



Johns Hopkins University, Dept. of Biostatistics Working Papers

10-22-2014

ENHANCED PRECISION IN THE ANALYSIS OF RANDOMIZED TRIALS WITH ORDINAL OUTCOMES

Iván Díaz

Johns Hopkins University, Johns Hopkins Bloomberg School of Public Health, Department of Biostatistics, idiaz@jhu.edu

Elizabeth Colantuoni

Johns Hopkins University, Johns Hopkins Bloomberg School of Public Health, Department of Biostatistics

Michael Rosenblum

Johns Hopkins University, Johns Hopkins Bloomberg School of Public Health, Department of Biostatistics

Suggested Citation

Díaz, Iván; Colantuoni, Elizabeth; and Rosenblum, Michael, "ENHANCED PRECISION IN THE ANALYSIS OF RANDOMIZED TRIALS WITH ORDINAL OUTCOMES" (October 2014). *Johns Hopkins University, Dept. of Biostatistics Working Papers*. Working Paper 270.

<http://biostats.bepress.com/jhubiostat/paper270>

This working paper is hosted by The Berkeley Electronic Press (bepress) and may not be commercially reproduced without the permission of the copyright holder.

Copyright © 2011 by the authors

Enhanced Precision in the Analysis of Randomized Trials with Ordinal Outcomes

Iván Díaz^{*1}, Elizabeth Colantuoni¹, and Michael Rosenblum¹

¹Department of Biostatistics, Johns Hopkins Bloomberg School of Public Health.

October 20, 2014

Abstract

We present a general method for estimating the effect of a treatment on an ordinal outcome in randomized trials. The method is robust in that it does not rely on the proportional odds assumption. Our estimator leverages information in prognostic baseline variables, and has all of the following properties: (i) it is consistent; (ii) it is locally efficient; (iii) it is guaranteed to match or improve the precision of the standard, unadjusted estimator. To the best of our knowledge, this is the first estimator of the causal relation between a treatment and an ordinal outcome to satisfy these properties. We demonstrate the estimator in simulations based on resampling from a completed randomized clinical trial of a new treatment for stroke; we show potential gains of up to 39% in relative efficiency compared to the unadjusted estimator. The proposed estimator could be a useful tool for analyzing randomized trials with ordinal outcomes, since existing methods either rely on model assumptions that are untenable in many practical applications, or lack the efficiency properties of the proposed estimator. We provide R code implementing the estimator.

1 Introduction

Our methods are motivated by questions that arose in a randomized trial of a new treatment for stroke. This trial is fully described in Section 2, and is summarized here. The outcome is an individual's score on the modified Rankin Scale (mRS), which takes integer values from 0 to 6 representing his/her degree of functional disability. The goal is to estimate the marginal effect of the treatment versus control on this ordinal outcome. We leverage information in prognostic baseline variables to improve precision in estimating the unconditional treatment effect. We propose a new class of estimators for this problem, which can be applied in randomized trials or more generally in observational studies.

In randomized trials with ordinal outcomes, a standard approach is to estimate the marginal treatment effect by assuming the proportional odds model. A downside is that this approach relies on the proportional odds assumption, which may be untenable in many practical applications. To address this, we introduce a nonparametric extension of the parameter in the proportional odds model. This parameter

^{*}idiaz@jhu.edu

is a summary measure for comparing any pair of ordered, multinomial distributions. It equals the average of the cumulative log odds ratios across categories. We propose an estimator of this parameter that incorporates information in prognostic baseline variables and is guaranteed to be consistent, without requiring parametric model assumptions. Furthermore, this estimator is guaranteed to have greater or equal asymptotic precision than the standard unadjusted estimator; our estimator is also locally semiparametric efficient. The above claims hold under the assumption that outcomes are either all observed, or are missing completely at random as defined in Section 3.2. We also consider other ways that outcomes could be missing, and describe how the estimator adjusts for informative loss to follow-up.

[13, 14, 17] laid the foundation for locally efficient estimation of causal effects, which has been extended to incorporate enhanced efficiency properties, e.g., by [19, 24, 28, 21, 5, 20, 15, 8]. Targeted minimum loss based estimation of the effect of treatment on binary, continuous, and time to event outcomes in randomized trials was first discussed in [10] and [18]. We build on the general ideas in [24] and [8] to construct a targeted minimum loss based estimator of the new parameter defined in Section 3.1, which represents the marginal effect of a treatment on an ordinal outcome, without making the proportional odds assumption. Our main contribution is to develop an estimator with variance guaranteed to be smaller or equal than the variance of the unadjusted estimator, asymptotically, for this parameter. Focusing on this goal allows us to develop estimators with computational advantages not available for the estimators of [8], which were constructed to achieve additional efficiency properties. Specifically, our estimators do not need to solve a non-convex optimization problem, which can be computationally challenging. We demonstrate our estimator in simulations based on resampling from a completed randomized trial, where we show potential efficiency gains of up to 39% compared to the unadjusted estimator.

In Section 2, we describe the MISTIE II randomized trial, which motivated our work. In Section 3, we define our estimation problem and present the nonparametric extension of the proportional odds model along with a corresponding unadjusted estimator. In Sections 4 and 5, we develop a new estimator of this nonparametric extension that is guaranteed to have asymptotic variance smaller or equal to the variance of the unadjusted estimator. We present a simulation study based on data from the MISTIE II trial in Section 6. We discuss directions for future research in Section 7.

2 Motivating Application: MISTIE II Trial

MISTIE II is a randomized, prospective Phase II trial comparing a novel, minimally invasive surgical procedure to standard medical care among patients with intracerebral hemorrhage (ICH) [1]. The primary outcome is the modified Rankin Scale (mRS) measured 180 days after randomization. The mRS measures an individual's degree of disability, and has seven ordinal levels. The following variables were collected at baseline and are strongly correlated with the outcome: age, ICH volume, and National Institutes of Health Stroke Scale (NIHSS) [4]. These variables are used in our adjusted estimator.

The treatment effect of interest is the average over levels of mRS of the cumulative log odds ratios comparing treatment versus control, as defined in (24) in Section 3.1. For both the unadjusted and adjusted estimators, we computed standard errors and confidence intervals based on the bias-corrected and accelerated (BCa) bootstrap [6]. In the MISTIE II trial, the unadjusted estimate of the treatment effect is 0.252 (standard error, SE = 0.401) with a 95% confidence interval of -0.558 to 1.047. The adjusted estimate of the treatment effect is 0.356 (SE = 0.346) with 95% confidence interval -0.298 to 1.02. The estimated relative efficiency, i.e., the ratio of estimated variances comparing the unadjusted estimator to the adjusted estimator, is 1.34. A relative efficiency of 1.34 means that the required sample

size is reduced by $1 - (1/1.34) \approx 25\%$ when using the adjusted estimator compared to the unadjusted estimator, in order to achieve a given power (e.g. 80% power) at a given alternative, asymptotically. This gives an initial indication that adjusting for baseline variables has potential to improve precision in the context of ordinal outcomes. Below, we present theoretical results and a simulation study comparing our adjusted estimator versus the unadjusted estimator.

3 Problem Definition

3.1 Notation

Let $O = (W, A, \Delta, \Delta Y) \sim P_0$ denote a random vector with unknown distribution P_0 , where Y denotes an ordinal outcome taking values in $\{1, \dots, K\}$ observed only when the censoring variable $\Delta = 1$; A denotes a binary treatment; W denotes a vector of pre-treatment variables. We refer to the distribution P_0 as the observed data distribution, in contrast to the distribution of potential outcomes defined next. Define the potential outcomes $Y_a : a \in \{0, 1\}$ to be the outcomes that would have been observed had treatment $A = a$ been assigned with probability one.

We consider randomized trials where all participants are randomly assigned to treatment or control independent of baseline variables W . Therefore, the distribution of A given W is set by design, and so is known. We assume that $P_0 \in \mathcal{M}$, where \mathcal{M} is the model defined as all continuous densities on $(W, A, \Delta, \Delta Y)$ with respect to a common measure ν , such that A is independent of W . The index naught in P_0 denotes the true distribution, and is added to the notation of any quantity whenever it is necessary to avoid confusion, but is omitted otherwise. Denote $\theta_a(k) = P(Y_a \leq k)$, and $\theta = (\theta_a(k) : k = 1, \dots, K - 1; a = 0, 1)$. Consider the following proportional odds model for ordinal outcomes:

$$\text{logit } \theta_a(k) = \alpha_k + \beta_{\text{par}} a, \text{ for } a = 0, 1; k = 1, \dots, K - 1, \quad (1)$$

where $\text{logit}(x) = \log\{x/(1 - x)\}$. Setting $a = 0$ in the above display, it follows that $\alpha_k = \text{logit } \theta_0(k)$. Therefore, the main assumption of the proportional odds model is that $\text{logit } \theta_1(k) - \text{logit } \theta_0(k)$ equals a common value β_{par} for all $k \leq K - 1$. The proportional odds assumption may be untenable in many practical applications, and it is desirable to use a method that is robust to this assumption being false. To address this, we define a univariate contrast between two ordered, multinomial distributions, using the general approach for constructing nonparametric extensions outlined by [11]. Define the nonparametric extension $\beta(\theta)$ of β_{par} as follows:

$$\beta(\theta) = \arg \min_{\beta' \in \mathbb{R}} \sum_{k=1}^{K-1} (\text{logit } \theta_1(k) - \text{logit } \theta_0(k) - \beta')^2.$$

The solution to the above minimization problem is

$$\beta(\theta) = \frac{1}{K-1} \sum_{k=1}^{K-1} \{\text{logit } \theta_1(k) - \text{logit } \theta_0(k)\} = \frac{1}{K-1} \sum_{k=1}^{K-1} \log \left\{ \frac{\theta_1(k)}{1 - \theta_1(k)} \bigg/ \frac{\theta_0(k)}{1 - \theta_0(k)} \right\}. \quad (2)$$

This implies that $\beta(\theta)$ is the average of the cumulative log odds ratios across $K - 1$ categories. We omit the dependence of β on θ whenever there is no confusion. The parameter β is well-defined without

having to assume the proportional odds model; in cases where the proportional odds assumption does hold, the above parameter coincides with the parameter β_{par} of the proportional odds model.

An estimator of β is called consistent if it converges to β in probability. We say an estimator of β is locally efficient at P_0 if it achieves the semiparametric efficiency bound in the model \mathcal{M} .

The above nonparametric extension of the proportional odds model has two main advantages compared to assuming model (1). First, the interpretation of β as a contrast between marginal distributions under treatment versus control averaged across categories does not depend on the correct specification of (1). Second, this extension allows construction of estimators of the causal contrast β (and corresponding standard errors) that are consistent regardless of whether assumption (1) holds.

Our goal is to estimate the parameter (24), which represents an unconditional treatment effect, i.e., a contrast between marginal distributions of the outcome under assignment to treatment versus control. An alternative goal, not considered here, is to estimate a conditional effect of treatment, i.e., a contrast between distributions conditioned on the values of certain baseline variables. Though we do not estimate conditional effects, we do harness information in baseline variables in our estimators of unconditional treatment effects.

3.2 Identification of β as a Function of P_0

Under the ignorability assumption $Y_a \perp\!\!\!\perp (A, \Delta) | W$ for each $a \in \{0, 1\}$, the positivity assumption $P(A = a, \Delta = 1 | W) > 0$ with probability 1 for each $a \in \{0, 1\}$, and the consistency assumption that $Y_a = Y$ on the event $A = a$, the treatment-specific probabilities $\theta_a(k)$ are identified as a function of P_0 as [see, e.g., 12]

$$\theta_a(k) = E\{p(k, a, W)\}, \quad (3)$$

for $k \in \{1, \dots, K - 1\}$, $a \in \{0, 1\}$, where $p(k, a, w) = P_0(Y \leq k | \Delta = 1, A = a, W = w)$ and the expectation in (3) is with respect to the marginal distribution of W under P_0 . This implies that under the above assumptions, the parameter β is identified as a function of observed data distribution P_0 of $O = (W, A, \Delta, \Delta Y)$, despite the definition of β being in terms of potential outcomes Y_a .

In a randomized trial, $(Y_a, W) \perp\!\!\!\perp A$ by design, which we call the randomization assumption. This assumption, combined with the assumption that outcomes are missing at random (i.e., $Y_a \perp\!\!\!\perp \Delta | A, W$, denoted MAR [16]), implies the ignorability assumption above. If, in addition, the positivity and consistency assumptions hold, then (3) follows. Throughout, we make the randomization and MAR assumptions.

We define the hazard representation of $p(k, a, w)$ as follows:

$$p(k, a, w) = \sum_{m=1}^k h(m, a, w) \prod_{j=0}^{m-1} (1 - h(j, a, w)), \quad (4)$$

where we define

$$h(m, a, w) = P(Y = m | Y \geq m, \Delta = 1, A = a, W = w),$$

whenever $P(Y \geq m | \Delta = 1, A = a, W = w) > 0$, and $h(m, a, w) = 0$ otherwise. It follows that $h(0, a, w) = 0$ and $h(K, a, w) = 1$ with probability 1. Under the assumptions in the first sentence of this subsection, we can express $\theta_a(k)$ in terms of the function h as follows:

$$\theta_a(k) = \sum_{m=1}^k \int h(m, a, w) \prod_{j=0}^{m-1} (1 - h(j, a, w)) dP_W(w), \quad (5)$$

where P_W is the marginal distribution of W . Substituting (5) into (24) gives an expression for β in terms of the function h , which is identified from the observed data distribution P_0 .

We use the notation $g_{A,0}(a|w) = P_0(A = a|W = w)$, $g_{\Delta,0}(a, w) = P_0(\Delta = 1|A = a, W = w)$, and $g_0 = (g_{A,0}, g_{\Delta,0})$. The randomization assumption implies $g_{A,0}(a|w)$ does not depend on w . Since the randomization probability is set by design in our context of a randomized trial, $g_{A,0}$ is known. However, our adjusted estimators involve fitting a parametric model for $g_{A,0}$; intuitively, this model fit captures chance imbalances of the baseline variables W between arms, for a given data set. The general approach of improving efficiency by estimating known nuisance parameters such as $g_{A,0}$ has been shown, e.g., by [14] and van der Laan and Robins [22, Theorem 2.3].

In the remainder of the paper our statistical parameter of interest is $\beta(\theta)$, where θ is defined in terms of h and the marginal distribution of W through (5). Then $\beta(\theta)$ is a function of P_0 , with the caveat that it only has a causal interpretation under the assumptions in the first sentence of this subsection. The estimators of $\theta_a(k)$ and β developed in Sections 4 and 5 are substitution estimators, that is, estimators constructed from (5) and (24) by replacing h and P_W by estimates of these quantities. Substitution estimators have the advantage that they remain within bounds of the parameter space, a characteristic that is particularly desirable in estimation of the multinomial probabilities $\theta_a(k)$. The goal of this paper is to construct a substitution estimator of β with the properties outlined in the abstract. We start by discussing simpler estimators: the unadjusted estimator and a regression-based estimator.

3.3 Unadjusted Estimator

Under the assumption that outcomes are missing completely at random (i.e., $(A, \Delta) \perp\!\!\!\perp (W, Y_a)$, denoted MCAR), it can be shown that $\theta_a(k) = P(Y \leq k|\Delta = 1, A = a)$. Then each $\theta_a(k)$ is consistently estimated by the following unadjusted estimator that ignores baseline variables:

$$\hat{\theta}_{\text{unadj},a}(k) = \frac{\sum_{i=1}^n \Delta_i \mathbf{1}\{A_i = a, Y_i \leq k\}}{\sum_{i=1}^n \Delta_i \mathbf{1}\{A_i = a\}}, \quad (6)$$

where $\mathbf{1}(X)$ is the indicator variable taking value 1 if X is true and 0 otherwise. Define the vector $\hat{\theta}_{\text{unadj}} = (\hat{\theta}_{\text{unadj},a}(k) : k = 1, \dots, K-1; a = 0, 1)$. The corresponding substitution estimator of β is $\hat{\beta}_{\text{unadj}} = \beta(\hat{\theta}_{\text{unadj}})$. The asymptotic behavior of $\hat{\beta}_{\text{unadj}}$ is obtained through the delta method as

$$\sqrt{n}(\hat{\beta}_{\text{unadj}} - \beta) = \frac{1}{\sqrt{n}} \sum_{i=1}^n S_{\hat{\beta}_{\text{unadj}}}(O_i) + o_P(1),$$

where

$$S_{\hat{\beta}_{\text{unadj}}}(O) = \frac{1}{K-1} \sum_{k=1}^{K-1} \left(\frac{\Delta \mathbf{1}\{A = 1\}}{\theta_1(k)(1 - \theta_1(k))g_{0,A}(1|W)g_{0,\Delta}(1, W)} [\mathbf{1}\{Y \leq k\} - \theta_1(k)] - \frac{\Delta \mathbf{1}\{A = 0\}}{\theta_0(k)(1 - \theta_0(k))g_{0,A}(0|W)g_{0,\Delta}(1, W)} [\mathbf{1}\{Y \leq k\} - \theta_0(k)] \right), \quad (7)$$

is the influence function of $\hat{\beta}_{\text{unadj}}$. As a consequence, $\hat{\beta}_{\text{unadj}}$ is asymptotically Gaussian with variance $V(S_{\hat{\beta}_{\text{unadj}}}(O))/n$. Under the MCAR assumption, this estimator is consistent and efficient for the case where only $(A, \Delta, \Delta Y)$ are observed; however, the unadjusted estimator is generally not efficient for

the case where $(W, A, \Delta, \Delta Y)$ are observed. Intuitively, this is because the unadjusted estimator fails to leverage prognostic information in baseline variables W . Furthermore, under the less restrictive missing at random (MAR) assumption, the unadjusted estimator will generally not even be consistent, while our proposed estimators are consistent under conditions given in Section 4.

3.4 Covariate Adjustment: Regression and Inverse Probability Weighted Estimators

Consider an estimator $\hat{p}(k, a, W)$ for the probability $p(k, a, W)$, for example, based on fitting an ordered logistic regression model. Under correct specification of that model, the estimator $\frac{1}{n} \sum_{i=1}^n \hat{p}(k, a, W_i)$ can be shown to be consistent and efficient for $\theta_a(k)$. However, if that model is misspecified, this regression estimator will generally be inconsistent, even in randomized trials in which an unadjusted, consistent estimator exists [7]. This undesirable property discourages us from using the aforementioned estimator.

As an alternative, an inverse probability weighted (IPW) estimator of $\theta_a(k)$ is given by

$$\hat{\theta}_{\text{ipw},a}(k) = \frac{\sum_{i=1}^n \Delta_i \mathbb{1}\{A_i = a, Y_i \leq k\} \{\hat{g}_A(a|W_i) \hat{g}_\Delta(a, W_i)\}^{-1}}{\sum_{i=1}^n \Delta_i \mathbb{1}\{A_i = a\} \{\hat{g}_A(a|W_i) \hat{g}_\Delta(a, W_i)\}^{-1}},$$

where $\hat{g}_A(a|w)$ and $\hat{g}_\Delta(a, w)$ are estimators of g_A and g_Δ , respectively. (An alternative definition of the IPW estimator replaces the denominator in the above display by n). The IPW estimator $\hat{\beta}_{\text{ipw}} = \beta(\hat{\theta}_{\text{ipw}})$ is consistent for β if the estimators $\hat{g}_A, \hat{g}_\Delta$ are consistent and converge to g_A, g_Δ at rate $n^{-1/2}$. However, the IPW estimator is generally not locally efficient.

4 Doubly Robust, Locally Efficient Estimation

We present a locally efficient, doubly robust estimator of θ . The corresponding substitution estimator of β can be shown to satisfy similar properties using the delta method. The efficient influence function for estimating $\theta_a(k)$ in the model \mathcal{M} is given by [see e.g., 2]

$$D_{ak}(O) = \frac{\Delta \mathbb{1}\{A = a\}}{g_A(a|W)g_\Delta(a, W)} \{\mathbb{1}(Y \leq k) - p(k, a, W)\} + p(k, a, W) - \theta_a(k), \quad (8)$$

where, for conciseness, we suppress the dependence of D_{ak} on g and p in our notation. The function D_{ak} has two important properties for estimation of $\theta_a(k)$. First, it is a doubly robust estimating function, i.e., for given estimators \hat{p}, \hat{g} of p and $g = (g_A, g_\Delta)$, respectively, the estimator of $\theta_a(k)$ formed by solving for $\theta_a(k)$ in the following estimating equation:

$$0 = \sum_{i=1}^n D_{ak}(O_i) = \sum_{i=1}^n \left[\frac{\Delta_i \mathbb{1}\{A_i = a\}}{\hat{g}_A(a|W_i) \hat{g}_\Delta(a, W_i)} \{\mathbb{1}(Y_i \leq k) - \hat{p}(k, a, W_i)\} + \hat{p}(k, a, W_i) - \theta_a(k) \right], \quad (9)$$

is consistent if at least one of p or g is estimated consistently at sufficient rate, as described below. This double robustness property is desirable since it guarantees that improper adjustment for covariates through a misspecified working model for p does not cause asymptotic bias in the estimator in randomized trials if outcomes satisfy MCAR, or more generally if outcomes satisfy MAR and g_Δ is consistently estimated at a sufficient rate. Second, the efficient influence function characterizes the efficiency bound

for estimation of $\theta_a(k)$ in the model \mathcal{M} [3]. Specifically, under consistent estimation of p and g at sufficient rate, the estimator $\tilde{\theta}_a(k)$ that solves the estimating equation (9) has variance smaller or equal to that of any regular, asymptotically linear estimator of $\theta_a(k)$ in \mathcal{M} . If we also assume MCAR, then under consistent estimation of p and g at sufficient rate, the corresponding substitution estimator $\beta(\tilde{\theta})$ is equal or superior to $\hat{\beta}_{\text{unadj}}$ in terms of asymptotic precision.

Various estimators using D_{ak} have been proposed in the context of a binary outcome, including the augmented inverse probability weighted estimator, the doubly robust weighted least squares estimator, and the targeted maximum likelihood estimator [see 2, 9, for a review]. In principle, those methods may be used to estimate the multinomial probabilities $\theta_a(k)$ by considering K binary outcomes $\mathbb{1}\{Y \leq k\}$: $k = 1, \dots, K$ separately. This approach, however, does not guarantee that the resulting estimates yield a well defined multinomial distribution, i.e., there is no guarantee that the estimate of $\theta_a(k)$ will be less than or equal to the estimate of $\theta_a(k+1)$ for all k .

We next express the efficient influence function (8) in terms of the hazard functions h . Define $I(k) = \mathbb{1}\{Y = k\}$ and $\bar{I}(k) = \mathbb{1}\{Y \geq k\}$. Intuitively, $\bar{I}(k)$ is an indicator of being at risk for having $Y = k$, given only the information $I(1), \dots, I(k-1)$. The following is proved in Appendix 1:

Theorem 1 (Hazard Representation of Efficient Influence Function). *The efficient influence function D_{ak} for the vector $(\theta_a(k) : k = 1, \dots, K-1; a = 0, 1)$ in the model \mathcal{M} can be expressed as*

$$D_{ak}(O) = \left[\sum_{j=1}^{K-1} \bar{I}(j) Z_{ak}(j, A, W) \{I(j) - h(j, A, W)\} \right] + p(k, a, W) - \theta_a(k), \quad (10)$$

where

$$Z_{ak}(j, A, W) = \frac{\Delta \mathbb{1}\{A = a\}}{g_A(a|W) g_\Delta(a, W)} \frac{1 - p(k, a, W)}{1 - p(j, a, W)} \mathbb{1}(j \leq k), \quad (11)$$

and p is defined as a function of h in (4).

We now describe a doubly robust, substitution estimator of θ , using targeted minimum loss based estimation (TMLE).

Targeted Minimum Loss Based Estimator It will be convenient to sometimes use a modified data set in which only observations O_i with $\Delta_i = 1$ are included, and each such observation O_i is repeated for each $j = 1, \dots, K-1$ with the indicators $I_i(j) = \mathbb{1}\{Y_i = j\}$ and $\bar{I}_i(j) = \mathbb{1}\{Y_i \geq j\}$ replacing Y_i as the outcome, i.e., the following data set:

$$\{(j, W_i, A_i, \bar{I}_i(j), I_i(j)) : \Delta_i = 1, j = 1, \dots, K-1; i = 1, \dots, n\}. \quad (12)$$

This data set is referred to as the long form, and the original data set is referred to as the short form. An observation in the long form data set is a vector of the form $(j, W_i, A_i, \bar{I}_i(j), I_i(j))$ where $\Delta_i = 1$.

A TMLE for θ may be computed through the following iterative procedure:

Step 1. *Initial estimators.* Obtain initial estimators \hat{g}_A , \hat{g}_Δ , and \hat{h} of g_A , g_Δ , and h , respectively.

Initial estimator of h . To estimate h we use an approach borrowed from the survival analysis literature: we pool the data and smooth across categories. An estimator \hat{h} may be obtained by running a prediction algorithm of $I(j)$ as a function of A , W , and j among observations with $\bar{I}(j) = 1$ in the long form data set (12).

Initial estimators of g_A and g_Δ . We estimate g_A by fitting a parametric model for the probability of $\{A = 1\}$ as a function of W in the short form data set. We estimate g_Δ analogously. In a randomized trial, g_A is set by design; however, it can still improve efficiency of the TMLE to estimate g_A using, e.g., the proportion of individuals in the treatment group, or a parametric model that contains baseline variables and an intercept. In general, the functional relation between Δ and (A, W) is unknown to the researcher; one may apply flexible prediction techniques such as model stacking [27] or super learning [25] to estimate g_Δ .

Step 2. *Iteratively update estimate of h by maximum likelihood in least favorable parametric model.* Initialize $l = 0$, and let $\hat{h}^l = \hat{h}$. In the long form data set, for each observation, augment it with $2(K - 1)$ covariates $Z_{ak}^l(j, A, W) : k = 1, \dots, K - 1; a = 0, 1$ by substituting the estimates \hat{h}^l and \hat{g} in (11). This augmented long form data set is defined as

$$\{(j, Z_{ak}^l(j, A_i, W_i), W_i, A_i, \bar{I}_i(j), I_i(j)) : \Delta_i = 1, j \leq K - 1; i \leq n; a \in \{0, 1\}; k \leq K - 1\}. \quad (13)$$

Estimate the parameter vector $\epsilon = \{\epsilon_{ak} : k = 1, \dots, K - 1; a = 0, 1\}$ in the logistic hazard submodel for $h(j, a, w)$:

$$\text{logit } h_\epsilon^l(j, a, w) = \text{logit } \hat{h}^l(j, a, w) + \sum_{k=1}^{K-1} \epsilon_{1k} Z_{1k}^l(j, a, w) + \sum_{k=1}^{K-1} \epsilon_{0k} Z_{0k}^l(j, a, w), \quad (14)$$

by computing the following maximum likelihood estimator:

$$\hat{\epsilon} = \arg \max_{\epsilon} \sum_{i=1}^n \sum_{j=1}^{K-1} \bar{I}_i(j) \log \{h_\epsilon^l(j, A_i, W_i)^{I_i(j)} (1 - h_\epsilon^l(j, A_i, W_i))^{1-I_i(j)}\}. \quad (15)$$

The maximizer $\hat{\epsilon}$ can be computed using standard statistical software by a logistic regression of $I(j)$ on the $2(K - 1)$ variables $\{Z_{ak}^l(j, A, W) : k = 1, \dots, K - 1; a = 0, 1\}$ among observations with $\bar{I}(j) = 1$ in the augmented long form data set (13), and using $\text{logit } \hat{h}^l(j, a, w)$ as an offset. Define $\hat{h}^{l+1} = h_{\hat{\epsilon}}^l$. Since the gradient of the expression on right side of (15) is 0 at the maximizer $\hat{\epsilon}$, it follows that \hat{h}^{l+1} solves the estimating equations

$$\sum_{i=1}^n \sum_{j=1}^{K-1} \bar{I}_i(j) Z_{ak}^l(j, A_i, W_i) \{I_i(j) - \hat{h}^{l+1}(j, A_i, W_i)\} = 0, \quad (16)$$

for each $k = 1, \dots, K - 1, a = 0, 1$.

Step 3. *Iterate.* Update $l = l + 1$ and iterate Step 2 until convergence. We use

$$\frac{1}{n(K-1)} \sum_i \sum_j \{\hat{h}^{l+1}(j, A_i, W_i) - \hat{h}^l(j, A_i, W_i)\}^2 \leq \frac{10^{-4}}{n}, \quad (17)$$

as a stopping criterion.

Let \hat{h}^* denote \hat{h}^l corresponding to the final iteration l of the above procedure. Define the TMLE of $\theta_a(k)$ as

$$\hat{\theta}_{\text{adj},a}(k) = \frac{1}{n} \sum_{i=1}^n \hat{h}^*(k, a, W_i) \prod_{j=0}^{k-1} (1 - \hat{h}^*(j, a, W_i)).$$

Analogously, define the TMLE of β as $\hat{\beta}_{\text{adj}} = \beta(\hat{\theta}_{\text{adj}})$. This is a substitution estimator, i.e., it is the result of substituting estimates of h and P_W into the parameter definitions (5) and (24). It follows from (16), (17), and the representation of D_{ak} in (10), that the efficient influence function estimating equation is approximately solved at the last iteration, i.e.,

$$\frac{1}{n} \sum_{i=1}^n \hat{D}_{ak}^*(O_i) \approx 0,$$

where \hat{D}_{ak}^* is obtained by substituting \hat{g}_A , \hat{g}_Δ , and \hat{h}^* into (10) and (11).

First, consider the case where \hat{h} and \hat{g} converge to h and g , respectively, each at rate faster than $n^{-1/4}$. It follows from the general arguments in van der Laan and Rose [23, Appendix 18] that

$$\sqrt{n}(\hat{\theta}_{\text{adj},a}(k) - \theta_a(k)) = \frac{1}{\sqrt{n}} \sum_{i=1}^n D_{ak}(O_i) + o_P(1),$$

which implies that $\hat{\theta}_{\text{adj},a}(k)$ is efficient in the semiparametric model \mathcal{M} . The delta method implies

$$\sqrt{n}(\hat{\beta}_{\text{adj}} - \beta) = \frac{1}{\sqrt{n}} \sum_{i=1}^n S(O_i) + o_P(1),$$

where

$$S(O) = \frac{1}{K-1} \sum_{k=1}^{K-1} \left\{ \frac{D_{1k}(O)}{[1 - \theta_1(k)]\theta_1(k)} - \frac{D_{0k}(O)}{[1 - \theta_0(k)]\theta_0(k)} \right\}, \quad (18)$$

is the efficient influence function for estimation of β in \mathcal{M} . In this case, if we additionally assume MCAR, local efficiency of $\hat{\beta}_{\text{adj}}$ implies that $\hat{\beta}_{\text{adj}}$ has asymptotic variance smaller or equal to the asymptotic variance of $\hat{\beta}_{\text{unadj}}$.

Second, consider the case where \hat{h} does not converge to h (which could occur if the model for h is misspecified) but does converge to some limit, and g is estimated consistently at rate $n^{-1/2}$ (which holds in our context of a randomized trial if g_A, g_Δ are estimated using parametric models that include intercepts, and if either (i) outcomes satisfy MCAR or (ii) outcomes satisfy MAR and the parametric model for g_Δ is correct). Then $\hat{\beta}_{\text{adj}}$ is consistent. However, in this case and under MCAR, there is no guarantee that the asymptotic variance of $\hat{\beta}_{\text{adj}}$ is at most that of $\hat{\beta}_{\text{unadj}}$. In the next section we present a substitution estimator of β that has this guarantee.

5 Guaranteed Equal or Better Asymptotic Efficiency Compared to the Unadjusted Estimator

We develop an adjusted estimator, denoted $\hat{\beta}_{\text{adj,eff}}$, that is guaranteed to be consistent and to have asymptotic variance smaller or equal to the asymptotic variance of $\hat{\beta}_{\text{unadj}}$, under MCAR and if the initial estimator for g is based on parametric models that include intercept terms. The corresponding estimator

$\hat{\beta}_{\text{adj,eff}}$ provides protection against biased estimates of h , allowing researchers to safely adjust for covariates in the analysis of a randomized trial, with the goal of improving precision for estimating β , under MCAR. Additionally, this estimator has the same double robustness and local efficiency properties as the estimator in the previous section. In particular, if the outcome is MAR and the parametric model used to estimate g_Δ is correctly specified, then the convergence of $\hat{\beta}_{\text{adj,eff}}$ to β still holds whereas $\hat{\beta}_{\text{unadj}}$ will generally be inconsistent.

The only difference between the TMLE algorithm for $\hat{\beta}_{\text{adj}}$ in the previous section and the TMLE algorithm for $\hat{\beta}_{\text{adj,eff}}$ below is that the latter updates the estimators \hat{g}_A and \hat{g}_Δ at each iteration. This is done by fitting logistic regression models for the expectation of A and Δ conditional on auxiliary variables that are functions of the following:

$$M_{ak}(W) = \frac{p(k, a, W) - \theta_a(k)}{g_A(a|W)}, \quad H_{ak}(A, W) = \frac{\mathbb{1}\{A = a\}}{g_\Delta(a, W)} M_{ak}(W). \quad (19)$$

The TMLE algorithm for $\hat{\beta}_{\text{adj,eff}}$ is given below:

Step 1. *Initial estimators.* Obtain initial estimators \hat{g}_A , \hat{g}_Δ , and \hat{h} of g_A , g_Δ , and h , respectively. This may be done as described in the previous section.

Step 2. *Iteratively update estimates of h and g .* Initialize $l = 0$, and let $\hat{h}^l = \hat{h}$, $\hat{g}_A^l = \hat{g}_A$, $\hat{g}_\Delta^l = \hat{g}_\Delta$. Compute $\hat{\theta}_a^l(k)$ by substituting the estimate \hat{h}^l in the definition of $\theta_a(k)$ in (5), for each $a \in \{0, 1\}$, $k \leq K - 1$.

- (a) *Update \hat{h}^l .* This is the same as step 2 in the TMLE algorithm from Section 4, except using the current updates of the estimators of g at each iteration, rather than always using the initial estimators of g . Augment each observation in the long form data set (12) by $2(K - 1)$ covariates $Z_{ak}^l(j, A, W) : k = 1, \dots, K - 1; a = 0, 1$, where each $Z_{ak}^l(j, A, W)$ is constructed by substituting the estimates \hat{h}^l , \hat{g}_A^l , and \hat{g}_Δ^l in (11). Estimate ϵ by $\hat{\epsilon}$ in the logistic hazard submodel (14) by maximum likelihood estimation as in (15).
- (b) *Update \hat{g}_A^l .* Let $\hat{\theta}_a^l(k)$ denote $\theta_a(k)$ with \hat{h}^l and the empirical distribution of W substituted for h and P_W , respectively, in (5). Let $\hat{p}(k, a, w)$ denote $p(k, a, w)$ with \hat{h}^l substituted for h in (4). Let $M_{ak}^l(W)$ denote $M_{ak}(W)$ with \hat{p} , \hat{g}_A^l , $\hat{\theta}_a^l(k)$ substituted for the corresponding components in (19). In the short form data set, compute the auxiliary covariate

$$M^l(W) = \sum_{k=1}^{K-1} \left\{ \frac{M_{1k}^l(W)}{[1 - \hat{\theta}_1^l(k)]\hat{\theta}_1^l(k)} - \frac{M_{0k}^l(W)}{[1 - \hat{\theta}_0^l(k)]\hat{\theta}_0^l(k)} \right\}.$$

Estimate the parameter ν in the following logistic regression submodel for $g_A(1|w)$:

$$\text{logit } g_{A,\nu}^l(1|w) = \text{logit } \hat{g}_A^l(1|w) + \nu M^l(w),$$

by regressing A on W among all participants $i = 1, \dots, n$ using a logistic regression model with single covariate $M^l(W)$ and offset $\text{logit } \hat{g}_A^l(1|W)$; denote the corresponding maximum likelihood estimate of ν by $\hat{\nu}$.

(c) Update \hat{g}_Δ^l . Let $H_{ak}^l(A, W)$ denote $H_{ak}(A, W)$ with \hat{g}_Δ^l and $M_{ak}^l(W)$ substituted for g_Δ and $M_{ak}(W)$, respectively, in (19). In the short form data set, compute the auxiliary covariate

$$H^l(A, W) = \sum_{k=1}^{K-1} \left\{ \frac{H_{1k}^l(A, W)}{[1 - \hat{\theta}_1^l(k)]\hat{\theta}_1^l(k)} - \frac{H_{0k}^l(A, W)}{[1 - \hat{\theta}_0^l(k)]\hat{\theta}_0^l(k)} \right\}.$$

Estimate the parameter γ in the following logistic regression submodel for $g_\Delta(a, w)$:

$$\text{logit } g_{\Delta, \gamma}^l(a, w) = \text{logit } \hat{g}_\Delta^l(a, w) + \gamma H^l(a, w),$$

by regressing Δ on A, W among all participants $i = 1, \dots, n$ using a logistic regression model with single covariate $H^l(A, W)$ and offset $\text{logit } \hat{g}_\Delta^l(A, W)$; denote the corresponding maximum likelihood estimate of γ by $\hat{\gamma}$.

Define $\hat{h}^{l+1} = h_{\hat{\epsilon}}^l$, $\hat{g}_A^{l+1} = g_{A, \hat{\nu}}^l$, and $\hat{g}_\Delta^{l+1} = g_{\Delta, \hat{\gamma}}^l$,

Step 3. Update $l = l + 1$ and iterate the previous step until convergence. Analogous to the TMLE algorithm in Section 4, we stop at the first iteration for which the mean of the squared difference of predictions between step l and step $l + 1$ is smaller or equal to $10^{-4}/n$, but this time adding the requirement that this holds for \hat{g}_Δ^l and \hat{g}_A^l as well.

Denote \hat{h}^* , \hat{g}_A^* , and \hat{g}_Δ^* the estimators obtained in the last iteration of the above algorithm, and define the corresponding TMLE estimator of $\theta_a(k)$ as

$$\hat{\theta}_{\text{adj,eff},a}(k) = \frac{1}{n} \sum_{i=1}^n \hat{h}^*(k, a, W_i) \prod_{j=0}^{k-1} (1 - \hat{h}^*(j, a, W_i)), \quad (20)$$

and the corresponding estimator $\hat{\beta}_{\text{adj,eff}} = \beta(\hat{\theta}_{\text{adj,eff}})$. Next, we present the main result of the paper giving conditions under which $\hat{\beta}_{\text{adj,eff}}$ is guaranteed to be at least as efficient as $\hat{\beta}_{\text{unadj}} = \beta(\hat{\theta}_{\text{unadj}})$.

Theorem 2 (Guaranteed equal or greater asymptotic efficiency of $\hat{\beta}_{\text{adj,eff}}$ compared to $\hat{\beta}_{\text{unadj}}$). *Assume that MCAR and the randomization assumption hold. Assume also that \hat{h} , \hat{g}_A , and \hat{g}_Δ are estimated using maximum likelihood estimation in parametric models containing at least an intercept term. Under the additional assumption that the empirical means of the estimated efficient influence functions $S(O)$ and $S_{\hat{\beta}_{\text{unadj}}}(O)$ are $o_P(1/\sqrt{n})$ (presented formally in Appendix 2), we have*

$$\begin{aligned} \sqrt{n}(\hat{\beta}_{\text{adj,eff}} - \beta) &\rightarrow N(0, \sigma_{\text{adj,eff}}^2) \\ \sqrt{n}(\hat{\beta}_{\text{unadj}} - \beta) &\rightarrow N(0, \sigma_{\text{unadj}}^2), \end{aligned}$$

where $\sigma_{\text{adj,eff}}^2 \leq \sigma_{\text{unadj}}^2$. In addition, if \hat{h}^* converges to h in $L^2(P_0)$, then $\hat{\beta}_{\text{adj,eff}}$ achieves the efficiency bound in the semiparametric model \mathcal{M} .

A sketch of the proof along with the technical conditions is provided in Appendix 2. By construction of the TMLE algorithm, it is expected that the empirical mean of the estimated efficient influence function will converge to zero, thereby satisfying the assumptions of the theorem.

6 Simulation Study: The MISTIE Trial of a New Surgical Treatment for Stroke

6.1 Simulation Design

Our simulated distributions are based on resampling data from the completed MISTIE II trial described in Section 2. This was done in order to make our simulations realistic, in that they mimic key features of a real trial. We compare the performance of the unadjusted, IPW, and proposed targeted minimum loss based estimator $\hat{\beta}_{\text{adj,eff}}$ (referred to as TMLE below). The baseline variables for each participant are $W = (W_1, W_2, W_3) = (\text{age, ICH volume, NIHSS})$ and we assume there is no missing data; i.e., $\Delta = 1$ for all subjects. Below, we refer to β as the treatment effect.

Simulations are conducted under the following four different types of data generating distributions (called scenarios):

- (i) Y and W dependent; zero treatment effect ($\beta = 0$).
- (ii) Y and W dependent; positive treatment effect ($\beta > 0$).
- (iii) Y and W independent; zero treatment effect ($\beta = 0$).
- (iv) Y and W independent; positive treatment effect ($\beta > 0$).

In scenarios (i) and (iii), we set $\beta = 0$ and in scenarios (ii) and (iv), the simulated data is generated to replicate the estimated treatment effect reported in the MISTIE II trial, which is 0.252 (which equals the unadjusted estimator directly applied to the trial data). In scenarios (i) and (ii), W is prognostic for Y , and there is potential for efficiency gains from adjusting for W . In scenarios (iii) and (iv), baseline variables are not prognostic for the outcome; we consider these scenarios in order to examine how much efficiency loss occurs for the adjusted estimators when the baseline variables are pure noise.

For each scenario, we generated 100,000 simulated randomized trials, each with $n = 412$ participants. This sample size was selected to mimic the projected size of the planned Phase III trial that is a follow-up to the MISTIE II trial.

The primary outcome, mRS at 180 days after enrollment, is ordinal and takes integer values from 0 to 6, with 0 representing no disability and 6 representing death. In the Phase II MISTIE trial, the primary outcome was always at least 1 and only one subject had mRS = 1. We combined the values 0,1,2 into a single category in our analysis, denoted by 0-2.

Each simulated randomized trial is based on resampling triples (W, A, Y) with replacement from the MISTIE II trial data, followed by modifications described below. These modifications were necessary since directly resampling with replacement corresponds to the empirical distribution from the MISTIE II trial, and in this data set the variables A and W have small, non-zero correlations. Therefore, directly resampling would correspond to a data generating distribution P_0 in which A and W are dependent, which violates the randomization assumption from Section 3.2 that under P_0 , the study arm A is assigned independent of W . We next describe modifications to the empirical distribution of the MISTIE II data that ensure the randomization assumption holds for the data generating distributions in our simulations.

For scenarios (i) and (ii), we resampled pairs (W, Y) with replacement from the MISTIE II trial data. This preserves the correlations between baseline variables and the outcome. In scenario (i), we generated A independent of (W, Y) , with probability 1/2 of being treatment or control. The resulting treatment effect for this data generating distribution is $\beta = 0$.

Each simulated data set in scenario (ii) consists of first generating an initial data set as in scenario (i), and then reassigning a subset of the outcome values using the algorithm described next. In order to induce a treatment effect that mimics the effect from the MISTIE II trial ($\beta = 0.252$), for each participant in the initial data set with $A = 1$ we randomly reassign the individual's outcome with probability that is a pre-computed function $q(Y)$ of his/her initial value of Y . If an individual's outcome is reassigned, it is replaced with a randomly generated outcome according to a pre-computed multinomial distribution. This procedure is done for each participant independent of the other participants. The pre-computed distributions, which were constructed to mimic features in the MISTIE II data, are given in Appendix 3.

For scenarios (iii) and (iv), baseline variables W for each participant were randomly drawn with replacement from the MISTIE II trial data. This results in the marginal distribution of W being the empirical distribution of the MISTIE II trial data. Study arm assignment A was generated independent of W , with probability $1/2$ of being treatment or control. In scenario (iii), Y is a random draw from a multinomial distribution with probabilities 0.12, 0.20, 0.27, 0.16, 0.25 for mRS of 0-2, 3, 4, 5, 6, respectively (which equal the corresponding empirical probabilities when pooling all participants in MISTIE II). Since A is generated independent of (W, Y) , this results in $\beta = 0$. In scenario (iv), the conditional distribution of Y given $A = a$ and W is set to be a multinomial distribution with probabilities 0.11, 0.14, 0.32, 0.16, 0.27 and 0.12, 0.24, 0.24, 0.16, 0.24, respectively for $a = 0, 1$, which correspond to the empirical probabilities in each arm in MISTIE II.

Logistic regression models are used to estimate g_A in the IPW estimator, and are used in the initial estimators in step 1 of the TMLE for both g_A and h . The logistic regression model for g_A includes an intercept and a main term for each component of W . Since treatment is randomized in all scenarios, we have $P(A = 1|W) = 1/2$, which implies the model for g_A is correctly specified. The logistic regression model for h includes intercepts for the first $K - 1$ levels of Y , main terms W and A , and interaction terms for the $K - 1$ intercepts with A and W .

In scenarios (i) and (ii), the simulation distribution of Y given A, W includes correlation between Y and W based on resampling from the MISTIE data; therefore, non-saturated parametric models for h such as those used in the TMLE estimator will generally be misspecified. We intentionally generate data in this way, to mimic the realistic feature that h will be at least somewhat misspecified in practice. In contrast, for scenarios (iii) and (iv), since Y and W are independent, the parametric model for h above is correctly specified.

6.2 Simulation Results

The results of the simulation are presented in Table 1. In all scenarios, each estimator has very small bias. In scenarios (i) and (ii), there are large gains in efficiency for both the IPW estimator and TMLE relative to the unadjusted estimator. These gains are due to the strong correlation between W and Y . The gains in efficiency are 4% greater for the TMLE relative to the IPW estimator. In scenarios (iii) and (iv), where W is not prognostic for Y , there are efficiency losses between 0.7% and 0.9% from using the adjusted estimators.

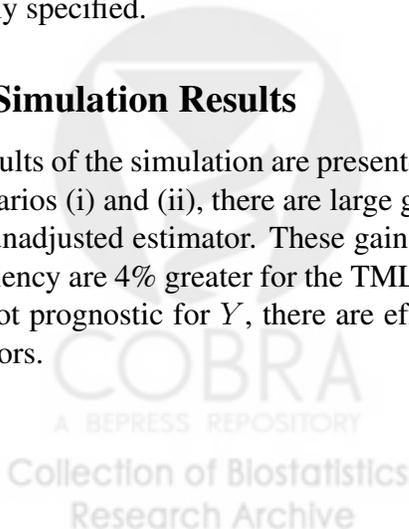


Table 1: Results of the simulation study where Y is multinomial taking on values $k = 1, \dots, 5$, (Y, W) are simulated by resampling with replacement from the MISTIE II trial for a study size of 412 patients. The relative efficiency is the ratio of the variance of the unadjusted estimator to the variance of the given estimator. Values greater than 1 indicate an efficiency gain compared to the unadjusted estimator.

Scenario	Estimator	Bias	Variance	Relative Efficiency
(i)	Unadjusted	0.000	0.033	1.000
	TMLE	-0.001	0.024	1.392
	IPW	0.000	0.024	1.353
(ii)	Unadjusted	0.004	0.035	1.000
	TMLE	-0.012	0.026	1.354
	IPW	0.004	0.027	1.317
(iii)	Unadjusted	0.001	0.033	1.000
	TMLE	0.001	0.034	0.991
	IPW	0.001	0.034	0.992
(iv)	Unadjusted	0.004	0.034	1.000
	TMLE	0.004	0.034	0.992
	IPW	0.004	0.034	0.993

7 Discussion

The method we present may be extended to other cases. In particular, observational data (where A and W are generally dependent) may be handled using the estimator $\hat{\beta}_{\text{adj}}$, under the ignorability assumption from Section 3.2. Under certain regularity conditions, $\hat{\beta}_{\text{adj}}$ is a doubly robust, locally efficient estimator of β , which identifies the causal effect. We conjecture $\hat{\beta}_{\text{adj,eff}}$ has asymptotic variance smaller or equal than an inverse probability weighted estimator computed with g_0 known, under stronger assumptions that include consistent estimation of g_0 at an appropriate rate, and other regularity conditions.

In addition, the natural formulation of the problem in terms of hazard functions allows the straightforward use of the methods to handle discrete time-to-event outcomes with no time varying covariates. This is done by considering a categorical outcome taking values in an ordered set of discrete time points $\{1, \dots, \tau\}$. In this case, our parameter defines a contrast between the two survival functions, and our method can be applied to estimate this parameter as well as the survival functions after accounting for baseline covariates. The method we describe may be generalized to estimate other parameters of interest relevant to time-to-event outcomes such as the relative hazard.

A Supplementary Materials

A.1 Derivation of Hazard Representation of the Efficient Influence Function in Theorem 1

First, the component of the tangent space generated by scores for $P(Y|\Delta = 1, A, W)$, which is the component of the likelihood that is relevant for the parameter $p(k, a, w)$, $T_Y = \{S(Y, A, W) : E(S|\Delta =$

$1, A, W) = 0\}$, can be decomposed as

$$\bigoplus_{k=1}^{K-1} T_{I(k)},$$

where $T_{I(k)} = \{S(I_1, \dots, I(k), A, W) : E(S|I(1), \dots, I(k-1), \Delta = 1, A, W) = 0\}$. Next, we take the projection of

$$D_Y(O) = \frac{\Delta \mathbf{1}\{A = a\}}{g_A(a|W)g_\Delta(a, W)} (\mathbf{1}\{Y \leq k\} - p(k, a, W)),$$

onto $T_{I(j)}$ for $j \leq k$. This projection is equal to

$$\begin{aligned} \Pi[D_Y(O)|T_{I(j)}] &= \frac{\Delta \mathbf{1}\{A = a\}}{g_A(a|W)g_\Delta(a, W)} \left(E(\mathbf{1}\{Y \leq k\}|I(1), \dots, I(j), A = a, W) \right. \\ &\quad \left. - E(\mathbf{1}\{Y \leq k\}|I(1), \dots, I(j-1), A = a, W) \right). \end{aligned}$$

Note that

$$\begin{aligned} E(\mathbf{1}\{Y \leq k\}|I(1), \dots, I(j), A = a, W) &= 1 - \bar{I}(j+1)P(Y > k|Y > j, A = a, W) \\ &= 1 - (1 - I(j))\bar{I}(j) \frac{P(Y > k|A = a, W)}{P(Y > j|A = a, W)}. \end{aligned}$$

This yields

$$E(\mathbf{1}\{Y \leq k\}|I(1), \dots, I(j-1), A = a, W) = 1 - (1 - h(j, a, W))\bar{I}(j) \frac{P(Y > k|A = a, W)}{P(Y > j|A = a, W)},$$

which implies

$$\Pi[D_Y(O)|T_{I(j)}] = \bar{I}(j)(I(j) - h(j, a, W)) \frac{\Delta \mathbf{1}\{A = a\}}{g_A(a|W)g_\Delta(a, W)} \frac{P(Y > k|A = a, W)}{P(Y > j|A = a, W)}.$$

For $j > k$ we have $\Pi[D_Y(O)|T_{I(j)}] = 0$. Because $T_Y = \bigoplus_{k=1}^{K-1} T_{I(k)}$, we have

$$D_Y(O) = \sum_{j=1}^K \Pi[D_Y(O)|T_{I(j)}],$$

and the result follows.

A.2 Sketch of Proof for Theorem 2

In this section we use the notation Pf to refer to $\int f(o)dP(o)$, for a given function f of the data. We also add an index naught to the notation to represent true quantities (e.g., $g_{A,0}$), and use no index to denote generic quantities (e.g., g_A). In addition, we augment the notation of the paper by explicitly including the dependence on the influence functions D_{ak} , D'_{ak} , and S on P . Thus, $D_{ak}(O)$ becomes $D_{ak}(P)(O)$, for example.

Define the following two quantities:

$$\begin{aligned}\phi_{ak}(p, g)(O) &= \frac{p(k, a, W) - \theta_a(k)}{g_{A,0}(a|W)g_{\Delta,0}(a|W)}(\Delta \mathbf{1}\{A = a\} - g_A(a|W)g_{\Delta}(a|W)) \\ &= H_{ak}(g_0, p)(A, W)(\Delta - g_{\Delta}(a, W)) + M_{ak}(g_0)(W)(\mathbf{1}\{A = a\} - g_A(a, W)); \\ \phi(p, g) &= \sum_{a \in \{0,1\}} \sum_{k=1}^{K-1} \frac{d\beta(\theta_0)}{d\theta_{a,0}(k)} \phi_{ak}(p, g).\end{aligned}$$

Assumptions in the theorem:

1. Randomization assumption and missing completely at random assumption (defined in Sections 3.2 and 3.3).
2. The following assumption:

$$P_n S_{\hat{\beta}_{\text{unadj}}}(\hat{g}^*, \hat{\theta}_{\text{adj,eff}}) = o_P(1/\sqrt{n}).$$

Intuitively, we expect that the above assumption will hold in many cases, due to the construction of the TMLE. This is because the efficient influence function S may be decomposed as $S_{\hat{\beta}_{\text{unadj}}} + \phi$, where ϕ is defined above. The TMLE is constructed to solve the estimating equations associated with S and with ϕ . As a result, the estimating equation associated with $S_{\hat{\beta}_{\text{unadj}}}$ is also solved and the above assumption may be expected to hold.

3. Consistency of initial estimators of outcome hazard. (This assumption is only used for the claim in the last sentence of the theorem.)

$$P_0(\hat{h}^* - h)^2 = o_P(1).$$

Lemmas:

1. Consistency of initial estimators of treatment/missingness. Under assumption (1) we have

$$\begin{aligned}\sqrt{n}P_0(\hat{g}_A^* - g_A)^2 &= o_P(1), \\ \sqrt{n}P_0(\hat{g}_{\Delta}^* - g_{\Delta})^2 &= o_P(1).\end{aligned}$$

This follows since assumption (1) implies that the parametric models used for g_A and g_{Δ} are correctly specified. Applying the delta method and the continuous mapping theorem, we have that $\sqrt{n}P(\hat{g}_A^*(a, w) - g_A(a, w)) = O_P(1)$, and $\hat{g}_A^*(a, w) - g_A(a, w) = o_P(1)$, for any fixed a and w . This, together with $O_P(1)o_P(1) = o_P(1)$ and the continuous mapping theorem proves the lemma.

2. Asymptotic linearity of smooth functional of g . Assume

$$P_0(\phi(p_0, \hat{g}^*) - \phi(p_0, g_0)) = (P_n - P_0)S^{\dagger} + o_P(1/\sqrt{n}),$$

for some function S^{\dagger} of O . This is a direct consequence of the asymptotic linearity of the MLE in a parametric model and the delta method.

Claims:

1. Under assumptions (1,2) we have

$$\begin{aligned}\sqrt{n}(\hat{\beta}_{\text{adj,eff}} - \beta) &\rightarrow N(0, \sigma_{\text{adj,eff}}^2), \\ \sqrt{n}(\hat{\beta}_{\text{unadj}} - \beta) &\rightarrow N(0, \sigma_{\text{unadj}}^2),\end{aligned}$$

where $\sigma_{\text{adj,eff}}^2 \leq \sigma_{\text{unadj}}^2$.

2. Under assumptions (1,2,3), we have

$$\sqrt{n}(\hat{\beta}_{\text{adj,eff}} - \beta) \rightarrow N(0, \sigma^2),$$

where σ^2 is the efficiency bound for estimation of β in the nonparametric model.

Sketch of proof:

1. Claim 1. Let

$$D_{\text{unadj},a,k}(g, \theta)(O) = \frac{\Delta \mathbf{1}\{A = a\}}{P(A = a, \Delta = 1)} [\mathbf{1}\{Y \leq k\} - \theta_a(k)],$$

be the influence function of $\hat{\theta}_{\text{unadj},a}(k)$, where $\theta_a(k)$ is the corresponding component of θ . It follows that

$$P_0 D_{\text{unadj},a,k}(g, \theta) = \theta_{a,0}(k) - \theta_a(k) + P_0 \phi_{ak}(p_0, g) + o(\|g - g_0\|^2). \quad (21)$$

Note that

$$\beta(\theta) - \beta(\theta_0) = \sum_{a \in \{0,1\}} \sum_{k=1}^{K-1} \frac{d\beta(\theta_0)}{d\theta_{a,0}(k)} (\theta_a(k) - \theta_{a,0}(k)) + o(\|\theta - \theta_0\|^2).$$

Solving (21) for $(\theta_a(k) - \theta_{a,0}(k))$, and noting that

$$S_{\hat{\beta}_{\text{unadj}}}(g, \theta) = \sum_{k=1}^{K-1} \frac{d\beta(\theta_0)}{d\theta_{a,0}(k)} D_{\text{unadj},a,k}(g, \theta)$$

we get

$$\beta(\theta) - \beta(\theta_0) = -P_0 S_{\hat{\beta}_{\text{unadj}}}(g, \theta) + P_0 \phi(p_0, g) + o(\|\theta - \theta_0\|^2) + o(\|g - g_0\|^2). \quad (22)$$

Using equation (22) with $g = \hat{g}^*$ and $\theta = \hat{\theta}_{\text{adj,eff}}$, we obtain

$$\hat{\beta}_{\text{adj,eff}} - \beta_0 = -P_0 S_{\hat{\beta}_{\text{unadj}}}(\hat{g}^*, \hat{\theta}_{\text{adj,eff}}) + P_0 \phi(p_0, \hat{g}^*) + o_P(\|\hat{\theta}_{\text{adj,eff}} - \theta_0\|^2) + o_P(\|\hat{g}^* - g_0\|^2).$$

By Lemma 1 above, we have $\|\hat{g}^* - g_0\|^2 = o_P(1/\sqrt{n})$. Since the TMLE is doubly robust, we have $\|\hat{\theta}_{\text{adj,eff}} - \theta_0\|^2 = o_P(1/\sqrt{n})$. In addition, by assumption (2), we have

$$P_n S_{\hat{\beta}_{\text{unadj}}}(\hat{g}^*, \hat{\theta}_{\text{adj,eff}}) = o_P(1/\sqrt{n}),$$

so that we can write

$$\begin{aligned}\hat{\beta}_{\text{adj,eff}} - \beta_0 &= (P_n - P_0)S_{\hat{\beta}_{\text{unadj}}}(g_0, \theta_0) + \\ &\quad (P_n - P_0)(S_{\hat{\beta}_{\text{unadj}}}(\hat{g}^*, \hat{\theta}_{\text{adj,eff}}) - S_{\hat{\beta}_{\text{unadj}}}(g_0, \theta_0)) + P_0\phi(p_0, \hat{g}^*) + o_P(1/\sqrt{n}).\end{aligned}$$

Since \hat{g}^* is a parametric model fit, it belongs to a Donsker class. In addition, since \hat{g}^* is correctly specified, it is consistent. As a result of the continuous mapping theorem we have $P_0(S_{\hat{\beta}_{\text{unadj}}}(\hat{g}^*, \hat{\theta}_{\text{adj,eff}}) - S_{\hat{\beta}_{\text{unadj}}}(g_0, \theta_0))^2 = o_P(1)$. The empirical process result in Theorem 19.24 of [26] then implies

$$\hat{\beta}_{\text{adj,eff}} - \beta_0 = (P_n - P_0)S_{\hat{\beta}_{\text{unadj}}}(g_0, \theta_0) + P_0\phi(p_0, \hat{g}^*) + o_P(1/\sqrt{n}).$$

By Lemma 2 above we have

$$\hat{\beta}_{\text{adj,eff}} - \beta_0 = (P_n - P_0)(S_{\hat{\beta}_{\text{unadj}}}(g_0, \theta_0) - S^\dagger) + o_P(1/\sqrt{n}).$$

We can now use the arguments in the proof of Theorem 2.3 of [22] to show that

$$S_{\hat{\beta}_{\text{unadj}}}(g_0, \theta_0) - S^\dagger = \Pi(S_{\hat{\beta}_{\text{unadj}}}(g_0, \theta_0) | T_{\Delta, A}^\perp), \quad (23)$$

where $T_{\Delta, A}^\perp$ is the orthogonal complement (in $L_0^2(P)$) of $T_{\Delta, A}$, the tangent space of the parametric model for $P(A = a, \Delta = \delta | W = w)$ used in the TMLE algorithm.

The argument in the proof of the aforementioned theorem is as follows. Let $S_{\hat{\beta}_{\text{unadj}}}(g_0, \theta_0)$ and S^\dagger be decomposed as $S_{\hat{\beta}_{\text{unadj}}}(g_0, \theta_0) = a + a^\perp$ and $S^\dagger = b + b^\perp$ according to the orthogonal decomposition $L_0^2(P_0) = T_{A, \Delta} + T_{A, \Delta}^\perp$. Since S^\dagger is an efficient influence function in the model with $(p_0, P_{W,0})$ known, $b^\perp = 0$. We also have $S_{\hat{\beta}_{\text{unadj}}}(g_0, \theta_0) - S^\dagger$ is an influence function of a regular, asymptotically linear estimator of β in the model \mathcal{M} , and hence is orthogonal to $T_{A, \Delta}$, so $a - b = 0$. Consequently, we have $S_{\hat{\beta}_{\text{unadj}}}(g_0, \theta_0) - S^\dagger = a^\perp$, completing the argument.

As a consequence of (23), we have

$$\|S_{\hat{\beta}_{\text{unadj}}}(g_0, \theta_0) - S^\dagger\|^2 = V(S_{\hat{\beta}_{\text{unadj}}}(g_0, \theta_0) - S^\dagger) \leq V(S_{\hat{\beta}_{\text{unadj}}}(g_0, \theta_0)) = \|S_{\hat{\beta}_{\text{unadj}}}(g_0, \theta_0)\|^2,$$

which together with

$$\sqrt{n}(\hat{\beta}_{\text{unadj}} - \beta) = \frac{1}{\sqrt{n}} \sum_{i=1}^n S_{\hat{\beta}_{\text{unadj}}}(O_i) + o_P(1),$$

and the central limit theorem completes the sketch of the proof.

2. Claim 2 follows from Appendix 18 of [23].

A.3 Data Generating Distribution for Simulation Scenario (ii)

In Section 6.1 of the main paper, it was stated that each simulated data set in scenario (ii) consists of first generating an initial data set as in scenario (i), and then reassigning a subset of the outcome values. Specifically, for each participant in the initial data set with $A = 1$ we randomly reassign the individual's outcome with probability that is a pre-computed function $q(Y)$ of his/her initial value of Y . We set $q(Y) = 0, 0, 0.44, 0.06, 0.10$ for Y in categories 0-2, 3, 4, 5, 6, respectively. If an individual's outcome is reassigned, it is replaced with a randomly generated outcome according to the pre-computed multinomial distribution, denoted m , on the same set of categories, with probabilities 0.10, 0.90, 0, 0, 0 for the categories 0-2, 3, 4, 5, 6, respectively.

The above values were selected in order that the simulated distribution in scenario (ii) has $\beta = 0.252$, which is the estimated treatment effect (unadjusted) in the MISTIE II trial. The above values were obtained by working backward from β to $\theta_a(k)$, as we describe next. Recall that

$$\beta(\theta) = \frac{1}{K-1} \sum_{k=1}^{K-1} \log \left\{ \frac{\theta_1(k)}{1-\theta_1(k)} \bigg/ \frac{\theta_0(k)}{1-\theta_0(k)} \right\}. \quad (24)$$

We select q and m in order that the resulting simulation distribution has

$$\log \left\{ \frac{\theta_1(k)}{1-\theta_1(k)} \bigg/ \frac{\theta_0(k)}{1-\theta_0(k)} \right\} \quad (25)$$

equal to the corresponding values where each $\theta_a(k)$ is replaced by the unadjusted estimator of that quantity in the MISTIE II trial. For categories 0-2, 3, 4, 5, the corresponding value of (25) is set to 0.129, 0.575, 0.159, 0.146. This implies $\beta(\theta) = (0.129 + 0.575 + 0.159 + 0.146)/4 = 0.252$. Since the distribution of $Y|A = 0, W$ is already determined (see the construction of the data generating mechanism for scenario (i)), we can compute $\theta_0(k)$, and then compute $\theta_1(k)$ using the aforementioned values of (25). We then experimented with values of q and m that produce precisely the aforementioned values of $\theta_a(k)$. Reproducing such values of $\theta_a(k)$ mimics the heterogenous log odds ratios across categories in the MISTIE II trial.

References

- [1] Emun Abdu, Daniel F Hanley, and David W Newell. Minimally invasive treatment for intracerebral hemorrhage. *Neurosurgical Focus*, 32(4):E3, 2012.
- [2] Heejung Bang and James M Robins. Doubly robust estimation in missing data and causal inference models. *Biometrics*, 61(4):962–973, 2005.
- [3] P.J. Bickel, C.A.J. Klaassen, Y. Ritov, and J. Wellner. *Efficient and Adaptive Estimation for Semiparametric Models*. Springer-Verlag, 1997.
- [4] T. Brott, H.P. Adams, C.P. Olinger, J. Marler, W.G. Barsan, J. Biller, J. Spilker, R. Holleran, R. Eberle, V. Hertzberg, M. Rorick, C.J. Moomaw, and M. Walker. Measurements of acute cerebral infarction: a clinical examination scale. *Stroke*, 20(7):864–70, 1989.

- [5] Weihua Cao, Anastasios A Tsiatis, and Marie Davidian. Improving efficiency and robustness of the doubly robust estimator for a population mean with incomplete data. *Biometrika*, 96(3):723–734, 2009.
- [6] Bradley Efron. Better bootstrap confidence intervals. *Journal of the American Statistical Association*, 82(397):171–185, 1987.
- [7] David A Freedman. On regression adjustments to experimental data. *Advances in Applied Mathematics*, 40(2):180–193, 2008.
- [8] Susan Gruber and Mark J. van der Laan. Targeted minimum loss based estimator that outperforms a given estimator. *The International Journal of Biostatistics*, 8(1):1–22, 2012.
- [9] J. Kang and J. Schafer. Demystifying double robustness: A comparison of alternative strategies for estimating a population mean from incomplete data (with discussion). *Statistical Science*, 22: 523–39, 2007.
- [10] Kelly L Moore and Mark J van der Laan. Covariate adjustment in randomized trials with binary outcomes: Targeted maximum likelihood estimation. *Statistics in Medicine*, 28(1):39–64, 2009.
- [11] R. Neugebauer and M. J. van der Laan. Nonparametric causal effects based on marginal structural models. *Journal of Statistical Planning & Inference*, 137(2):419 – 434, 2007. ISSN 0378-3758. doi: DOI: 10.1016/j.jspi.2005.12.008.
- [12] J. Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, Cambridge, 2000.
- [13] J.M. Robins and A. Rotnitzky. Recovery of information and adjustment for dependent censoring using surrogate markers. In *AIDS Epidemiology, Methodological issues*. Birkhäuser, 1992.
- [14] J.M. Robins, A. Rotnitzky, and L.P. Zhao. Estimation of regression coefficients when some regressors are not always observed. *Journal of the American Statistical Association*, 89(427):846–866, September 1994.
- [15] Andrea Rotnitzky, Quanhong Lei, Mariela Sued, and James M Robins. Improved double-robust estimation in missing data and causal inference models. *Biometrika*, 99(2):439–456, 2012.
- [16] Donald B Rubin. *Multiple Imputation for Nonresponse in Surveys*. John Wiley & Sons, 1987.
- [17] Daniel O. Scharfstein, Andrea Rotnitzky, and James M. Robins. Adjusting for nonignorable drop-out using semiparametric nonresponse models: Rejoinder. *Journal of the American Statistical Association*, 94(448):pp. 1135–1146, 1999. ISSN 01621459.
- [18] Ori M Stitelman, Victor De Gruttola, and Mark J van der Laan. A general implementation of tmle for longitudinal data applied to causal inference in survival analysis. *The International Journal of Biostatistics*, 8(1), 2011.
- [19] Zhiqiang Tan. A distributional approach for causal inference using propensity scores. *Journal of the American Statistical Association*, 101(476):1619–1637, 2006.

- [20] Zhiqiang Tan. Bounded, efficient and doubly robust estimation with inverse weighting. *Biometrika*, 97(3):661–682, 2010.
- [21] Anastasios A Tsiatis, Marie Davidian, Min Zhang, and Xiaomin Lu. Covariate adjustment for two-sample treatment comparisons in randomized clinical trials: A principled yet flexible approach. *Statistics in Medicine*, 27(23):4658–4677, 2008.
- [22] M.J. van der Laan and J.M. Robins. *Unified Methods for Censored Longitudinal Data and Causality*. Springer, New York, 2003.
- [23] M.J. van der Laan and S. Rose. *Targeted Learning: Causal Inference for Observational and Experimental Data*. Springer, New York, 2011.
- [24] M.J. van der Laan and D. Rubin. Targeted maximum likelihood learning. *The International Journal of Biostatistics*, 2(1):Article 11, 2006.
- [25] M.J. van der Laan, E. Polley, and A. Hubbard. Super learner. *Statistical Applications in Genetics & Molecular Biology*, 6(25):Article 25, 2007.
- [26] A. W. van der Vaart. *Asymptotic Statistics*. Cambridge University Press, 1998.
- [27] David H Wolpert. Stacked generalization. *Neural Networks*, 5(2):241–259, 1992.
- [28] Min Zhang, Anastasios A Tsiatis, and Marie Davidian. Improving efficiency of inferences in randomized clinical trials using auxiliary covariates. *Biometrics*, 64(3):707–715, 2008.

