ may be right-censored by a censoring time $C$. We assume independent censoring:

$$C \perp T | (A, \boldsymbol{W}), \tag{5}$$

where $\perp$ denotes independence. The observed data for an individual subject can be represented as $\boldsymbol{O} = (\boldsymbol{W}, A, X, \Delta)$, where $X = T \wedge C$ and $\Delta = I(T \leq C)$. The observed data for the whole study will be conceptualized as independent copies of $\boldsymbol{O}$ and denoted by $\boldsymbol{O}_i = (\boldsymbol{W}_i, A_i, X_i, \Delta_i)$, $i = 1, \ldots, n$.

Our goal is to use the observed data to estimate $\theta$ for a general (specified) function $h$. For identifiability, we assume that $h(t_1, t_0)$ is determined by $t_1 \wedge \tau$ and $t_0 \wedge \tau$ for some $\tau > 0$ such that $P(C > \tau|A, \boldsymbol{W}) > 0$ almost surely. Accordingly, we will restrict attention to the interval $(0, \tau]$ in discussions of distribution, survival and hazard functions. For theoretical reasons, we also assume that $E\{h(T_i(1), T_j(0))^2\} < \infty$, which is satisfied by all examples mentioned earlier. The efficient influence function in this estimation problem is given in Appendix A. In the rest of this section, we propose several estimators of $\theta$, whose asymptotic properties are studied in Appendix B. Asymptotic variance formulas can be derived but may be cumbersome to use in variance estimation, depending on the specific forms of the working models involved. For ease of implementation, we use bootstrap standard errors and confidence intervals for inference.

## 2.2 Outcome regression (OR) estimator

The OR method is based on the fact that

$$\theta = h(F_1, F_0) = E\{h(F_1(\cdot|\boldsymbol{W}_i), F_0(\cdot|\boldsymbol{W}_j))\},$$

where $i \neq j$ and $F_a(t|\boldsymbol{W}) = P\{T(a) \leq t|\boldsymbol{W}\}$ for $a = 0, 1$. Here and in the sequel, we write $h(\nu_1, \nu_0) = \iint h(t_1, t_0)d\nu_1(t_1)d\nu_0(t_0)$, $h(\nu_1, t_0) = \int h(t_1, t_0)d\nu_1(t_1)$ and $h(t_1, \nu_0) = \int h(t_1, t_0)d\nu_0(t_0)$ for any measures $(\nu_1, \nu_0)$. Assumption (3) implies that $F_a(t|\boldsymbol{W}) = P(T \leq t|A = a, \boldsymbol{W})$, which, by assumptions (4) and (5), can be identified from the observed data.

To deal with the curse of dimensionality, we will assume a parametric or semiparametric model for $F_a(t|\boldsymbol{W})$, say $F_a(t|\boldsymbol{W}; \boldsymbol{\alpha})$, where $\boldsymbol{\alpha}$ may be finite- or infinite-dimensional. This will be referred to as the OR model. A prominent example of a semiparametric model for $F_a(t|\boldsymbol{W})$ is the proportional hazards model [23, 24]:

$$\lambda(t|A, \boldsymbol{W}; \boldsymbol{\alpha}) = \lambda_0(t)\exp(\boldsymbol{\eta}'\boldsymbol{V}),$$

where $\lambda(\cdot|A, \boldsymbol{W}; \boldsymbol{\alpha})$ is the conditional hazard function of $T$ given $(A, \boldsymbol{W})$, $\boldsymbol{\alpha} = (\boldsymbol{\eta}, \lambda_0)$, $\lambda_0$ is the baseline hazard function, and $\boldsymbol{V}$ is a vector-valued function of $(A, \boldsymbol{W})$. Under this model,

$$F_a(t|\boldsymbol{W}) = 1 - \exp\{-\Lambda(t|a, \boldsymbol{W}; \boldsymbol{\alpha})\} = 1 - \exp\{-\Lambda_0(t)\exp(\boldsymbol{\eta}'\boldsymbol{V})\},$$

where $\Lambda$ and $\Lambda_0$ are cumulative versions of $\lambda$ and $\lambda_0$, respectively.

Whether $F_a(t|\boldsymbol{W}; \boldsymbol{\alpha})$ is parametric or semiparametric, it is usually convenient to work with the corresponding hazard function $\lambda(t|A, \boldsymbol{W}; \boldsymbol{\alpha})$. The likelihood for $\boldsymbol{\alpha}$ based on the observed data may be written as

$$\prod_{i=1}^{n} \lambda(X_i|A_i, \boldsymbol{W}_i; \boldsymbol{\alpha})^{\Delta_i} \exp\{-\Lambda(X_i|A_i, \boldsymbol{W}_i; \boldsymbol{\alpha})\}.$$

At least for parametric models and the proportional hazards model, $\boldsymbol{\alpha}$ can be estimated by maximizing the above likelihood. Let $\widehat{\boldsymbol{\alpha}}$ be an estimate of $\boldsymbol{\alpha}$, which may be obtained by maximum likelihood or other means. Then $\theta$ can be estimated by

$$\widehat{\theta}_{\mathrm{or}} = \frac{1}{n(n-1)} \sum_{i \neq j} h(F_1(\cdot|\boldsymbol{W}_i; \widehat{\boldsymbol{\alpha}}), F_0(\cdot|\boldsymbol{W}_j; \widehat{\boldsymbol{\alpha}})).$$

If the model $F_a(t|\boldsymbol{W}; \boldsymbol{\alpha})$ is correct and $\widehat{\boldsymbol{\alpha}}$ is consistent for $\boldsymbol{\alpha}$ (in a suitable sense), then $\widehat{\theta}_{\mathrm{or}}$ is consistent for $\theta$ under mild regularity conditions. In Appendix B, we show that $\sqrt{n}(\widehat{\theta}_{\mathrm{or}} - \theta)$ converges to a zero-mean normal distribution under general conditions. A key condition we assume is that $\widehat{\boldsymbol{\alpha}}$ is $\sqrt{n}$-consistent and asymptotically linear in a suitable sense.

## 2.3 Inverse probability weighted (IPW) estimator

It follows from assumptions (3)–(5) and a conditioning argument that

$$
\theta = \mathrm{E}\left[\frac{A_i(1 - A_j)\Delta_i\Delta_j h(X_i, X_j)}{p(\boldsymbol{W}_i)\{1 - p(\boldsymbol{W}_j)\}S^c(X_i|A_i, \boldsymbol{W}_i)S^c(X_j|A_j, \boldsymbol{W}_j)}\right] \qquad (i \neq j),
$$

where $S^c(t|A, \boldsymbol{W}) = \mathrm{P}(C > t|A, \boldsymbol{W})$. If $S^c(t|A, \boldsymbol{W})$ and $p(\boldsymbol{W})$ are known, then the above identity suggests that $\theta$ can be estimated by

$$
\frac{1}{n(n-1)}\sum_{i \neq j}\frac{A_i(1 - A_j)\Delta_i\Delta_j h(X_i, X_j)}{p(\boldsymbol{W}_i)\{1 - p(\boldsymbol{W}_j)\}S^c(X_i|A_i, \boldsymbol{W}_i)S^c(X_j|A_j, \boldsymbol{W}_j)}.
$$

In reality, although $p(\boldsymbol{W})$ is known by design in a clinical trial, $S^c(t|A, \boldsymbol{W})$ is generally unknown and must be estimated. We assume a model for $S^c(t|A, \boldsymbol{W})$, say $S^c(t|A, \boldsymbol{W}; \boldsymbol{\beta})$, which may be parametric or semi-parametric so $\boldsymbol{\beta}$ may be finite- or infinite-dimensional. The model $S^c(t|A, \boldsymbol{W}; \boldsymbol{\beta})$ may be specified using the same considerations for specifying $F_a(t|\boldsymbol{W}; \boldsymbol{\alpha})$ except that we are now modeling the censoring time instead of the survival time. Let $\widehat{\boldsymbol{\beta}}$ be an estimate of $\boldsymbol{\beta}$, which may be obtained by maximizing the likelihood:

$$
\prod_{i=1}^{n}\lambda^c(X_i|A_i, \boldsymbol{W}_i; \boldsymbol{\beta})^{1-\Delta_i}\exp\{-\Lambda^c(X_i|A_i, \boldsymbol{W}_i; \boldsymbol{\beta})\},
$$

where $\lambda^c(\cdot|A, \boldsymbol{W}; \boldsymbol{\beta})$ is the conditional hazard function of $C$ given $(A, \boldsymbol{W})$ under the model $S^c(t|A, \boldsymbol{W}; \boldsymbol{\beta})$, and $\Lambda^c(\cdot|A, \boldsymbol{W}; \boldsymbol{\beta})$ is the cumulative version of $\lambda^c(\cdot|A, \boldsymbol{W}; \boldsymbol{\beta})$.

In an observational study, we also need to estimate $p(\boldsymbol{W})$. In fact, even if $p(\boldsymbol{W})$ is known, using an estimate of $p(\boldsymbol{W})$ instead of the known value usually leads to better efficiency in the IPW approach [25]. Let $p(\boldsymbol{W})$ be modeled as $p(\boldsymbol{W}; \boldsymbol{\gamma})$, where $\boldsymbol{\gamma}$ is a finite- or infinite-dimensional parameter. Because $A$ is binary, a typical choice for $p(\boldsymbol{W}; \boldsymbol{\gamma})$ would be a logistic regression model. Let $\widehat{\boldsymbol{\gamma}}$ be an estimate of $\boldsymbol{\gamma}$, which may be obtained by maximizing the likelihood:

$$
\prod_{i=1}^{n}p(\boldsymbol{W}_i; \boldsymbol{\gamma})^{A_i}\{1 - p(\boldsymbol{W}_i; \boldsymbol{\gamma})\}^{1-A_i}.
$$

Once $\widehat{\boldsymbol{\beta}}$ and $\widehat{\boldsymbol{\gamma}}$ are obtained, the IPW estimator of $\theta$ is readily available as a weighted average:

$$
\widehat{\theta}_{\mathrm{ipw}} = \frac{\sum_{i \neq j}\frac{A_i(1-A_j)\Delta_i\Delta_j h(X_i, X_j)}{p(\boldsymbol{W}_i;\widehat{\boldsymbol{\gamma}})\{1-p(\boldsymbol{W}_j;\widehat{\boldsymbol{\gamma}})\}S^c(X_i|A_i, \boldsymbol{W}_i;\widehat{\boldsymbol{\beta}})S^c(X_j|A_j, \boldsymbol{W}_j;\widehat{\boldsymbol{\beta}})}}{\sum_{i \neq j}\frac{A_i(1-A_j)\Delta_i\Delta_j}{p(\boldsymbol{W}_i;\widehat{\boldsymbol{\gamma}})\{1-p(\boldsymbol{W}_j;\widehat{\boldsymbol{\gamma}})\}S^c(X_i|A_i, \boldsymbol{W}_i;\widehat{\boldsymbol{\beta}})S^c(X_j|A_j, \boldsymbol{W}_j;\widehat{\boldsymbol{\beta}})}}.
$$

The denominator here serves to normalize the inverse probability weights. If the models $S^c(t|A, \boldsymbol{W}; \boldsymbol{\beta})$ and $p(\boldsymbol{W}; \boldsymbol{\gamma})$ are correct and $(\widehat{\boldsymbol{\beta}}, \widehat{\boldsymbol{\gamma}})$ estimate $(\boldsymbol{\beta}, \boldsymbol{\gamma})$ consistently, then $\widehat{\theta}_{\mathrm{ipw}}$ is consistent for $\theta$ and asymptotically linear under appropriate conditions (see Appendix B).

## 2.4 First doubly robust (DR1) estimator

The DR1 method is adapted from the DR method of Zhang and Schaubel [19] for estimating the difference in mean restricted survival time. As an important by-product, Zhang and Schaubel [19] propose a DR estimator of $F_a(t)$, which can be described as follows. Let us define

$$
N_i(t) = \Delta_i I(X_i \leq t),
$$
$$
Y_i(t) = I(X_i \geq t),
$$
$$
\widehat{\Lambda}_{ai}(t) = \Lambda(t|a, \boldsymbol{W}_i; \widehat{\boldsymbol{\alpha}}),
$$

$$\widehat{\Lambda}_{ai}^c(t) = \Lambda^c(t|a, \boldsymbol{W}_i; \widehat{\boldsymbol{\beta}}),$$

$$\widehat{M}_{ai}^c(t) = I(A_i = a)\{(1 - \Delta_i)I(X_i \le t) - \widehat{\Lambda}_{ai}^c(X_i \wedge t)\},$$

$$\widehat{G}_{ai}(t) = 1 - \int_0^t e^{\widehat{\Lambda}_{ai}(u) + \widehat{\Lambda}_{ai}^c(u)} d\widehat{M}_{ai}^c(u),$$

$$\omega_{ai} = I(A_i = a)p(\boldsymbol{W}_i; \widehat{\boldsymbol{\gamma}})^{-a}\{1 - p(\boldsymbol{W}_i; \widehat{\boldsymbol{\gamma}})\}^{a-1}.$$

Then the Zhang-Schaubel estimator of $\Lambda_a(t)$, the cumulative hazard function of $T(a)$, is given by

$$\widehat{\Lambda}_a^{\mathrm{dr}}(t) = \int_0^t \frac{\sum_{i=1}^n \left[\omega_{ai}e^{\widehat{\Lambda}_{ai}^c(u)}dN_i(u) + e^{-\widehat{\Lambda}_{ai}(u)}d\widehat{\Lambda}_{ai}(u)\{1 - \omega_{ai}\widehat{G}_{ai}(u)\}\right]}{\sum_{i=1}^n \left[\omega_{ai}e^{\widehat{\Lambda}_{ai}^c(u)}Y_i(u) + e^{-\widehat{\Lambda}_{ai}(u)}\{1 - \omega_{ai}\widehat{G}_{ai}(u)\}\right]},$$

and the corresponding estimator of $F_a(t)$ is $\widehat{F}_a^{\mathrm{dr}}(t) = 1 - \exp\{-\widehat{\Lambda}_a^{\mathrm{dr}}(t)\}$. The superscript "dr" in these estimators indicates that the estimators are doubly robust in the sense of being consistent and asymptotically normal if (i) $F_a(t|\boldsymbol{W}; \boldsymbol{\alpha})$ is correctly specified, (ii) $S^c(t|A, \boldsymbol{W}; \boldsymbol{\beta})$ and $p(\boldsymbol{W}; \boldsymbol{\gamma})$ are both correctly specified, or (iii) both (i) and (ii) hold.

Now $\theta$ can be estimated by $\widehat{\theta}_{\mathrm{dr1}} = h(\widehat{F}_1^{\mathrm{dr}}, \widehat{F}_0^{\mathrm{dr}})$, and it is easy to see that $\widehat{\theta}_{\mathrm{dr1}}$ is doubly robust in the same sense described above. Because the $\widehat{F}_a^{\mathrm{dr}}$ ($a = 0, 1$) are not guaranteed to be probability measures, $\widehat{\theta}_{\mathrm{dr1}}$ is not guaranteed to lie in the range of $h$. When $\widehat{\theta}_{\mathrm{dr1}}$ falls outside the range of $h$, it can be truncated into the range of $h$ without changing its consistency or asymptotic normality. It is not immediately clear whether $\widehat{\theta}_{\mathrm{dr1}}$ is locally efficient, that is, whether it attains the nonparametric information bound when all three working models are correct. This is the main motivation for our developing a second DR estimator of $\theta$ based directly on the efficient influence function for estimating $\theta$.

## 2.5 Second doubly robust (DR2) estimator

In Appendix A, we show that the efficient influence function for estimating $\theta$ is

$$
\begin{aligned}
\phi_{\mathrm{eff}}(\boldsymbol{O}) = {} & \frac{A\Delta h(X, F_0)}{p(\boldsymbol{W})S^c(X|A, \boldsymbol{W})} + \frac{(1 - A)\Delta h(F_1, X)}{\{1 - p(\boldsymbol{W})\}S^c(X|A, \boldsymbol{W})} - 2\theta \\
& + \{A - p(\boldsymbol{W})\}\left\{\frac{h(F_1, F_0(\cdot|\boldsymbol{W}))}{1 - p(\boldsymbol{W})} - \frac{h(F_1(\cdot|\boldsymbol{W}), F_0)}{p(\boldsymbol{W})}\right\} \\
& + \frac{A}{p(\boldsymbol{W})}\int \frac{\hbar_1(t, F_0|\boldsymbol{W})}{S^c(t|A, \boldsymbol{W})}dM^c(t) + \frac{1 - A}{1 - p(\boldsymbol{W})}\int \frac{\hbar_0(F_1, t|\boldsymbol{W})}{S^c(t|A, \boldsymbol{W})}dM^c(t),
\end{aligned}
\tag{6}
$$

where $\hbar_1(t, F_0|\boldsymbol{W}) = \mathrm{E}\{h(T, F_0)|A = 1, \boldsymbol{W}, T \ge t\}$, $\hbar_0(F_1, t|\boldsymbol{W}) = \mathrm{E}\{h(F_1, T)|A = 0, \boldsymbol{W}, T \ge t\}$, and $M^c(t) = (1 - \Delta)I(X \le t) - \Lambda^c(X \wedge t|A, \boldsymbol{W})$. Motivated by this result, we propose to estimate $\theta$ by setting the sample average of an empirical version of $\phi_{\mathrm{eff}}(\boldsymbol{O})$ to 0. The resulting estimator is

$$
\begin{aligned}
\widehat{\theta}_{\mathrm{dr2}} = \frac{1}{2n}\sum_{i=1}^n \Bigg[ & \frac{A_i\Delta_i h(X_i, \widehat{F}_0^{\mathrm{dr}})}{p(\boldsymbol{W}_i; \widehat{\boldsymbol{\gamma}})S^c(X_i|A_i, \boldsymbol{W}_i; \widehat{\boldsymbol{\beta}})} + \frac{(1 - A_i)\Delta_i h(\widehat{F}_1^{\mathrm{dr}}, X_i)}{\{1 - p(\boldsymbol{W}_i; \widehat{\boldsymbol{\gamma}})\}S^c(X_i|A_i, \boldsymbol{W}_i; \widehat{\boldsymbol{\beta}})} \\
& + \{A_i - p(\boldsymbol{W}_i; \widehat{\boldsymbol{\gamma}})\}\left\{\frac{h(\widehat{F}_1^{\mathrm{dr}}, F_0(\cdot|\boldsymbol{W}_i; \widehat{\boldsymbol{\alpha}}))}{1 - p(\boldsymbol{W}_i; \widehat{\boldsymbol{\gamma}})} - \frac{h(F_1(\cdot|\boldsymbol{W}_i; \widehat{\boldsymbol{\alpha}}), \widehat{F}_0^{\mathrm{dr}})}{p(\boldsymbol{W}_i; \widehat{\boldsymbol{\gamma}})}\right\} \\
& + \frac{A_i}{p(\boldsymbol{W}_i; \widehat{\boldsymbol{\gamma}})}\int \frac{\hbar_1(t, \widehat{F}_0^{\mathrm{dr}}|\boldsymbol{W}_i; \widehat{\boldsymbol{\alpha}})}{S^c(t|A_i, \boldsymbol{W}_i; \widehat{\boldsymbol{\beta}})}d\widehat{M}_{1i}^c(t) \\
& + \frac{1 - A_i}{1 - p(\boldsymbol{W}_i; \widehat{\boldsymbol{\gamma}})}\int \frac{\hbar_0(\widehat{F}_1^{\mathrm{dr}}, t|\boldsymbol{W}_i; \widehat{\boldsymbol{\alpha}})}{S^c(t|A_i, \boldsymbol{W}_i; \widehat{\boldsymbol{\beta}})}d\widehat{M}_{0i}^c(t) \Bigg],
\end{aligned}
$$

where

$$\hbar_1(t, \widehat{F}_0^{\mathrm{dr}}|\boldsymbol{W}_i; \widehat{\boldsymbol{\alpha}}) = \frac{\int_t^\infty h(u, \widehat{F}_0^{\mathrm{dr}}) dF_1(u|\boldsymbol{W}_i; \widehat{\boldsymbol{\alpha}})}{1 - F_1(t|\boldsymbol{W}_i; \widehat{\boldsymbol{\alpha}})},$$

$$\hbar_0(\widehat{F}_1^{\mathrm{dr}}, t|\boldsymbol{W}_i; \widehat{\boldsymbol{\alpha}}) = \frac{\int_t^\infty h(\widehat{F}_0^{\mathrm{dr}}, u) dF_0(u|\boldsymbol{W}_i; \widehat{\boldsymbol{\alpha}})}{1 - F_0(t|\boldsymbol{W}_i; \widehat{\boldsymbol{\alpha}})}.$$

The use of $\widehat{F}_a^{\mathrm{dr}}$ in $\widehat{\theta}_{\mathrm{dr}2}$ to replace $F_a$ in (6) is an attempt to provide maximal protection against model misspecification. In Appendix B, we show that $\widehat{\theta}_{\mathrm{dr}2}$ is indeed doubly robust and, furthermore, locally efficient. Like $\widehat{\theta}_{\mathrm{dr}1}$, $\widehat{\theta}_{\mathrm{dr}2}$ may fall outside the range of $h$ but can be truncated back into the range of $h$.

# 3 Numerical results

## 3.1 Simulation

Here we report a simulation study of the finite-sample performance of the methods described in Section 2. Data are generated according to the following mechanism:

$$\boldsymbol{W} = (W_1, W_2)' \sim N(\boldsymbol{0}, \boldsymbol{I}),$$
$$p(\boldsymbol{W}) = \mathrm{expit}(W_1 - W_2),$$
$$\lambda(t|A, \boldsymbol{W}) = \exp(\eta_A A - W_1 + W_2),$$
$$\lambda^c(t|A, \boldsymbol{W}) = \exp(W_1 + W_2),$$

where $\boldsymbol{0} = (0, 0)'$, $\boldsymbol{I}$ is the 2-by-2 identity matrix, $\mathrm{expit}(x) = 1/\{1 + \exp(-x)\}$, and $\eta_A$ is either 0 (no treatment effect) or $-0.5$ (protective treatment effect). The function $h$ in the definition of $\theta$ is take to be (2) with $\tau = 2$. The true value of $\theta$ is 0 when $\eta_A = 0$ and approximately 0.15 when $\eta_A = -0.5$. We consider two sample sizes: $n = 200, 500$. In each case, 1,000 replicate samples are generated.

Each simulated sample is analyzed using the four methods (OR, IPW, DR1 and DR2) described in Section 2 as well as a naive method that ignores censoring and possible confounding and estimates $\theta$ with the $U$-statistic

$$\widehat{\theta}_{\mathrm{nv}} = \frac{1}{n_0 n_1} \sum_{i=1}^n \sum_{j=1}^n A_i(1 - A_j) h(X_i, X_j),$$

where $n_1 = \sum_{i=1}^n A_i$ and $n_0 = n - n_1$. The OR, IPW, DR1 and DR2 methods are implemented with correct and incorrect working models. The correct models are:

$$\lambda(t|A = a, \boldsymbol{W}; \boldsymbol{\alpha}_a) = \lambda_{a0}(t) \exp\left((W_1, W_2)\boldsymbol{\eta}_a\right),$$
$$\lambda^c(t|A = a, \boldsymbol{W}; \boldsymbol{\beta}_a) = \lambda_{a0}^c(t) \exp\left((W_1, W_2)\boldsymbol{\eta}_a^c\right),$$
$$p(\boldsymbol{W}; \boldsymbol{\gamma}) = \mathrm{expit}\left((1, W_1, W_2)\boldsymbol{\gamma}\right),$$

where $\boldsymbol{\alpha}_a = (\boldsymbol{\eta}_a, \lambda_{a0})$, $\boldsymbol{\beta}_a = (\boldsymbol{\eta}_a^c, \lambda_{a0}^c)$, and $\lambda_a$ and $\lambda_a^c$ are unspecified baseline hazard functions. The incorrect models result from replacing $(W_1, W_2)$ with $(I(W_1 > 0), I(W_2 > 0))$ in the correct models. All models are estimated using the maximum likelihood approach.

Table 1 shows a summary of the simulation results: empirical bias and standard deviation for estimating $\theta$. As expected, the naive method is severely biased. The OR estimator becomes biased when the model $\lambda(t|A, \boldsymbol{W}; \boldsymbol{\alpha})$ is misspecified, as does the IPW estimator when the models $\lambda^c(t|A, \boldsymbol{W}; \boldsymbol{\beta})$ and $p(\boldsymbol{W}; \boldsymbol{\gamma})$ are misspecified. In contrast, the two DR estimators are nearly unbiased unless all models are misspecified, demonstrating double robustness. Regardless of model (in)correctness, the estimators that adjust for confounding and censoring generally follow the pattern OR < DR < IPW in terms of variability. The efficiency comparison of the two DR estimators does not seem to follow a clear pattern. At $n = 200$, DR2 appears more efficient than DR1 when all models are correct, but this difference appears to diminish with increasing sample size.

**Table 1:** Simulation results: empirical bias and standard deviation (SD) for the naive, OR, IPW, DR1 and DR2 estimators with correct and incorrect working models (see Section 3.1 for details).

| Method | Models | | | $\theta = 0$ | | $\theta \approx 0.15$ | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| | $(T\|A, W)$ | $(C\|A, W)$ | $(A\|W)$ | Bias | SD | Bias | SD |
| $n = 200$ | | | | | | | |
| naive | | | | 0.255 | 0.077 | 0.173 | 0.078 |
| OR | correct | | | 0.001 | 0.080 | −0.011 | 0.087 |
| OR | incorrect | | | 0.157 | 0.097 | 0.141 | 0.103 |
| IPW | | correct | correct | 0.048 | 0.174 | 0.014 | 0.177 |
| IPW | | incorrect | incorrect | 0.180 | 0.190 | 0.142 | 0.197 |
| DR1 | correct | correct | correct | 0.008 | 0.138 | 0.001 | 0.146 |
| DR1 | correct | incorrect | incorrect | −0.004 | 0.100 | −0.013 | 0.106 |
| DR1 | incorrect | correct | correct | 0.017 | 0.143 | 0.010 | 0.136 |
| DR1 | incorrect | incorrect | incorrect | 0.155 | 0.120 | 0.146 | 0.112 |
| DR2 | correct | correct | correct | 0.008 | 0.118 | −0.045 | 0.116 |
| DR2 | correct | incorrect | incorrect | −0.018 | 0.105 | −0.056 | 0.117 |
| DR2 | incorrect | correct | correct | 0.014 | 0.120 | −0.038 | 0.108 |
| DR2 | incorrect | incorrect | incorrect | 0.138 | 0.122 | 0.096 | 0.121 |
| $n = 500$ | | | | | | | |
| naive | | | | 0.256 | 0.050 | 0.172 | 0.048 |
| OR | correct | | | 0.003 | 0.051 | −0.009 | 0.055 |
| OR | incorrect | | | 0.162 | 0.061 | 0.147 | 0.064 |
| IPW | | correct | correct | 0.044 | 0.141 | 0.011 | 0.147 |
| IPW | | incorrect | incorrect | 0.163 | 0.159 | 0.137 | 0.173 |
| DR1 | correct | correct | correct | 0.014 | 0.082 | 0.002 | 0.104 |
| DR1 | correct | incorrect | incorrect | −0.008 | 0.068 | −0.019 | 0.077 |
| DR1 | incorrect | correct | correct | 0.018 | 0.077 | 0.007 | 0.088 |
| DR1 | incorrect | incorrect | incorrect | 0.157 | 0.080 | 0.146 | 0.078 |
| DR2 | correct | correct | correct | 0.017 | 0.077 | −0.032 | 0.107 |
| DR2 | correct | incorrect | incorrect | −0.020 | 0.099 | −0.046 | 0.140 |
| DR2 | incorrect | correct | correct | 0.020 | 0.078 | −0.028 | 0.091 |
| DR2 | incorrect | incorrect | incorrect | 0.146 | 0.110 | 0.117 | 0.141 |

## 3.2 Application

We now apply the methods compared in Section 3.1 to a study of hospitalized pneumonia in young children, described in Section 1.13 of Klein and Moeschberger [26]. The study is based on 3,470 newborn children from the National Longitudinal Survey of Youth [27]. The research question is whether the mother's feeding choice (breast feeding or not) protects the infant against hospitalized pneumonia in the first year of life. The outcome variable of interest is the time to hospitalized pneumonia, which is restricted to one year and possibly censored before one year. The available baseline covariates are infant race (black, white or other), indicators of normal birthweight (at least 5.5 pounds) and having at least one sibling, and some maternal and family characteristics: age, years of schooling, alcohol use, cigarette use, region of the country (Northeast, Northcentral, South and West), poverty, and urban environment.

  We are interested in the causal effect of breast feeding on hospitalized pneumonia as measured by $\theta$ with $h$ defined by (2). This effect measure can be interpreted as a win-lose probability difference in a hypothetical comparison of the restricted times to hospitalized pneumonia of two randomly chosen infants, who are randomly assigned to breast feeding and no breast feeding as in a clinical trial. The propensity score model in our analysis is a logistic regression model based on all infant characteristics as well as maternal age, cigarette use, and years of schooling. The OR model is specified as a proportional hazards model based on infant characteristics, maternal cigarette use, and urban environment, with separate parameters for each treatment group. The model for censoring is also a separate proportional hazards model for each treatment group, with infant race and maternal age as covariates.

**Table 2:** Analysis of pneumonia data: point estimates (PE) of $\theta$ and bootstrap standard errors (SE) obtained using different methods (see Section 3.2 for details).

| Method | PE | SE |
|--------|-------:|-------:|
| naive | −0.0098 | 0.0169 |
| OR | −0.0020 | 0.0015 |
| IPW | 0.1373 | 0.1902 |
| DR1 | −0.0028 | 0.0018 |
| DR2 | 0.0076 | 0.0023 |

Table 2 shows the point estimates of $\theta$ obtained using the five methods compared in Section 3.1 together with nonparametric bootstrap standard errors based on 1,000 bootstrap samples. The IPW method clearly stands out with a large positive point estimate and a large standard error, both of which are likely due to the large variability of the IPW estimator. Considering the large standard error, the IPW estimate of $\theta$ may well be due to random variation and does not establish a positive effect of breast feeding. The naive method, on the other hand, has the smallest (i. e., most negative) point estimate and also a relatively large standard error. The results are somewhat similar for the other three methods (OR, DR1 and DR2). DR2 is the only method that reaches statistical significance (with or without adjusting for multiplicity); however, the estimated effect is rather small in magnitude under the DR2 method. Taken together, the point estimates and standard errors in Table 2 provide no clear evidence for a substantial effect of breast feeding on hospitalized pneumonia.

# 4 Discussion

Mann–Whitney-type causal effects are clinically relevant, easy to interpret, and readily applicable to a wide range of study settings. This article considers estimation of such effects when the outcome variable is a survival time subject to right censoring. We have derived four methods for doing this and compared them theoretically and empirically. Among these methods, the OR and DR methods have unique advantages, and their suitability to a given application will depend on the available information and the relative importance of bias versus efficiency. If one is primarily concerned about bias, then the DR methods may be preferable. The OR method is more efficient when the OR model is correctly specified. No clear advantage has been observed for the IPW method.

The relationship between the two DR estimators remains an open question. The DR1 estimator adapted from Zhang and Schaubel [19] is known to be doubly robust. The proposed DR2 estimator, based directly on the efficient influence function for estimating $\theta$, is doubly robust and locally efficient. In our simulation study, the two estimators appear to perform similarly in large samples when all working models are correctly specified, which suggests that the DR1 estimator might be locally efficient as well. Further research is warranted to clarify how the two DR estimators might relate to each other.

All of these methods assume that treatment assignment is ignorable in the sense of (3). The ignorability assumption cannot be validated with observed data and must be based on background knowledge. If there is insufficient background knowledge to justify the ignorability assumption, it is important to assess the robustness of study results in the presence of unmeasured confounders. Methods for conducting such a sensitivity analysis have been developed for some causal effects [e. g., 28, 29] but not for the Mann–Whitney-type effects considered here. Further research is needed to close this methodological gap.

# Appendix A.  Semiparametric theory

Here we derive the efficient influence function for estimating $\theta$ using the monotone coarsening argument of Tsiatis [21]. Recall that the observed data for an individual subject is $\boldsymbol{O} = (\boldsymbol{W}, A, X, \Delta)$, which is considered a coarsened version of the "full data" $\boldsymbol{Z} = (\boldsymbol{W}, A, T)$. Unless otherwise stated, the distributions of $\boldsymbol{Z}$ and $\boldsymbol{O}$ are restricted by assumptions (3)–(5) only. The full-data Hilbert space $\mathcal{H}^F$ is the space of functions of $\boldsymbol{Z}$ with zero mean and finite variance, equipped with the covariance inner product. The observed-data Hilbert space $\mathcal{H}$ is the space of functions of $\boldsymbol{O}$ with zero mean and finite variance, equipped with the covariance inner product.

First, we characterize the influence function of any regular, asymptotically linear (RAL) estimator of $\theta$ based on the full data. For argument's sake, we assume for a moment that $p(\boldsymbol{W}) = \mathrm{P}(A = 1|\boldsymbol{W})$ is known, as in a randomized clinical trial. In this special case, a nonparametric RAL estimator of $\theta$ based on the full data is given by

$$\frac{1}{n(n-1)} \sum_{i \neq j} \frac{A_i(1 - A_j)h(T_i, T_j)}{p(\boldsymbol{W}_i)\{1 - p(\boldsymbol{W}_j)\}}.$$

Using the theory of $U$-statistics [e. g., 30, Chapter 12], it is easy to see that the above estimator is RAL with influence function

$$\phi_0^F(\boldsymbol{Z}) = \frac{Ah(T, F_0)}{p(\boldsymbol{W})} + \frac{(1 - A)h(F_1, T)}{1 - p(\boldsymbol{W})} - 2\theta.$$

With $p(\boldsymbol{W})$ known, the orthogonal complement of the full-data tangent space is

$$\Psi_1 = \left\{ \{A - p(\boldsymbol{W})\}b(\boldsymbol{W}) : \mathrm{E}\{b(\boldsymbol{W})^2\} < \infty \right\},$$

and the influence function of any full-data RAL estimator of $\theta$ takes the form

$$\phi_b^F(\boldsymbol{Z}) = \phi_0^F(\boldsymbol{Z}) + \{A - p(\boldsymbol{W})\}b(\boldsymbol{W})$$

for some function $b$ such that $\mathrm{E}\{b(\boldsymbol{W})^2\} < \infty$. Now we remove the assumption that $p(\boldsymbol{W})$ is known. An RAL estimator of $\theta$ in this larger model is certainly an RAL estimator of $\theta$ in the smaller model with $p(\boldsymbol{W})$ known. Therefore, the influence function of an RAL estimator of $\theta$ in the larger model must also belong to $\phi_0^F(\boldsymbol{Z}) + \Psi_1$.

Next, we characterize the influence function of any RAL estimator of $\theta$ based on the observed data. According to Tsiatis [21, Chapter 9], such an influence function can be represented as $\phi_b(\boldsymbol{O}) + \psi(\boldsymbol{O})$, where $\phi_b(\boldsymbol{O})$ satisfies

$$\mathrm{E}\{\phi_b(\boldsymbol{O})|\boldsymbol{Z}\} = \phi_b^F(\boldsymbol{Z}) \tag{A.1}$$

and $\psi(\boldsymbol{O})$ is an element of $\Psi_2$, the augmentation space for censoring. It is easy to verify that equation (A.1) is satisfied by $\phi_b(\boldsymbol{O}) = \phi_0(\boldsymbol{O}) + \{A - p(\boldsymbol{W})\}b(\boldsymbol{W})$, where

$$\phi_0(\boldsymbol{O}) = \frac{A\Delta h(X, F_0)}{p(\boldsymbol{W})S^c(X|A, \boldsymbol{W})} + \frac{(1 - A)\Delta h(F_1, X)}{\{1 - p(\boldsymbol{W})\}S^c(X|A, \boldsymbol{W})} - 2\theta.$$

As shown in Tsiatis [21, page 217], the augmentation space $\Psi_2$ consists of martingale integrals of the form

$$\int \frac{q(t|A, \boldsymbol{W})}{S^c(t|A, \boldsymbol{W})}dM^c(t),$$

where $q$ is an arbitrary function, $M^c(t) = (1 - \Delta)I(X \leq t) - \Lambda^c(X \wedge t|A, \boldsymbol{W})$, and $S^c(\cdot|A, \boldsymbol{W})$ and $\Lambda^c(\cdot|A, \boldsymbol{W})$ are, respectively, the conditional survival and cumulative hazard functions of $C$ given $(A, \boldsymbol{W})$. Thus, the influence function of any RAL estimator of $\theta$ based on the observed data must be an element of $\phi_0(\boldsymbol{O}) + \Psi_1 + \Psi_2$, which can be written as

$$\phi_{b,q}(\boldsymbol{O}) = \phi_0(\boldsymbol{O}) + \{A - p(\boldsymbol{W})\}b(\boldsymbol{W}) + \int \frac{q(t|A, \boldsymbol{W})}{S^c(t|A, \boldsymbol{W})}dM^c(t)$$

for some functions $b$ and $q$.

To find the efficient influence function, the influence function with the smallest possible variance, we need to minimize the variance of $\phi_{b,q}(\boldsymbol{O})$ with respect to $b$ and $q$, which amounts to projecting $\phi_0(\boldsymbol{O})$ into the orthogonal complement of $\Psi_1 + \Psi_2$. To this end, we note that $\Psi_1$ and $\Psi_2$ are orthogonal to each other, so it suffices to project $\phi_0(\boldsymbol{O})$ into $\Psi_1$ and $\Psi_2$ separately. The projection of $\phi_0(\boldsymbol{O})$ into $\Psi_1$ is easily seen to be

$$\Pi(\phi_0(\boldsymbol{O})|\Psi_1) = \mathrm{E}\left(\phi_0(\boldsymbol{O})|A, \boldsymbol{W}\right) - \mathrm{E}\left(\phi_0(\boldsymbol{O})|\boldsymbol{W}\right)$$
$$= \{A - p(\boldsymbol{W})\}\left\{\frac{h(F_1(\cdot|\boldsymbol{W}), F_0)}{p(\boldsymbol{W})} - \frac{h(F_1, F_0(\cdot|\boldsymbol{W}))}{1 - p(\boldsymbol{W})}\right\}.$$

Section 10.4 of Tsiatis [21] indicates that the projection of $\phi_0(\boldsymbol{O})$ into $\Psi_2$ is

$$\Pi(\phi_0(\boldsymbol{O})|\Psi_2) = -\frac{A}{p(\boldsymbol{W})}\int \frac{\mathrm{E}\{h(T, F_0)|A, \boldsymbol{W}, T \geq t\}}{S^c(t|A, \boldsymbol{W})} dM^c(t)$$
$$- \frac{1 - A}{1 - p(\boldsymbol{W})}\int \frac{\mathrm{E}\{h(F_1, T)|A, \boldsymbol{W}, T \geq t\}}{S^c(t|A, \boldsymbol{W})} dM^c(t).$$

The efficient influence function is therefore

$$\phi_{\mathrm{eff}}(\boldsymbol{O}) = \phi_0(\boldsymbol{O}) - \Pi(\phi_0(\boldsymbol{O})|\Psi_1) - \Pi(\phi_0(\boldsymbol{O})|\Psi_2),$$

which is equivalent to equation (6) in the main paper.

# Appendix B.  Asymptotic theory

Standard regularity conditions in the $M$-estimation theory [e. g., 30, Chapter 5] are assumed. These include identifiability and smoothness (in parameters) of working models, existence of integrable envelopes that permit use of the dominated convergence theorem, and certain Donsker properties that help deal with random functions. Techniques for verifying the Donsker property can be found in van der Vaart and Wellner [31].

Let $\mathbb{P}_0$ denote the true distribution of $\boldsymbol{O} = (\boldsymbol{W}, A, X, \Delta)$, and let $\mathbb{P}_n$ denote the empirical distribution of $\boldsymbol{O}_i = (\boldsymbol{W}_i, A_i, X_i, \Delta_i)$, $i = 1, \ldots, n$. Write $\mathbb{G}_n = \sqrt{n}(\mathbb{P}_n - \mathbb{P}_0)$ for the empirical processes. For any measures $(v_1, v_0)$, we write

$$h(v_1, y_0) = \int h(y_1, y_0) dv_1(y_1),$$
$$h(y_1, v_0) = \int h(y_1, y_0) dv_0(y_0),$$
$$h(v_1, v_0) = \iint h(y_1, y_0) dv_1(y_1) dv_0(y_0).$$

## B.1  Asymptotics for $\widehat{\boldsymbol{\theta}}_{\mathrm{or}}$

Here we assume that the model $\lambda(t|A, \boldsymbol{W}; \boldsymbol{\alpha})$ is correct and that $\boldsymbol{\alpha}$ and $\widehat{\boldsymbol{\alpha}}$ take values in a suitable Banach space with

$$\sqrt{n}(\widehat{\boldsymbol{\alpha}} - \boldsymbol{\alpha}) = \mathbb{G}_n\boldsymbol{\phi}_{\boldsymbol{\alpha}}(\boldsymbol{O}) + o_p(1).$$

For any $\boldsymbol{a}$ in the space for $\boldsymbol{\alpha}$, we define

$$k_{\mathrm{or}}(\boldsymbol{W}_i, \boldsymbol{W}_j; \boldsymbol{a}) = \{h(F_1(\cdot|\boldsymbol{W}_i; \boldsymbol{a}), F_0(\cdot|\boldsymbol{W}_j; \boldsymbol{a})) + h(F_1(\cdot|\boldsymbol{W}_j; \boldsymbol{a}), F_0(\cdot|\boldsymbol{W}_i; \boldsymbol{a}))\}/2,$$
$$U_{\mathrm{or}}(\boldsymbol{a}) = \frac{2}{n(n-1)}\sum_{i<j} k_{\mathrm{or}}(\boldsymbol{W}_i, \boldsymbol{W}_j; \boldsymbol{a}),$$
$$u_{\mathrm{or}}(\boldsymbol{a}) = \mathrm{E}\{U_{\mathrm{or}}(\boldsymbol{a})\} = \mathrm{E}\{k_{\mathrm{or}}(\boldsymbol{W}_i, \boldsymbol{W}_j; \boldsymbol{a})\}, \qquad i \neq j.$$

Then $\widehat{\theta}_{\mathrm{or}} = U_{\mathrm{or}}(\widehat{\boldsymbol{\alpha}})$, $\theta = u_{\mathrm{or}}(\boldsymbol{\alpha})$, and

$$\sqrt{n}(\widehat{\theta}_{\mathrm{or}} - \theta) = \sqrt{n}\{U_{\mathrm{or}}(\widehat{\boldsymbol{\alpha}}) - u_{\mathrm{or}}(\boldsymbol{\alpha})\} = \sqrt{n}\{U_{\mathrm{or}}(\widehat{\boldsymbol{\alpha}}) - u_{\mathrm{or}}(\widehat{\boldsymbol{\alpha}})\} + \sqrt{n}\{u_{\mathrm{or}}(\widehat{\boldsymbol{\alpha}}) - u_{\mathrm{or}}(\boldsymbol{\alpha})\}. \tag{B.1}$$

It follows from the theory of $U$-statistics [e. g., 30, Theorem 12.3] that, for any fixed $\boldsymbol{a}$,

$$\sqrt{n}\{U_{\mathrm{or}}(\boldsymbol{a}) - u_{\mathrm{or}}(\boldsymbol{a})\} = \mathbb{G}_n\{2k'_{\mathrm{or}}(\boldsymbol{W};\boldsymbol{a}) - 2\theta\} + o_p(1), \tag{B.2}$$

where $k'_{\mathrm{or}}(\boldsymbol{w};\boldsymbol{a}) = \mathrm{E}\{k_{\mathrm{or}}(\boldsymbol{w},\boldsymbol{W};\boldsymbol{a})\}$. It can be argued as in Nolan and Pollard [32, proof of Theorem 5] that the $o_p(1)$ term in (B.2) is uniformly negligible for $\boldsymbol{a}$ in a neighborhood of $\boldsymbol{\alpha}$. This, together with Theorem 19.24 of van der Vaart [30], implies that

$$\sqrt{n}\{U_{\mathrm{or}}(\widehat{\boldsymbol{\alpha}}) - u_{\mathrm{or}}(\widehat{\boldsymbol{\alpha}})\} = \mathbb{G}_n\{2k'_{\mathrm{or}}(\boldsymbol{W};\widehat{\boldsymbol{\alpha}}) - 2\theta\} + o_p(1) = \mathbb{G}_n\{2k'_{\mathrm{or}}(\boldsymbol{W};\boldsymbol{\alpha}) - 2\theta\} + o_p(1). \tag{B.3}$$

We assume that the map $\boldsymbol{a} \mapsto u_{\mathrm{or}}(\boldsymbol{a})$ is differentiable at $\boldsymbol{\alpha}$ with derivative $\mathbf{D}_{\alpha}^{\mathrm{or}}$. By the delta method,

$$\sqrt{n}\{u_{\mathrm{or}}(\widehat{\boldsymbol{\alpha}}) - u_{\mathrm{or}}(\boldsymbol{\alpha})\} = \mathbf{D}_{\alpha}^{\mathrm{or}}\mathbb{G}_n\boldsymbol{\phi}_{\alpha}(\boldsymbol{O}) + o_p(1). \tag{B.4}$$

Substituting (B.3) and (B.4) into (B.1) yields

$$\sqrt{n}(\widehat{\theta}_{\mathrm{or}} - \theta) = \mathbb{G}_n\left\{2k'_{\mathrm{or}}(\boldsymbol{W};\boldsymbol{\alpha}) - 2\theta + \mathbf{D}_{\alpha}^{\mathrm{or}}\boldsymbol{\phi}_{\alpha}(\boldsymbol{O})\right\} + o_p(1).$$

## B.2 Asymptotics for $\widehat{\theta}_{\mathrm{ipw}}$

Here we assume that the models $\lambda^c(t|\boldsymbol{W};\boldsymbol{\beta})$ and $p(\boldsymbol{W};\boldsymbol{\gamma})$ are correct and that $\boldsymbol{\beta}$ and $\boldsymbol{\gamma}$ take values in suitable Banach spaces with

$$\sqrt{n}(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}) = \mathbb{G}_n\boldsymbol{\phi}_{\beta}(\boldsymbol{O}) + o_p(1),$$
$$\sqrt{n}(\widehat{\boldsymbol{\gamma}} - \boldsymbol{\gamma}) = \mathbb{G}_n\boldsymbol{\phi}_{\gamma}(\boldsymbol{O}) + o_p(1).$$

Let us define

$$k_1^{\dagger}(\boldsymbol{O}_i, \boldsymbol{O}_j; \boldsymbol{b}, \boldsymbol{g}) = \frac{A_i(1 - A_j)\Delta_i\Delta_j h(X_i, X_j)}{p(\boldsymbol{W}_i;\boldsymbol{g})\{1 - p(\boldsymbol{W}_j;\boldsymbol{g})\}S^c(X_i|A_i, \boldsymbol{W}_i;\boldsymbol{b})S^c(X_j|A_j, \boldsymbol{W}_j;\boldsymbol{b})},$$
$$k_2^{\dagger}(\boldsymbol{O}_i, \boldsymbol{O}_j; \boldsymbol{b}, \boldsymbol{g}) = \frac{A_i(1 - A_j)\Delta_i\Delta_j}{p(\boldsymbol{W}_i;\boldsymbol{g})\{1 - p(\boldsymbol{W}_j;\boldsymbol{g})\}S^c(X_i|A_i, \boldsymbol{W}_i;\boldsymbol{b})S^c(X_j|A_j, \boldsymbol{W}_j;\boldsymbol{b})},$$

and, for $r = 1, 2$,

$$k_r(\boldsymbol{O}_i, \boldsymbol{O}_j; \boldsymbol{b}, \boldsymbol{g}) = \{k_r^{\dagger}(\boldsymbol{O}_i, \boldsymbol{O}_j; \boldsymbol{b}, \boldsymbol{g}) + k_r^{\dagger}(\boldsymbol{O}_j, \boldsymbol{O}_i; \boldsymbol{b}, \boldsymbol{g})\}/2,$$
$$U_r(\boldsymbol{b}, \boldsymbol{g}) = \frac{2}{n(n-1)}\sum_{i<j} k_r(\boldsymbol{O}_i, \boldsymbol{O}_j; \boldsymbol{b}, \boldsymbol{g}),$$
$$u_r(\boldsymbol{b}, \boldsymbol{g}) = \mathrm{E}\{U_r(\boldsymbol{b}, \boldsymbol{g})\} = \mathrm{E}\{k_r(\boldsymbol{O}_i, \boldsymbol{O}_j; \boldsymbol{b}, \boldsymbol{g})\}, \qquad i \neq j.$$

Then $\widehat{\theta}_{\mathrm{ipw}} = U_1(\widehat{\boldsymbol{\beta}}, \widehat{\boldsymbol{\gamma}})/U_2(\widehat{\boldsymbol{\beta}}, \widehat{\boldsymbol{\gamma}})$, $u_1(\boldsymbol{\beta}, \boldsymbol{\gamma}) = \theta$, and $u_2(\boldsymbol{\beta}, \boldsymbol{\gamma}) = 1$. For $r = 1, 2$, it can be argued as before that

$$\sqrt{n}\{U_r(\widehat{\boldsymbol{\beta}}, \widehat{\boldsymbol{\gamma}}) - u_r(\boldsymbol{\beta}, \boldsymbol{\gamma})\} = \sqrt{n}\{U_r(\widehat{\boldsymbol{\beta}}, \widehat{\boldsymbol{\gamma}}) - u_r(\widehat{\boldsymbol{\beta}}, \widehat{\boldsymbol{\gamma}})\} + \sqrt{n}\{u_r(\widehat{\boldsymbol{\beta}}, \widehat{\boldsymbol{\gamma}}) - u_r(\boldsymbol{\beta}, \boldsymbol{\gamma})\}$$
$$= \sqrt{n}\{U_r(\boldsymbol{\beta}, \boldsymbol{\gamma}) - u_r(\boldsymbol{\beta}, \boldsymbol{\gamma})\} + \sqrt{n}\{u_r(\widehat{\boldsymbol{\beta}}, \widehat{\boldsymbol{\gamma}}) - u_r(\boldsymbol{\beta}, \boldsymbol{\gamma})\} + o_p(1)$$
$$= \mathbb{G}_n k'_r(\boldsymbol{O}; \boldsymbol{\beta}, \boldsymbol{\gamma}) + \sqrt{n}\{u_r(\widehat{\boldsymbol{\beta}}, \widehat{\boldsymbol{\gamma}}) - u_r(\boldsymbol{\beta}, \boldsymbol{\gamma})\} + o_p(1),$$

where $k'_r(\boldsymbol{o}; \boldsymbol{b}, \boldsymbol{g}) = \mathrm{E}\{k_r(\boldsymbol{o}, \boldsymbol{O}; \boldsymbol{b}, \boldsymbol{g})\}$. It is straightforward to see that

$$k_1'(\boldsymbol{O}; \boldsymbol{\beta}, \boldsymbol{\gamma}) = \frac{\Delta}{S^c(X|A, \boldsymbol{W};\boldsymbol{\beta})}\left\{\frac{Ah(X, F_0)}{p(\boldsymbol{W};\boldsymbol{\gamma})} + \frac{(1-A)h(F_1, X)}{1 - p(\boldsymbol{W};\boldsymbol{\gamma})}\right\},$$
$$k_2'(\boldsymbol{O}; \boldsymbol{\beta}, \boldsymbol{\gamma}) = \frac{\Delta}{S^c(X|A, \boldsymbol{W};\boldsymbol{\beta})}\left\{\frac{A}{p(\boldsymbol{W};\boldsymbol{\gamma})} + \frac{1-A}{1 - p(\boldsymbol{W};\boldsymbol{\gamma})}\right\}.$$

For $r = 1, 2$, we assume that the map $(\boldsymbol{b}, \boldsymbol{g}) \mapsto u_r(\boldsymbol{b}, \boldsymbol{g})$ is differentiable at $(\boldsymbol{\beta}, \boldsymbol{\gamma})$ with derivative $\mathbf{D}_r^{\mathrm{ipw}} = (\mathbf{D}_{r,\beta}^{\mathrm{ipw}}, \mathbf{D}_{r,\gamma}^{\mathrm{ipw}})$, and use the delta method to deduce that

$$\sqrt{n}\{u_r(\widehat{\boldsymbol{\beta}}, \widehat{\boldsymbol{\gamma}}) - u_r(\boldsymbol{\beta}, \boldsymbol{\gamma})\} = \mathbf{D}_{r,\beta}^{\mathrm{ipw}} \mathbb{G}_n \boldsymbol{\phi}_\beta(\boldsymbol{O}) + \mathbf{D}_{r,\gamma}^{\mathrm{ipw}} \mathbb{G}_n \boldsymbol{\phi}_\gamma(\boldsymbol{O}) + o_p(1),$$

which implies that

$$\sqrt{n}\{U_r(\widehat{\boldsymbol{\beta}}, \widehat{\boldsymbol{\gamma}}) - u_r(\boldsymbol{\beta}, \boldsymbol{\gamma})\} = \mathbb{G}_n\{k_r'(\boldsymbol{O}; \boldsymbol{\beta}, \boldsymbol{\gamma}) + \mathbf{D}_{r,\beta}^{\mathrm{ipw}} \boldsymbol{\phi}_\beta(\boldsymbol{O}) + \mathbf{D}_{r,\gamma}^{\mathrm{ipw}} \boldsymbol{\phi}_\gamma(\boldsymbol{O})\} + o_p(1).$$

Applying the delta method once again to $\widehat{\theta}_{\mathrm{ipw}} = U_1(\widehat{\boldsymbol{\beta}}, \widehat{\boldsymbol{\gamma}})/U_2(\widehat{\boldsymbol{\beta}}, \widehat{\boldsymbol{\gamma}})$, we see that

$$\sqrt{n}(\widehat{\theta}_{\mathrm{ipw}} - \theta)$$
$$= \mathbb{G}_n\left[k_1'(\boldsymbol{O}; \boldsymbol{\beta}, \boldsymbol{\gamma}) - \theta k_2'(\boldsymbol{O}; \boldsymbol{\beta}, \boldsymbol{\gamma}) + (\mathbf{D}_{1,\beta}^{\mathrm{ipw}} - \theta \mathbf{D}_{2,\beta}^{\mathrm{ipw}}) \boldsymbol{\phi}_\beta(\boldsymbol{O}) + (\mathbf{D}_{1,\gamma}^{\mathrm{ipw}} - \theta \mathbf{D}_{2,\gamma}^{\mathrm{ipw}}) \boldsymbol{\phi}_\gamma(\boldsymbol{O})\right] + o_p(1).$$

## B.3 Asymptotics for $\widehat{\theta}_{\mathrm{dr1}}$

Here we outline a general approach to the analysis of $\widehat{\theta}_{\mathrm{dr1}} = h(\widehat{F}_1^{\mathrm{dr}}, \widehat{F}_0^{\mathrm{dr}})$ using the functional delta method. We identify $F_1, F_0$ and their estimates as elements of $BV$, the space of real-valued functions of bounded variation equipped with the total variation norm. Assuming that $h(t_1, t_0)$ is bounded, we have

$$h(F_1 + \delta F_1, F_0 + \delta F_0) - h(F_1, F_0) = h(\delta F_1, F_0) + h(F_1, \delta F_0) + h(\delta F_1, \delta F_0)$$
$$= h(\delta F_1, F_0) + h(F_1, \delta F_0) + o(\|\delta F_1\| + \|\delta F_0\|).$$

Therefore, the mapping of $(\nu_1, \nu_0) \in BV^2$ to $h(\nu_1, \nu_0) \in \mathbb{R}$ is Frechet-differentiable at $(F_1, F_0)$ with derivative

$$(\delta F_1, \delta F_0) \mapsto h(\delta F_1, F_0) + h(F_1, \delta F_0).$$

For specific working models, Zhang and Schaubel [19] show that $\widehat{F}_a^{\mathrm{dr}}(t)$ ($a = 0, 1$) converges in probability to $F_a(t)$ uniformly in $t \in (0, \tau]$ if (i) $\lambda(t|A, \boldsymbol{W}; \boldsymbol{\alpha})$ is correctly specified, (ii) $\lambda^c(t|A, \boldsymbol{W}; \boldsymbol{\beta})$ and $p(\boldsymbol{W}; \boldsymbol{\gamma})$ are both correctly specified, or (iii) both (i) and (ii) hold. For general working models that satisfy (i) or (ii), we assume that $\sqrt{n}(\widehat{F}_a^{\mathrm{dr}} - F_a)$ converges weakly to a random element in $BV$ with

$$\sqrt{n}(\widehat{F}_a^{\mathrm{dr}} - F_a)(t) = \mathbb{G}_n \phi_a(\boldsymbol{O}, t) + o_p(1)$$

for some influence function $\phi_a(\boldsymbol{O}, t)$. This is generally true for parametric working models. With $\lambda(t|A, \boldsymbol{W}; \boldsymbol{\alpha})$ and $\lambda^c(t|A, \boldsymbol{W}; \boldsymbol{\beta})$ specified as semiparametric proportional hazards models, an explicit expression for $\phi_a(\boldsymbol{O}, t)$ is given in Zhang and Schaubel [19, Web-Based Supplementary Materials]. Now it follows from Theorem 20.8 of van der Vaart [30] and simple algebra that

$$\sqrt{n}(\widehat{\theta} - \theta) = \mathbb{G}_n\{h(\phi_1(\boldsymbol{O}, \cdot), F_0) + h(F_1, \phi_0(\boldsymbol{O}, \cdot))\} + o_p(1).$$

## B.4 Asymptotics for $\widehat{\theta}_{\mathrm{dr2}}$

Here we assume that (i) $\lambda(t|A, \boldsymbol{W}; \boldsymbol{\alpha})$ is correctly specified, (ii) $\lambda^c(t|A, \boldsymbol{W}; \boldsymbol{\beta})$ and $p(\boldsymbol{W}; \boldsymbol{\gamma})$ are both correctly specified, or (iii) both (i) and (ii) hold. With possibly misspecified models, we assume that $(\widehat{\boldsymbol{\alpha}}, \widehat{\boldsymbol{\beta}}, \widehat{\boldsymbol{\gamma}})$ converges in probability to some $(\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*, \boldsymbol{\gamma}^*)$ with

$$\sqrt{n}(\widehat{\boldsymbol{\alpha}} - \boldsymbol{\alpha}^*) = \mathbb{G}_n \boldsymbol{\phi}_\alpha(\boldsymbol{O}) + o_p(1),$$
$$\sqrt{n}(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*) = \mathbb{G}_n \boldsymbol{\phi}_\beta(\boldsymbol{O}) + o_p(1),$$
$$\sqrt{n}(\widehat{\boldsymbol{\gamma}} - \boldsymbol{\gamma}^*) = \mathbb{G}_n \boldsymbol{\phi}_\gamma(\boldsymbol{O}) + o_p(1).$$

We first demonstrate the consistency of $\widehat{\theta}_{\mathrm{dr2}}$ under condition (i) or (ii). Write $\widehat{F}_a^{\mathrm{dr}} = F_a^{\mathrm{dr}}(\cdot; \widehat{\boldsymbol{\alpha}}, \widehat{\boldsymbol{\beta}}, \widehat{\boldsymbol{\gamma}})$. Under mild regularity conditions, $\widehat{\theta}_{\mathrm{dr2}}$ converges in probability to $\vartheta(\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*, \boldsymbol{\gamma}^*)$, where

$$
\begin{aligned}
\vartheta(\boldsymbol{a}, \boldsymbol{b}, \boldsymbol{g}) = \frac{1}{2}\,\mathrm{E}\Bigg[ & \frac{A\Delta h(X, F_0^{\mathrm{dr}}(\cdot; \boldsymbol{a}, \boldsymbol{b}, \boldsymbol{g}))}{p(\boldsymbol{W}; \boldsymbol{g}) S^c(X|A, \boldsymbol{W}; \boldsymbol{b})} + \frac{(1-A)\Delta h(F_1^{\mathrm{dr}}(\cdot; \boldsymbol{a}, \boldsymbol{b}, \boldsymbol{g}), X)}{\{1 - p(\boldsymbol{W}; \boldsymbol{g})\} S^c(X|A, \boldsymbol{W}; \boldsymbol{b})} \\
& + \{A - p(\boldsymbol{W}; \boldsymbol{g})\}\left\{ \frac{h(F_1^{\mathrm{dr}}(\cdot; \boldsymbol{a}, \boldsymbol{b}, \boldsymbol{g}), F_0(\cdot|\boldsymbol{W}; \boldsymbol{a}))}{1 - p(\boldsymbol{W}; \boldsymbol{g})} - \frac{h(F_1(\cdot|\boldsymbol{W}; \boldsymbol{a}), F_0)}{p(\boldsymbol{W}; \boldsymbol{g})} \right\} \\
& + \frac{A}{p(\boldsymbol{W}; \boldsymbol{g})} \int \frac{\hbar_1(t, F_0^{\mathrm{dr}}(\cdot; \boldsymbol{a}, \boldsymbol{b}, \boldsymbol{g})|\boldsymbol{W}; \boldsymbol{a})}{S^c(t|A, \boldsymbol{W}; \boldsymbol{b})} dM^c(t; \boldsymbol{b}) \\
& + \frac{1-A}{1 - p(\boldsymbol{W}; \boldsymbol{g})} \int \frac{\hbar_0(F_1^{\mathrm{dr}}(\cdot; \boldsymbol{a}, \boldsymbol{b}, \boldsymbol{g}), t|\boldsymbol{W}; \boldsymbol{a})}{S^c(t|A, \boldsymbol{W}; \boldsymbol{b})} dM^c(t; \boldsymbol{b}) \Bigg],
\end{aligned}
$$

with $M^c(t; \boldsymbol{b}) = (1-\Delta)I(X \le t) - \Lambda^c(X \wedge t|A, \boldsymbol{W}; \boldsymbol{b})$. Zhang and Schaubel [19] have shown that $F_a^{\mathrm{dr}}(\cdot; \boldsymbol{\alpha}^*, \boldsymbol{\beta}^*, \boldsymbol{\gamma}^*) = F_a$ under condition (i) or (ii). Under condition (ii), $(\boldsymbol{\beta}^*, \boldsymbol{\gamma}^*) = (\boldsymbol{\beta}, \boldsymbol{\gamma})$, $p(\boldsymbol{W}; \boldsymbol{\gamma}^*) = p(\boldsymbol{W})$, $S^c(X|A, \boldsymbol{W}; \boldsymbol{\beta}^*) = S^c(X|A, \boldsymbol{W})$, $M^c(t; \boldsymbol{\beta}^*) = M^c(t)$, and $\widehat{\theta}_{\mathrm{dr2}}$ converges to

$$
\vartheta(\boldsymbol{\alpha}^*, \boldsymbol{\beta}, \boldsymbol{\gamma}) = \frac{1}{2}\,\mathrm{E}\left[ \frac{A\Delta h(X, F_0)}{p(\boldsymbol{W}) S^c(X|A, \boldsymbol{W})} + \frac{(1-A)\Delta h(F_1, X)}{\{1 - p(\boldsymbol{W})\} S^c(X|A, \boldsymbol{W})} \right] \tag{B.5}
$$

$$
+ \frac{1}{2}\,\mathrm{E}\left[ \{A - p(\boldsymbol{W})\}\left\{ \frac{h(F_1, F_0(\cdot|\boldsymbol{W}; \boldsymbol{\alpha}^*))}{1 - p(\boldsymbol{W})} \right\} \right] \tag{B.6}
$$

$$
- \frac{1}{2}\,\mathrm{E}\left[ \{A - p(\boldsymbol{W})\}\left\{ \frac{h(F_1(\cdot|\boldsymbol{W}; \boldsymbol{\alpha}^*), F_0)}{p(\boldsymbol{W})} \right\} \right] \tag{B.7}
$$

$$
+ \frac{1}{2}\,\mathrm{E}\left[ \frac{A}{p(\boldsymbol{W})} \int \frac{\hbar_1(t, F_0|\boldsymbol{W}; \boldsymbol{\alpha}^*)}{S^c(t|A, \boldsymbol{W})} dM^c(t) \right] \tag{B.8}
$$

$$
+ \frac{1}{2}\,\mathrm{E}\left[ \frac{1-A}{1 - p(\boldsymbol{W})} \int \frac{\hbar_0(F_1, t|\boldsymbol{W}; \boldsymbol{\alpha}^*)}{S^c(t|A, \boldsymbol{W})} dM^c(t) \right], \tag{B.9}
$$

It follows from a conditioning argument that term (B.5) is equal to $\theta$. Terms (B.6) and (B.7) are 0 because $A - p(\boldsymbol{W})$ times any function of $\boldsymbol{W}$ alone has mean 0. Terms (B.8) and (B.9) are 0 because $M^c$ is a martingale. Thus, $\widehat{\theta}_{\mathrm{dr2}}$ is consistent under condition (ii). Next, suppose condition (i) holds, so that $\boldsymbol{\alpha}^* = \boldsymbol{\alpha}$, $F_a(\cdot|\boldsymbol{W}; \boldsymbol{\alpha}^*) = F_a(\cdot|\boldsymbol{W})$ ($a = 0, 1$), $\hbar_1(t, F_0|\boldsymbol{W}; \boldsymbol{\alpha}^*) = \hbar_1(t, F_0|\boldsymbol{W})$, $\hbar_0(F_1, t|\boldsymbol{W}; \boldsymbol{\alpha}^*) = \hbar_0(F_1, t|\boldsymbol{W})$, and $\widehat{\theta}_{\mathrm{dr2}}$ converges to

$$
\begin{aligned}
\vartheta(\boldsymbol{\alpha}, \boldsymbol{\beta}^*, \boldsymbol{\gamma}^*) = \frac{1}{2}\,\mathrm{E}\Bigg[ & \frac{A\Delta h(X, F_0)}{p(\boldsymbol{W}; \boldsymbol{\gamma}^*) S^c(X|A, \boldsymbol{W}; \boldsymbol{\beta}^*)} + \frac{(1-A)\Delta h(F_1, X)}{\{1 - p(\boldsymbol{W}; \boldsymbol{\gamma}^*)\} S^c(X|A, \boldsymbol{W}; \boldsymbol{\beta}^*)} \\
& - \frac{A h(F_1(\cdot|\boldsymbol{W}), F_0)}{p(\boldsymbol{W}; \boldsymbol{\gamma}^*)} - \frac{(1-A) h(F_1, F_0(\cdot|\boldsymbol{W}))}{1 - p(\boldsymbol{W}; \boldsymbol{\gamma}^*)} \\
& + h(F_1(\cdot|\boldsymbol{W}), F_0) + h(F_1, F_0(\cdot|\boldsymbol{W})) \\
& + \frac{A}{p(\boldsymbol{W}; \boldsymbol{\gamma}^*)} \int \frac{\hbar_1(t, F_0|\boldsymbol{W})}{S^c(t|A, \boldsymbol{W}; \boldsymbol{\beta}^*)} dM^c(t; \boldsymbol{\beta}^*) \\
& + \frac{1-A}{1 - p(\boldsymbol{W}; \boldsymbol{\gamma}^*)} \int \frac{\hbar_0(F_1, t|\boldsymbol{W})}{S^c(t|A, \boldsymbol{W}; \boldsymbol{\beta}^*)} dM^c(t; \boldsymbol{\beta}^*) \Bigg].
\end{aligned} \tag{B.10}
$$

Note that

$$
\mathrm{E}\{h(F_1(\cdot|\boldsymbol{W}), F_0)\} = \mathrm{E}\{h(F_1, F_0(\cdot|\boldsymbol{W}))\} = \theta. \tag{B.11}
$$

It follows from Tsiatis [21, Section 9.3 and Lemma 10.4] that

$$
\frac{\Delta}{S^c(X|A, \boldsymbol{W}; \boldsymbol{\beta}^*)} = 1 - \int \frac{dM^c(t; \boldsymbol{\beta}^*)}{S^c(t|A, \boldsymbol{W}; \boldsymbol{\beta}^*)}. \tag{B.12}
$$

Substituting (B.11) and (B.12) into (B.10) leads to

$$\vartheta(\boldsymbol{\alpha}, \boldsymbol{\beta}^*, \boldsymbol{\gamma}^*) - \theta = \frac{1}{2} \mathrm{E}\left[ \frac{A\{h(X, F_0) - h(F_1(\cdot|\boldsymbol{W}), F_0)\}}{p(\boldsymbol{W}; \boldsymbol{\gamma}^*)} \right] \tag{B.13}$$

$$+ \frac{1}{2} \mathrm{E}\left[ \frac{(1-A)\{h(F_1, X) - h(F_1, F_0(\cdot|\boldsymbol{W}))\}}{1 - p(\boldsymbol{W}; \boldsymbol{\gamma}^*)} \right] \tag{B.14}$$

$$+ \frac{1}{2} \mathrm{E}\left[ \frac{A}{p(\boldsymbol{W}; \boldsymbol{\gamma}^*)} \int \frac{\hbar_1(t, F_0|\boldsymbol{W}) - h(X, F_0)}{S^c(t|A, \boldsymbol{W}; \boldsymbol{\beta}^*)} dM^c(t; \boldsymbol{\beta}^*) \right] \tag{B.15}$$

$$+ \frac{1}{2} \mathrm{E}\left[ \frac{1-A}{1 - p(\boldsymbol{W}; \boldsymbol{\gamma}^*)} \int \frac{\hbar_0(F_1, t|\boldsymbol{W}) - h(F_1, X)}{S^c(t|A, \boldsymbol{W}; \boldsymbol{\beta}^*)} dM^c(t; \boldsymbol{\beta}^*) \right]. \tag{B.16}$$

It follows from a conditioning argument that terms (B.13) and (B.14) are 0. Terms (B.15) and (B.16) can be shown to be 0 using the arguments of Zhang and Schaubel [19, Web-Based Supplementary Materials]. Thus, $\hat{\theta}_{\mathrm{dr2}}$ is consistent under condition (i).

Next, we show that $\hat{\theta}_{\mathrm{dr2}}$ is asymptotically normal. Define

$$\varphi(\boldsymbol{O}; \boldsymbol{a}, \boldsymbol{b}, \boldsymbol{g}) = \frac{1}{2}\Bigg[ \frac{A\Delta h(X, F_0^{\mathrm{dr}}(\cdot; \boldsymbol{a}, \boldsymbol{b}, \boldsymbol{g}))}{p(\boldsymbol{W}; \boldsymbol{g})S^c(X|A, \boldsymbol{W}; \boldsymbol{b})} + \frac{(1-A)\Delta h(F_1^{\mathrm{dr}}(\cdot; \boldsymbol{a}, \boldsymbol{b}, \boldsymbol{g}), X)}{\{1 - p(\boldsymbol{W}; \boldsymbol{g})\}S^c(X|A, \boldsymbol{W}; \boldsymbol{b})}$$

$$+ \{A - p(\boldsymbol{W}; \boldsymbol{g})\}\left\{ \frac{h(F_1^{\mathrm{dr}}(\cdot; \boldsymbol{a}, \boldsymbol{b}, \boldsymbol{g}), F_0(\cdot|\boldsymbol{W}; \boldsymbol{a}))}{1 - p(\boldsymbol{W}; \boldsymbol{g})} - \frac{h(F_1(\cdot|\boldsymbol{W}; \boldsymbol{a}), F_0^{\mathrm{dr}}(\cdot; \boldsymbol{a}, \boldsymbol{b}, \boldsymbol{g}))}{p(\boldsymbol{W}; \boldsymbol{g})} \right\}$$

$$+ \frac{A}{p(\boldsymbol{W}; \boldsymbol{g})} \int \frac{\hbar_1(t, F_0|\boldsymbol{W}; \boldsymbol{a})}{S^c(t|A, \boldsymbol{W}; \boldsymbol{b})} dM^c(t; \boldsymbol{b})$$

$$+ \frac{1-A}{1 - p(\boldsymbol{W}; \boldsymbol{g})} \int \frac{\hbar_0(F_1, t|\boldsymbol{W}; \boldsymbol{a})}{S^c(t|A, \boldsymbol{W}; \boldsymbol{b})} dM^c(t; \boldsymbol{b}) \Bigg].$$

Then $\hat{\theta}_{\mathrm{dr2}} = \mathbb{P}_n \varphi(\boldsymbol{O}; \hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\gamma}})$, $\theta = \mathbb{P}_0 \varphi(\boldsymbol{O}; \boldsymbol{\alpha}^*, \boldsymbol{\beta}^*, \boldsymbol{\gamma}^*)$, and it can be argued as before that

$$\sqrt{n}(\hat{\theta}_{\mathrm{dr2}} - \theta) = \mathbb{G}_n \varphi(\boldsymbol{O}; \hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\gamma}}) + \sqrt{n}\{\mathbb{P}_0 \varphi(\boldsymbol{O}; \hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\gamma}}) - \mathbb{P}_0 \varphi(\boldsymbol{O}; \boldsymbol{\alpha}^*, \boldsymbol{\beta}^*, \boldsymbol{\gamma}^*)\}$$

$$= \mathbb{G}_n \varphi(\boldsymbol{O}; \boldsymbol{\alpha}^*, \boldsymbol{\beta}^*, \boldsymbol{\gamma}^*) + \sqrt{n}\{\vartheta(\hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\gamma}}) - \vartheta(\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*, \boldsymbol{\gamma}^*)\} + o_p(1).$$

We assume that $\vartheta$ is differentiable at $(\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*, \boldsymbol{\gamma}^*)$ with derivative $\mathbf{D}^{\mathrm{dr2}} = (\mathbf{D}_\alpha^{\mathrm{dr2}}, \mathbf{D}_\beta^{\mathrm{dr2}}, \mathbf{D}_\gamma^{\mathrm{dr2}})$, which implies that

$$\sqrt{n}\{\vartheta(\hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\gamma}}) - \vartheta(\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*, \boldsymbol{\gamma}^*)\} = \mathbf{D}^{\mathrm{dr2}}\sqrt{n}(\hat{\boldsymbol{\alpha}} - \boldsymbol{\alpha}^*, \hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*, \hat{\boldsymbol{\gamma}} - \boldsymbol{\gamma}^*) + o_p(1)$$

$$= \mathbb{G}_n\left\{ \mathbf{D}_\alpha^{\mathrm{dr2}}\boldsymbol{\phi}_\alpha(\boldsymbol{O}) + \mathbf{D}_\beta^{\mathrm{dr2}}\boldsymbol{\phi}_\beta(\boldsymbol{O}) + \mathbf{D}_\gamma^{\mathrm{dr2}}\boldsymbol{\phi}_\gamma(\boldsymbol{O}) \right\} + o_p(1).$$

It follows that

$$\sqrt{n}(\hat{\theta}_{\mathrm{dr2}} - \theta) = \mathbb{G}_n\left\{ \varphi(\boldsymbol{O}; \boldsymbol{\alpha}^*, \boldsymbol{\beta}^*, \boldsymbol{\gamma}^*) + \mathbf{D}_\alpha^{\mathrm{dr2}}\boldsymbol{\phi}_\alpha(\boldsymbol{O}) + \mathbf{D}_\beta^{\mathrm{dr2}}\boldsymbol{\phi}_\beta(\boldsymbol{O}) + \mathbf{D}_\gamma^{\mathrm{dr2}}\boldsymbol{\phi}_\gamma(\boldsymbol{O}) \right\} + o_p(1). \tag{B.17}$$

Finally, we show that $\hat{\theta}_{\mathrm{dr2}}$ is locally efficient, that is, the influence function of $\hat{\theta}_{\mathrm{dr2}}$ is equal to $\phi_{\mathrm{eff}}(\boldsymbol{O}) = \varphi(\boldsymbol{O}; \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\gamma})$ when all three models are correct. We do so by showing that $\mathbf{D}^{\mathrm{dr2}} = \mathbf{0}$ under conditions (i) and (ii). Under condition (ii), $\vartheta(\boldsymbol{a}, \boldsymbol{\beta}, \boldsymbol{\gamma}) = \theta$ for any $\boldsymbol{a}$, and its partial derivative, $\mathbf{D}_\alpha^{\mathrm{dr2}}$, must be zero. Similarly, we can show that $\mathbf{D}_\beta^{\mathrm{dr2}}$ and $\mathbf{D}_\gamma^{\mathrm{dr2}}$ are both zero under condition (i). Therefore, when all three models are correct, $\mathbf{D}^{\mathrm{dr2}} = \mathbf{0}$ and

$$\sqrt{n}(\hat{\theta}_{\mathrm{dr2}} - \theta) = \mathbb{G}_n \varphi(\boldsymbol{O}; \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\gamma}) + o_p(1) = \mathbb{G}_n \phi_{\mathrm{eff}}(\boldsymbol{O}) + o_p(1).$$

# References

1.  Rubin DB. Estimating causal effects of treatments in randomized and nonrandomized studies. J Educ Psychol. 1974;66:688–701.

2.  Agresti A. Categorical data analysis. 3rd ed. Hoboken, NJ: John Wiley and Sons; 2013.
3.  Wilcoxon F. Individual comparisons by ranking methods. Biometrics. 1945;1:80–3.
4.  Mann HB, Whitney DR. On a test of whether one of two random variables is stochastically larger than the other. Ann Math Stat. 1947;18:50–60.
5.  Acion L, Peterson JJ, Temple S, Arndt S. Probabilistic index: an intuitive non-parametric approach to measuring the size of treatment effects. Stat Med. 2006;25:591–602.
6.  Brumback LC, Pepe MS, Alonzo TA. Using the ROC curve for gauging treatment effect in clinical trials. Stat Med. 2006;25:575–90.
7.  Newcombe RG. Confidence intervals for an effect size measure based on the Mann–Whitney statistic. Part 1: General issues and tail-area-based methods. Stat Med. 2006;25:543–57.
8.  Wang C, Scharfstein DO, Colantuoni E, Girard TD, Yan Y. Inference in randomized trials with death and missingness. Biometrics. 2017;73:431–40.
9.  Liu W, Zhang Z, Nie L, Soon G. A case study in personalized medicine: rilpivirine versus efavirenz for treatment-naive HIV patients. J Am Stat Assoc. 2017;112:1381–92.
10. Lumley T. Good, better, worst: what do rank tests really test? Presented at Canterbury Statistics Day, 2012. Available online at http://www.math.canterbury.ac.nz/canterbury-tails/download/68/Lumley---Plant-and-Food.pdf. 2012.
11. Chen SX, Qin J, Tang CY. Mann-Whitney test with adjustments to pretreatment variables for missing values and observational study. J R Stat Soc, Ser B, Stat Methodol. 2013;75:81–102.
12. Vermeulen K, Thas O, Vansteelandt S. Increasing the power of the Mann–Whitney test in randomized experiments through flexible covariate adjustment. Stat Med. 2015;34:1012–30.
13. Zhang Z, Ma S, Shen C, Liu C. Estimating Mann–Whitney-type causal effects. 2018. Under review.
14. Hubbard A, van der Laan MJ, Robins JM. Nonparametric locally efficient estimation of the treatment specific survival distribution with right censored data and covariates in observational studies. In: Halloran E, Berry D, editors. Statistical models in epidemiology: the environment and clinical trials. New York: Springer; 1999. p. 135–78.
15. Zhang M, Schaubel DE. Contrasting treatment-specific survival using double-robust estimators. Stat Med. 2012;31:4255–68.
16. Zhang M. Robust methods to improve efficiency and reduce bias in estimating survival curves in randomized clinical trials. Lifetime Data Anal. 2015;21:119–37.
17. Chen P, Tsiatis AA. Causal inference on the difference of the restricted mean life between two groups. Biometrics. 2001;57:1030–8.
18. Zhang M, Schaubel DE. Estimating differences in restricted mean lifetime using observational data subject to dependent censoring. Biometrics. 2011;67:740–9.
19. Zhang M, Schaubel DE. Double-robust semiparametric estimator for differences in restricted mean lifetimes using observational data. Biometrics. 2012;68:999–1009.
20. Bickel PJ, Klaassen CAJ, Ritov Y, Wellner JA. Efficient and adaptive estimation for semiparametric models. Baltimore, MD: Johns Hopkins University Press; 1993.
21. Tsiatis AA. Semiparametric theory and missing data. New York: Springer; 2006.
22. Rosenbaum PR, Rubin DB. The central role of the propensity score in observational studies for causal effects. Biometrika. 1983;70:41–55.
23. Cox DR. Regression models and life tables (with discussion). J R Stat Soc, Ser B, Stat Methodol. 1972;34:187–200.
24. Cox DR. Partial likelihood. Biometrika. 1975;62:269–75.
25. Robins JM, Rotnitzky A, Zhao LP. Estimation of regression coefficients when some regressors are not always observed. J Am Stat Assoc. 1994;89:846–66.
26. Klein JP, Moeschberger ML. Survival analysis: techniques for censored and truncated data. 2nd ed. New York: Springer; 2003.
27. National Longitudinal Survey of Youth. NLS Handbook. Columbus, Ohio: Center for Human Resource Research, Ohio State University; 1995.
28. Rosenbaum PR, Rubin DB. Assessing sensitivity to an unobserved binary covariate in an observational study with binary outcome. J R Stat Soc, Ser B, Stat Methodol. 1983;45:212–8.
29. VanderWeele TJ, Ding P. Sensitivity analysis in observational research: introducing the E-value. Ann Intern Med. 2017;167:268–74.
30. van der Vaart AW. Asymptotic statistics. Cambridge, UK: Cambridge University Press; 1998.
31. van der Vaart AW, Wellner JA. Weak convergence and empirical processes with applications to statistics. New York: Springer; 1996.
32. Nolan D, Pollard D. Functional limit theorems for U-processes. Ann Probab. 1988;16:1291–8.