# Optimizing Precision and Power by Machine Learning in Randomized Trials, with an Application to COVID-19

Nicholas Williams[1], Michael Rosenblum[2], and Iván Díaz[*1]

[1]Division of Biostatistics, Department of Population Health Sciences, Weill Cornell Medicine.
[2]Department of Biostatistics, Johns Hopkins Bloomberg School of Public Health.

September 10, 2021

## Abstract

The rapid finding of effective therapeutics requires the efficient use of available resources in clinical trials. The use of covariate adjustment can yield statistical estimates with improved precision, resulting in a reduction in the number of participants required to draw futility or efficacy conclusions. We focus on time-to-event and ordinal outcomes. When more than a few baseline covariates are available, a key question for covariate adjustment in randomized studies is how to fit a model relating the outcome and the baseline covariates to maximize precision. We present a novel theoretical result establishing conditions for asymptotic normality of a variety of covariate-adjusted estimators that rely on machine learning (e.g., $\ell_1$-regularization, Random Forests, XGBoost, and Multivariate Adaptive Regression Splines), under the assumption that outcome data is missing completely at random. We further present a consistent estimator of the asymptotic variance. Importantly, the conditions do not require the machine learning methods to converge to the true outcome distribution conditional on baseline variables, as long as they converge to some (possibly incorrect) limit. We conducted a simulation study to evaluate the performance of the aforementioned prediction methods in COVID-19 trials using longitudinal data from over 1,500 patients hospitalized with COVID-19 at Weill Cornell Medicine New York Presbyterian Hospital. We found that using $\ell_1$-regularization led to estimators and corresponding hypothesis tests that control type 1 error and are more precise than an unadjusted estimator across all sample sizes tested. We also show that when covariates are not prognostic of the outcome, $\ell_1$-regularization remains as precise as the unadjusted estimator, even at small sample sizes ($n = 100$). We give an R package adjrct that performs model-robust covariate adjustment for ordinal and time-to-event outcomes.

[*]corresponding author: ild2005@med.cornell.edu

# 1 Introduction

Coronavirus disease 2019 (COVID-19) has affected more than 125 million people and caused more than 2.7 million deaths worldwide (World Health Organization 2021). Governments and scientists around the globe have deployed an enormous amount of resources to combat the pandemic with remarkable success, such as the development in record time of highly effective vaccines to prevent disease (e.g., Polack et al. 2020; Baden et al. 2021). Global and local organizations are launching large-scale collaborations to collect robust scientific data to test potential COVID-19 treatments, including the testing of drugs re-purposed from other diseases as well as new compounds (Kupferschmidt and Cohen 2020). For example, the World Health Organization launched the SOLIDARITY trial, enrolling almost 12,000 patients in 500 hospital sites in over 30 countries (WHO Solidarity Trial Consortium 2021). Other large initiatives include the RECOVERY trial (The RECOVERY Collaborative Group 2021) and the ACTIV initiative (Collins and Stoffels 2020). To date, there are approximately 2,400 randomized trials for the treatment of COVID-19 registered in `clinicaltrials.gov`.

The rapid finding of effective therapeutics for COVID-19 requires the efficient use of available resources. One area where such efficiency is achievable at little cost is in the statistical design and analysis of the clinical trials. Specifically, a statistical technique known as *covariate adjustment* may yield estimates with increased precision (compared to unadjusted estimators), and may result in a reduction of the time, number of participants, and resources required to draw futility or efficacy conclusions. This results in faster trial designs, which may help accelerate the delivery of effective treatments to patients who need them (and may help rule out ineffective treatments faster).

Covariate adjustment refers to pre-planned analysis methods that use data on patient baseline characteristics to correct for chance imbalances across study arms, thereby yielding more precise treatment effect estimates. The ICH E9 Guidance on Statistical Methods for Analyzing Clinical Trials (FDA and EMA 1998) states that "Pretrial deliberations should identify those covariates and factors expected to have an important influence on the primary variable(s), and should consider how to account for these in the analysis to improve precision and to compensate for any lack of balance between treatment groups." Even though its benefits can be substantial, covariate adjustment is underutilized; only 24%-34% of trials use covariate adjustment (Kahan et al. 2014).

We focus on estimation of marginal treatment effects, defined as a contrast between study arms in the marginal distribution of the outcome. Many approaches for estimation of marginal treatment effects using covariate adjustment in randomized trials invoke a model relating the outcome and the baseline covariates within strata of treatment. Recent decades have seen a surge in research on the development of *model-robust* methods for estimating marginal effects that remain consistent even if this outcome regression model is arbitrarily misspecified (e.g., Yang and Tsiatis 2001; Tsiatis et al. 2008; Zhang et al. 2008; Moore and van der Laan 2009a; Austin et al. 2010; Zhang and Gilbert 2010; Benkeser et al. 2020). We focus on a study of the model-robust covariate adjusted estimators for time-to-event and ordinal outcomes developed by Moore and van der Laan (2009a), Díaz et al. (2019), and Díaz et al. (2016).

All potential adjustment covariates must be pre-specified in the statistical analysis plan. At the end of the trial, a prespecified prediction algorithm (e.g., random forests, or using

regularization for variable selection) will be run and its output used to construct a model-robust, covariate adjusted estimator of the marginal treatment effect for the trial's primary efficacy analysis. We aim to address the question of how to do this in a model-robust way that guarantees consistency and asymptotic normality, under some weaker regularity conditions than related work (described below). We also aim to demonstrate the potential value added by covariate adjustment combined with machine learning, through a simulation study based on COVID-19 data.

As a standard regression method for high-dimensional data, $\ell_1$-regularization has been studied by several authors in the context of covariate selection for randomized studies. For example, Wager et al. (2016) present estimators that are asymptotically normal under strong assumptions that include linearity of the outcome-covariate relationship. Bloniarz et al. (2016) present estimators under a randomization inference framework, and show asymptotic normality of the estimators under assumptions similar to the assumptions made in this paper. Both of these papers present results only for continuous outcomes. The method of Tian et al. (2012) is general and can be applied to continuous, ordinal, binary, and time-to-event data, and its asymptotic properties are similar to the properties of the methods we discuss for the case of $\ell_1$-regularization, under similar assumptions.

More related to our general approach, Wager et al. (2016) also present a cross-validation procedure that can be used with arbitrary non-parametric prediction methods (e.g., $\ell_1$-regularization, random forests, etc.) in the estimation procedure. Their proposal amounts to computation of a cross-fitted augmented inverse probability weighted estimator (Chernozhukov et al. 2018). Their asymptotic normality results, unlike ours, require that that their predictor of the outcome given baseline variables converges to the true regression function. Wu and Gagnon-Bartsch (2018) proposed a "leave-one-out-potential outcomes" estimator where automatic prediction can also be performed using any regression procedure such as linear regression or random forests, and they propose a conservative variance estimator. It is unclear as of yet whether Wald-type confidence intervals based on the normal distribution are appropriate for this estimator. As in the above related work that compares the precision of covariate adjusted estimators to the unadjusted estimator, we assume that outcomes are missing completely at random (since otherwise the unadjusted estimator is generally inconsistent).

In Section 3.3, we present our main theorem. It shows that any of a large class of prediction algorithms (e.g., $\ell_1$-regularization, Random Forests, XGBoost, and Multivariate Adaptive Regression Splines) can be combined with the covariate adjusted estimator of Moore and van der Laan (2009b) to produce a consistent, asymptotically normal estimator of the marginal treatment effect, under regularity conditions. These conditions do not require consistent estimation of the outcome regression function (as in key related work described above); instead, our theorem requires the weaker condition of convergence to some (possibly incorrect) limit. We also give a consistent, easy to compute variance estimator. This has important practical implications because it allows the use machine learning coupled with Wald-type confidence intervals and hypothesis tests, under the conditions of the theorem. The above estimator can be used with ordinal or time-to-event outcomes.

We next conduct a simulation study to evaluate the performance of the aforementioned machine learning algorithms for covariate adjustment in the context of COVID-19 trials. We simulate two-arm trials comparing a hypothetical COVID-19 treatment to standard of care. The simulated data distributions are generated from longitudinal data on approximately

3

1,500 patients hospitalized at Weill Cornell Medicine New York Presbyterian Hospital prior to 15 May 2020. We present results for two types of endpoints: time-to-event (e.g., time to intubation or death) and ordinal (e.g., WHO scale, see Marshall et al. 2020) outcomes. For survival outcomes, we present results for two different estimands (i.e., targets of inference): the survival probability at any given time and the restricted mean survival time. For ordinal outcomes we present results for the average log-odds ratio, and for the Mann-Whitney estimand, interpreted as the probability that a randomly chosen treated patient has a better outcome than a randomly chosen control patient (with ties broken at random).

Benkeser et al. (2020) used simulations based on the above data source to illustrate the efficiency gains achievable by covariate adjustment with parametric models including a small number of adjustment variables (and not using machine learning to improve efficiency). In this paper we evaluate the performance of four machine learning algorithms ($\ell_1$-regularization, Random Forests, XGBoost, and Multivariate Adaptive Regression Splines) in several sample sizes, and compare them in terms of their bias, mean squared error, and type-1 and type-2 errors, to unadjusted estimators and to fully adjusted main terms logistic regression with all available variables included. Furthermore, we introduce a new R package `adjrct` (Díaz and Williams 2021) that can be used to perform model-robust covariate adjustment for ordinal and time-to-event outcomes, and provide R code that can be used to replicate our simulation analyses with other data sources.

# 2   Estimands

In what follows, we focus on estimating *intention-to-treat* effects and refer to study arm assignment simply as *treatment*. We focus on estimation of marginal treatment effects, defined as a contrast between study arms in the marginal distribution of the outcome. We further assume that we have data on $n$ trial participants, represented by $n$ independent and identically distributed copies of data $O_i : i = 1, \ldots, n$. We assume $O_i$ is distributed as $\mathsf{P}$, where we make no assumptions about the functional form of $\mathsf{P}$ except that treatment is independent of baseline covariates (by randomization). We denote a generic draw from the distribution $\mathsf{P}$ by $O$. We use the terms "baseline covariate" and "baseline variable" interchangeably to indicate a measurement made before randomization.

We are interested in making inferences about a feature of the distribution $\mathsf{P}$. We use the word *estimand* to refer to such a feature. We describe example estimands, which include those used in our simulations studies, below.

## 2.1   Ordinal Outcomes

For ordinal outcomes, assume the observed data is $O = (W, A, Y)$, where $W$ is a vector of baseline covariates, $A$ is the treatment arm, and $Y$ is an ordinal variable that can take values in $\{1, \ldots, K\}$. Let $F(k, a) = \mathsf{P}(Y \leq k \mid A = a)$ denote the cumulative distribution function for patients in arm $A = a$, and let $f(k, a) = F(k, a) - F(k-1, a)$ denote the corresponding probability mass function. For notational convenience we will sometimes use the "survival"

function instead: $S(k, a) = 1 - F(k, a)$. The average log-odds ratio is then equal to

$$\mathsf{LOR} = \frac{1}{K-1} \sum_{k=1}^{K-1} \log \left[ \frac{F(k,1)/\{1 - F(k,1)\}}{F(k,0)/\{1 - F(k,0)\}} \right],$$

and the Mann-Whitney estimand is equal to

$$\mathsf{MW} = \sum_{k=1}^{K} \left\{ F(k-1, 0) + \frac{1}{2} f(k, 0) \right\} f(k, 1).$$

The Mann-Whitney estimand can be interpreted as the probability that a randomly drawn patient from the treated arm has a better outcome than a randomly drawn patient from the control arm, with ties broken at random (Ahmad 1996). The average log-odds ratio is more difficult to interpret and we discourage its use, but we include it in our comparisons because it is a non-parametric extension of the parameter $\beta$ estimated by the commonly used proportional odds model logit$\{F(k, a)\} = \alpha_k + \beta a$ (Díaz et al. 2016).

## 2.2 Time to Event Outcomes

For time to event outcomes, we assume the observed data is $O = (W, A, \Delta = \mathbb{1}\{Y \leq C\}, \widetilde{Y} = \min(C, Y))$, where $C$ is a right-censoring time denoting the time that a patient is last seen, and $\mathbb{1}\{E\}$ is the indicator variable taking the value 1 on the event $E$ and 0 otherwise. We further assume that events are observed at discrete time points $\{1, \ldots, K\}$ (e.g., days) as is typical in clinical trials. The difference in restricted mean survival time is given by

$$\mathsf{RMST} = \sum_{k=1}^{K-1} \{S(k, 1) - S(k, 0)\},$$

and can be interpreted as a contrast comparing the expected survival time within the first $K$ time units for the treated arm minus the control arm (Chen and Tsiatis 2001; Royston and Parmar 2011). The risk difference at a user-given time point $k$ is defined as

$$\mathsf{RD} = S(k, 1) - S(k, 0),$$

and is interpreted as the difference in survival probability for a patient in the treated arm minus the control arm. We note that the MW and RD parameters may be meaningful for both ordinal and time-to-event outcomes.

## 3 Estimators

For the sake of generality, in what follows we use a common data structure $O = (W, A, \Delta = \mathbb{1}\{Y \leq C\}, \widetilde{Y})$ for both ordinal and survival outcomes, where for ordinal outcomes $C = K$ if the outcome is observed and $C = 0$ if it is missing.

Many approaches for estimation of marginal treatment effects using covariate adjustment in randomized trials invoke a model relating the outcome and the baseline covariates within

5

strata of treatment. It is important that the consistency and interpretability of the treatment effect estimates do not rely on the ability to correctly posit such a model. Specifically, in a recent draft guidance (U.S. Food and Drug Administration 2021), the FDA states: "Sponsors can perform covariate adjusted estimation and inference for an unconditional treatment effect ... in the primary analysis of data from a randomized trial. The method used should provide valid inference under approximately the same minimal statistical assumptions that would be needed for unadjusted estimation in a randomized trial." The assumption of a correctly specified model is not typically part of the assumptions needed for an unadjusted analysis, and should therefore be avoided when possible.

All estimands described in this paper can be computed from the cumulative distribution functions (CDF) $F(\cdot, a)$ for $a \in \{0, 1\}$, which can be estimated using the empirical cumulative distribution function (ECDF) or the Kaplan-Meier estimator. Model-robust, covariate adjusted estimators have been developed for the CDF, including, e.g., Chen and Tsiatis (2001); Rubin and van der Laan (2008); Moore and van der Laan (2009b); Stitelman et al. (2011); Lu and Tsiatis (2011); Brooks et al. (2013); Zhang (2014); Parast et al. (2014); Benkeser et al. (2018); Díaz (2019).

We focus on the model-robust, covariate adjusted estimators of Moore and van der Laan (2009b), Díaz et al. (2016), and Díaz et al. (2019). These estimators have at least two advantages compared to unadjusted estimators based on the ECDF or the Kaplan-Meier estimator. First, with time-to-event outcomes, the adjusted estimator can achieve consistency under an assumption of censoring being independent of the outcome given study arm and baseline covariates ($C \perp\!\!\!\perp Y \mid A, W$), rather than the assumption of censoring in each arm being independent of the outcome marginally ($C \perp\!\!\!\perp Y \mid A$) required by unadjusted estimators. The former assumption is arguably more likely to hold in typical situations where patients are lost to follow-up due to reasons correlated with their baseline variables. Second, in large samples and under regularity conditions, the adjusted estimators of Díaz et al. (2016) and Díaz et al. (2019) can be at least as precise as the unadjusted estimator (this requires that missingness/censoring is completely at random, i.e., that in each arm $a \in \{0, 1\}$, $C \perp\!\!\!\perp (Y, W)|A = a$), under additional assumptions.

Additionally, under regularity conditions, the three aforementioned adjusted estimators are asymptotically normal. This allows the construction of Wald-type confidence intervals and corresponding tests of the null hypothesis of no treatment effect.

## 3.1 Prediction algorithms

While we make no assumption on the functional form of the distribution $\mathsf{P}$ (except that treatment is independent of baseline variables by randomization), implementation of our estimators requires a *working model* for the following conditional probability

$$m(k, a, W) = \mathsf{P}(\widetilde{Y} = k, \Delta = 1 \mid \widetilde{Y} \geq k, A = a, W).$$

In time-to-event analysis, this probability is known as the conditional hazard. The expression *working model* here means that we do not assume that the model represents the true relationship between the outcome and the treatment/covariates. Fitting a working model for $m$ is equivalent to training a prediction model for $m$ (specifically, a prediction model for

the probability of $\widetilde{Y} = k, \Delta = 1$ given $\widetilde{Y} \geq k, A = a, W$), and we sometimes refer to the model fit as a predictor.

In our simulation studies, we will use the following working models, fitted in a dataset where each participant contributes a row of data corresponding to each time $k = 1$ through $k = \widetilde{Y}$:

- The following pooled main terms logistic regression (LR) $\text{logit}\{m_\beta(k, a, W)\} = \beta_{a,0,k} + \beta_{a,1}^\top W$ estimated with maximum likelihood estimation. Note that this model has (i) separate parameters for each study arm, and (ii) in each arm, intercepts for each possible outcome level $k$.

- The above model fitted with an $\ell_1$ penalty on the parameter $\beta_{a,1}$ ($\ell_1$-LR, Tibshirani 1996; Park and Hastie 2007).

- A random forest classification model (RF, Breiman 2001).

- An extreme gradient boosting tree ensemble (XGBoost, Friedman 2001).

- Multivariate adaptive regression splines (MARS, Friedman 1991).

For RF, XGBoost, and MARS, the algorithms are trained in the whole sample $\{1, \ldots, n\}$. For these algorithms, we also assessed the performance of cross-fitted versions of the estimators. Cross-fitting is sometimes necessary to guarantee that the regularity assumptions required for asymptotic normality of the estimators hold when using data-adaptive regression methods (Klaassen 1987; Zheng and van der Laan 2011; Chernozhukov et al. 2018), and is performed as follows. Let $\mathcal{V}_1, \ldots, \mathcal{V}_J$ denote a random partition of the index set $\{1, \ldots, n\}$ into $J$ prediction sets of approximately the same size. That is, $\mathcal{V}_j \subset \{1, \ldots, n\}$; $\bigcup_{j=1}^J \mathcal{V}_j = \{1, \ldots, n\}$; and $\mathcal{V}_j \cap \mathcal{V}_{j'} = \emptyset$. In addition, for each $j$, the associated training sample is given by $\mathcal{T}_j = \{1, \ldots, n\} \backslash \mathcal{V}_j$. Let $\widehat{m}_j$ denote the prediction algorithm trained in $\mathcal{T}_j$. Letting $j(i)$ denote the index of the prediction set which contains observation $i$, cross-fitting entails using only observations in $\mathcal{T}_{j(i)}$ for fitting models when making predictions about observation $i$. That is, the outcome predictions for each subject $i$ are given by $\widehat{m}_{j(i)}(u, a, W_i)$. We let $\widehat{\eta}_{j(i)} = (\widehat{m}_{j(i)}, \widehat{\pi}_A, \widehat{\pi}_C)$ for cross-fitted estimators and $\widehat{\eta}_{j(i)} = (\widehat{m}, \widehat{\pi}_A, \widehat{\pi}_C)$ for non-cross-fitted ones. RF, XGBoost, and MARS were fit using the *ranger* (Wright and Ziegler 2017), *xgboost* (Chen et al. 2021), and *earth* (Milborrow 2020) R packages, respectively. Hyperparameter tuning was performed using cross-validation with the *origami* (Coyle and Hejazi 2020) R package.

## 3.2 Targeted minimum loss based estimation (TMLE)

Our simulation studies use the TMLE procedure presented in Díaz et al. (2019). We will refer to that estimator as TMLE with improved efficiency, or IE-TMLE. We will first present the TMLE of (Moore and van der Laan 2009b), which constitutes the basis for the construction of the IE-TMLE.

In the supplementary materials we present some of the efficiency theory underlying the construction of the TMLE. Briefly, TMLE is a framework to construct estimators $\widehat{\eta}_{j(i)}$ that solve the efficient influence function estimating equation $n^{-1} \sum_{i=1}^n D_{\widehat{\eta}_{j(i)}}(O_i) = 0$, where

$D_\eta(O)$ is the efficient influence function for $S(k, a)$ in the non-parametric model that only assumes treatment $A$ is independent of baseline variables $W$ (which holds by design), defined in the supplementary materials. TMLE enjoys desirable properties such as local efficiency, outcome model robustness under censoring completely at random, and asymptotic normality, under regularity assumptions.

**TMLE estimator definition:** Given a predictor $\widehat{m}$ constructed as in the previous subsection and any $k, a$, the corresponding TMLE estimation procedure for $F(k, a)$ can be summarized in the next steps:

1. Create a long-form dataset where each participant $i$ contributes the following row of data corresponding to each time $u = 0$ through $k$:

$$\left(u, W_i, A_i, 1\{\widetilde{Y} \geq u\}, 1\{\widetilde{Y} = u, \Delta = 0\}, 1\{\widetilde{Y} = u, \Delta = 1\}\right),$$

where $1\{X\}$ is the indicator variable taking value 1 if $X$ is true and 0 otherwise.

2. For each individual $i$, obtain a prediction $\widehat{m}(u, a, W_i)$ for each pair in the set $\{(u, a) : a = 0, 1; u = 0, \ldots, k\}$.

3. Fit a model $\pi_A(a, W)$ for the probability $\mathsf{P}(A = a \mid W)$. Note that, in randomized trials, this model may be correctly specified by a logistic regression $\operatorname{logit} \pi_A(1, W) = \alpha_0 + \alpha_1^\top W$. Let $\widehat{\pi}_A(a, W_i)$ denote the prediction of the model for individual $i$.

4. Fit a model $\pi_C(u, a, W)$ for the censoring probabilities $\mathsf{P}(\widetilde{Y} = u, \Delta = 0 \mid \widetilde{Y} \geq u, A = a, W)$. For time-to-event outcomes, this is a model for the censoring probabilities. For ordinal outcomes, the only possibilities are that $C = 0$ (outcome missing) or $C = K$ (outcome observed); in this case we only fit the aforementioned model at $u = 0$ and we set $\pi_C(u, a, W) = 0$ for each $u > 0$. For either outcome type, if there is no censoring (i.e., if $P(\Delta = 1) = 1$), then we set $\pi_C(u, a, W) = 0$ for all $u$. Let $\widehat{\pi}_C(u, a, W_i)$ denote the prediction of this model for individual $i$, i.e., using the baseline variable values from individual $i$.

5. For each individual $i$ and each $u \leq k$, compute a "clever" covariate $H_{Y,k,u}$ as a function of $\widehat{m}$, $\widehat{\pi}_A$, and $\widehat{\pi}_C$ as detailed in the supplementary materials. The outcome model fit $\widehat{m}$ is then updated by fitting the following logistic regression "tilting" model with single parameter $\epsilon$ and offset based on $\widehat{m}$:

$$\mathsf{P}(\widetilde{Y} = u, \Delta = 1 \mid \widetilde{Y} \geq u, A = a, W) = \operatorname{logit}^{-1}\left\{\operatorname{logit} \widehat{m}(u, a, W) + \varepsilon H_{Y,k,u}\right\}.$$

This can be done using standard statistical software for fitting a logistic regression of the indicator variable $1\{\widetilde{Y} = u, \Delta = 1\}$ on the variable $H_{Y,k,u}$ using offset $\operatorname{logit} \widehat{m}(u, a, W)$ among observations with $\widetilde{Y} \geq u$ and $A = a$ in the long-form dataset from step 1. The above model fitting process is iterated where at the beginning of each iteration we replace $\widehat{m}$ in the above display and in the definition of $H_{Y,k,u}$ by the updated model fit. We denote the maximum number of iterations that we allow by $i_{\max}$.

6. Let $\widetilde{m}(u, a, W_i)$ denote the estimate of $m(u, a, W_i)$ for individual $i$ at the final iteration of the previous step. Note that this estimator is specific to the value $k$ under consideration.

7. Compute the estimate of $S(k, a) = 1 - F(k, a)$ as the following standardized estimator

$$\tilde{S}_{\text{TMLE}}(k, a) = \frac{1}{n} \sum_{i=1}^{n} \prod_{u=1}^{k} \{1 - \widetilde{m}(u, a, W_i)\}, \tag{1}$$

and let the estimator of $F(k, a)$ be $1 - \tilde{S}_{\text{TMLE}}(k, a)$.

This estimator was originally proposed by Moore and van der Laan (2009b). The role of the clever covariate $H_{Y,k,u}$ is to endow the resulting estimator $\widetilde{S}(k, a)$ with properties such as model-robustness compared to unadjusted estimators. In particular, it can be shown that this estimator is efficient when the working model for $m$ is correctly specified. The specific form of the covariate $H_{Y,k,u}$ is given in the supplementary materials. Throughout, the notation $\widehat{m}$ is used to represent the predictor constructed as in Section 3.1 and which is an input to the above TMLE algorithm, while $\widetilde{m}$ denotes the updated version of this predictor that is output by the above TMLE algorithm at step 6.

**IE-TMLE estimator definition:** In Section 4 we will compare several machine learning procedures for estimating $m$ in finite samples. The estimators used in the simulation study are the IE-TMLE of Díaz et al. (2019), where in addition to updating the initial estimator for the outcome regression $m$, we also update the estimators of the treatment and censoring mechanisms. Specifically, we replace step 5 of the above procedure with the following:

5. For each individual $i$ construct "clever" covariates $H_{Y,k,u}$, $H_A$, and $H_{C,k,u}$ (defined in the supplementary materials) as a function of $\widehat{m}$, $\widehat{\pi}_A$, and $\widehat{\pi}_C$. For each $k = 1, \ldots, K$, the model fits are then iteratively updated using logistic regression "tilting" models:

$$\text{logit } m_\varepsilon(u, a, W) = \text{logit } \widehat{m}(u, a, W) + \varepsilon H_{Y,k,u}$$
$$\text{logit } \pi_{\gamma,A}(1, W) = \text{logit } \widehat{\pi}_A(1, W) + \gamma H_A$$
$$\text{logit } \pi_{\upsilon,C}(u, a, W) = \text{logit } \widehat{\pi}_C(u, a, W) + \upsilon H_{C,k,u}$$

where the iteration is necessary because $H_{Y,k,u}$, $H_A$, and $H_{C,k,u}$ are functions of $\widehat{m}$, $\widehat{\pi}_A$, and $\widehat{\pi}_C$ that must be updated at each step. As before, for ordinal outcomes we only fit the aforementioned model at $u = 0$ and we set $\pi_C(u, a, W) = 0$ for each $u > 0$.

We use $\tilde{S}_{\text{IE}-\text{TMLE}}$ to denote this estimator. The updating step above combines ideas from Moore and van der Laan (2009b), Gruber and van der Laan (2012), and Rotnitzky et al. (2012) to produce an estimator with the following properties:

(i) Consistency and at least as precise as the Kaplan-Meier and inverse probability weighted estimators;

(ii) Consistency under violations of independent censoring (unlike the Kaplan-Meier estimator) when either the censoring or survival distributions, conditional on covariates, are estimated consistently and censoring is such that $C \perp\!\!\!\perp Y \mid W, A$; and

9

(iii) Nonparametric efficiency when both of these distributions are consistently estimated at rate $n^{1/4}$.

Please see Díaz et al. (2019) for more details on these estimators, which are implemented in the R package `adjrct` (Díaz and Williams 2021).

Next, we present a result (Theorem 1) stating asymptotic normality of $\tilde{S}_{\text{TMLE}}$ using machine learning for prediction that avoids some limitations of existing methods, and present a consistent estimator of its variance. In Section 4 we present simulation results where we evaluate the performance of $\tilde{S}_{\text{IE-TMLE}}$ for covariate adjustment in COVID-19 trials for hospitalized patients. We favor $\tilde{S}_{\text{IE-TMLE}}$ in our numerical studies because, unlike $\tilde{S}_{\text{TMLE}}$, it satisfies property (i) above. The simulation uses Wald-type hypothesis tests based on the asymptotic approximation of Theorem 1, where we note that the variance estimator in the theorem is consistent for $\tilde{S}_{\text{TMLE}}$ but it is conservative for $\tilde{S}_{\text{IE-TMLE}}$ (Moore and van der Laan 2009b).

## 3.3 Asymptotically correct confidence intervals and hypothesis tests for TMLE combined with machine learning

Most available methods to construct confidence intervals and hypothesis tests in the statistics literature are based on the sampling distribution of the estimator. While using the exact finite-sample distribution would be ideal for this task, such distributions are notoriously difficult to derive for our problem in the absence of strong and unrealistic assumptions (such as linear models with Gaussian noise). Thus, here we focus on methods that rely on approximating the finite-sample distribution using asymptotic results as $n$ goes to infinity.

In order to discuss existing methods, it will be useful to introduce and compare the following assumptions:

**A1.** Censoring is completely at random, i.e., $C \perp\!\!\!\perp (Y, W) \mid A = a$ for each treatment arm $a$.

**A2.** Let $||f||^2$ denote the $L_2(\mathsf{P})$ norm $\int f^2(o)\mathrm{d}\mathsf{P}(o)$, for $O = (W, A, \Delta = \mathbb{1}\{Y \leq C\}, \widetilde{Y})$. We abbreviate $m(k, a, W)$ and $\widehat{m}(k, a, W)$ by $m$ and $\widehat{m}$, respectively. Assume the estimator $\widehat{m}$ is consistent in the sense that $||\widehat{m} - m|| = o_P(1)$ for all $k \in \{1, \ldots, K\}$ and $a \in \{0, 1\}$. We also assume that there exists a $\delta > 0$ such that $\delta < m < 1 - \delta$ with probability 1.

**A3.** Assume the estimator $\widehat{m}$ converges to a possibly misspecified limit $m_1$ in the sense that $||\widehat{m} - m_1|| = o_P(1)$ for all $k \in \{1, \ldots, K\}$ and $a \in \{0, 1\}$, where we emphasize that $m_1$ can be different from the true regression function $m$. We also assume that there exists a $\delta > 0$ such that $\delta < m_1 < 1 - \delta$ with probability 1.

For estimators $\widehat{m}$ of $m$ that use cross-fitting, the function $\widehat{m}$ consists of $J$ maps (one for each training set) from the sample space of $O$ to the interval $[0, 1]$. In this case, by convention we define $||\widehat{m} - m||$ in A2 as the average across the $J$ maps of the $L_2(\mathsf{P})$ norm of each such map minus $m$. Convergence of $||\widehat{m} - m||$ to 0 in probability is then equivalent to the same convergence where $\widehat{m}$ is replaced by the corresponding map before cross-fitting is applied. The same convention is used in A3.

There are at least two results on asymptotic normality for $\tilde{S}_{\text{TMLE}}$ relevant to the problem we are studying. The first result is a general theorem for TMLE (see Appendix A.1 of

[van der Laan and Rose 2011](#)), stating that the estimator is asymptotically normal and efficient under regularity assumptions which include A2. Among other important implications, this asymptotic normality implies that the variance of the estimators can be consistently estimated by the empirical variance of the efficient influence function. This means that asymptotically correct confidence intervals and hypothesis tests can be constructed using a Wald-type procedure. As stated above, it is often undesirable to assume A2 in the setting of a randomized trial, as it is a much stronger assumption than what would be required for an unadjusted estimator.

The second result of relevance to this paper establishes asymptotic normality of $\widetilde{S}(k, a)$ under assumptions that include A3 ([Moore and van der Laan 2009a](#)). The asymptotic variance derived by these authors depends on the true outcome regression function $m$, and is thus difficult to estimate. As a solution, the authors propose to use a conservative estimate of the variance whose computation does not rely on the true regression function $m$. While this conservative method yields correct type 1 error control, its use is not guaranteed to fully covert precision gains from covariate adjustment into power gains.

We note that the above asymptotic normality results from related works rely on the additional condition that the estimator $\widehat{m}$ lies in a Donsker class. This assumption may be violated by some of the data-adaptive regression techniques that we consider. Furthermore, we note that resampling methods such as the bootstrap cannot be safely used for variance estimation in this setting. Their correctness is currently unknown when the working model for $m$ is based on data-adaptive regression procedures such as those described in Section 3.1 and used in our simulation studies.

In what follows, we build on recent literature on estimation of causal effects using machine learning to improve upon the aforementioned asymptotic normality results on two fronts. First, we introduce cross-fitting ([Klaassen 1987](#); [Zheng and van der Laan 2011](#); [Chernozhukov et al. 2018](#)) to avoid the Donsker condition. Second, and most importantly, we present a novel asymptotic normality result that avoids the above limitations of existing methods regarding strong assumptions (specifically A2) and conservative variance estimators (that may sacrifice power).

The following are a set of assumptions about how components of the TMLE are implemented, which we'll use in our theorem below:

**A4.** The initial estimator of $\pi_A(1)$ is set to be the empirical mean $n^{-1} \sum_{i=1}^n A_i$.

**A5.** For time-to-event outcomes, the initial estimator $\widehat{\Pi}_C(a, u)$ is set to be the Kaplan-Meier estimator estimated separately within each treatment arm $a$. For ordinal outcomes, $\widehat{\Pi}_C(a, 0)$ is the proportion of missing outcomes in treatment arm $a$ and $\widehat{\Pi}_C(a, u) = 0$ for $u > 0$.

**A6.** The initial estimator $\widehat{m}(u, a, W)$ is constructed using one of the following:

1. Any estimator in a parametric working model (i.e., a model that can be indexed by a Euclidean parameter) such as maximum likelihood, $\ell_1$ regularization, etc.

2. Any data-adaptive regression method (e.g., random forests, MARS, XGBoost, etc.) estimated using cross-fitting as described above.

**A7.** The regularity conditions in Theorem 5.7 of (van der Vaart 1998, p.45) hold for the maximum likelihood estimator corresponding to each logistic regression model fit in step (5) of the TMLE algorithm.

**Theorem 1.** *Assume A1 and A3–A7 above. Define the variance estimator*

$$\widetilde{\sigma}^2 = \frac{1}{n}\sum_{i=1}^{n}[D_{\widetilde{\eta}_{j(i)}}(O_i)]^2.$$

*Then we have for all $k \in \{1,\ldots,K\}$ and $a \in \{0,1\}$ that*

$$\sqrt{n}\{\tilde{S}_{\mathrm{TMLE}}(k,a) - S(k,a)\}/\widetilde{\sigma} \rightsquigarrow N(0,1).$$

Theorem 1 is a novel result establishing the asymptotic correctness of Wald-type confidence intervals and hypothesis tests for the covariate-adjusted estimator $\tilde{S}_{\mathrm{TMLE}}(k,a)$ based on machine learning regression procedures constructed as stated in A6. For example, the confidence interval $\tilde{S}_{\mathrm{TMLE}}(k,a) \pm 1.96 \times \widetilde{\sigma}/\sqrt{n}$ has approximately 95% coverage at large sample sizes, under the assumptions of the theorem. The theorem licenses the large sample use of any regression procedure for $m$ when combined with the TMLE of Section 3.2, as long as the regression procedure is either (i) based on a parametric model (such as $\ell_1$-regularization) or (ii) based on cross-fitted data-adaptive regression, and the assumptions of the theorem hold. The theorem states sufficient assumptions under which Wald-type tests from such a procedure will be asymptotically correct.

Assumption A3 states that the predictions given by the regression method used to construct the adjusted estimator converge to some arbitrary function (i.e., not assumed to be equal to the true regression function). This assumption is akin to Condition 3 assumed by Bloniarz et al. (2016) in the context of establishing asymptotic normality of a covariate-adjusted estimator based on $\ell_1$-regularization. We note that this is an assumption on the predictions themselves and not on the functional form of the predictors. Therefore, issues like collinearity do not necessarily cause problems. While this assumption can hold for many off-the-shelf machine learning regression methods under assumptions on the data-generating mechanism, general conditions have not been established and the assumption must be checked on a case-by-case basis.

We note that assumption A1 is stronger than the assumption $C \perp\!\!\!\perp Y \mid A = a$ required by unadjusted estimators such as the Kaplan-Meier estimator. However, if $W$ is prognostic (meaning that $W \not\perp\!\!\!\perp Y \mid A = a$), then the assumption $C \perp\!\!\!\perp Y \mid A = a$ required by the Kaplan-Meier estimator cannot generally be guaranteed to hold, unless A1 also holds. Thus, our theorem aligns with the recent FDA draft guidance on covariate adjustment in the sense that "it provides valid inference under approximately the same minimal statistical assumptions that would be needed for unadjusted estimation in a randomized trial" (U.S. Food and Drug Administration 2021).

The construction of estimators based on A5 should be avoided if A1 does not hold. Confidence that A1 holds is typically warranted in trials where the only form of right censoring is administrative. When applied to ordinal outcomes, A1 is trivially satisfied if there is no missing outcome data.

Consider the case where censoring is informative such that A1 does not hold, but censoring at random holds (i.e., $C \perp\!\!\!\perp Y \mid W, A$). Then consistency of the estimators $\tilde{S}_{\text{TMLE}}$ and $\tilde{S}_{\text{IE-TMLE}}$ will typically require that *at least one* of two assumptions hold: (a) that the censoring probabilities $\pi_C(u, a, w)$ are estimated consistently, or that (b) the outcome regression $m(u, a, w)$ is estimated consistently. To maximize the chances of either of these conditions being true, we recommend the use of flexible machine learning for both of these regressions, including model selection and ensembling techniques such as the Super Learner (van der Laan et al. 2007). The conditions for asymptotic normality of $\tilde{S}_{\text{TMLE}}$ and $\tilde{S}_{\text{IE-TMLE}}$ under these circumstances are much stronger than those for Theorem 1, and typically include consistent estimation of *both* $\pi_C(u, a, w)$ and $m(u, a, w)$ at certain rates (e.g., each of them converging at $n^{1/4}$-rate is sufficient, see Appendix A.1 of van der Laan and Rose 2011).

# 4 Simulation methods

Our data generating distribution is based on a database of over 1,500 patients hospitalized at Weill Cornell Medicine New York Presbyterian Hospital prior to 15 May 2020. The database includes information on patients 18 years of age and older with COVID-19 confirmed through reverse-transcriptase–polymerase-chain-reaction assays. For a full description of the clinical characteristics and data collection methods of the initial cohort sampling, see Goyal et al. (2020).

We evaluate the potential to improve efficiency by adjustment for subsets of the following baseline variables: age, sex, BMI, smoking status, whether the patient required supplemental oxygen within three-hours of presenting to the emergency department, number of comorbidities (diabetes, hypertension, COPD, CKD, ESRD, asthma, interstitial lung disease, obstructive sleep apnea, any rheumatological disease, any pulmonary disease, hepatitis or HIV, renal disease, stroke, cirrhosis, coronary artery disease, active cancer), number of relevant symptoms, presence of bilateral infiltrates on chest x-ray, dyspnea, and hypertension. These variables were chosen because they have been previously identified as risk factors for severe disease (Guan et al. 2020; Goyal et al. 2020; Gupta et al. 2020), and therefore are likely to improve efficiency of covariate-adjusted effect estimators in randomized trials in hospitalized patients.

Code to reproduce our simulations may be found at https://github.com/nt-williams/covid-RCT-co

## 4.1 Data generating mechanisms

We consider two types of outcomes: a time-to-event outcome defined as the time from hospitalization to intubation or death, and a six-level ordinal outcome at 14 days post-hospitalization based on the WHO Ordinal Scale for Clinical Improvement (Marshall et al. 2020). The categories are as follows: 0, discharged from hospital; 1, hospitalized with no oxygen therapy; 2, hospitalized with oxygen by mask or nasal prong; 3, hospitalized with non-invasive ventilation or high-flow oxygen; 4, hospitalized with intubation and mechanical ventilation; 5, dead. For time to event outcomes, we focus on evaluating the effect of treatment on the RMST at 14 days and the RD at 7 days after hospitalization, and for ordinal outcomes we evaluate results for both the LOR and the Mann-Whitney statistic.

We simulate datasets for four scenarios where we consider two effect sizes (null versus positive) and two baseline variable settings (prognostic versus not prognostic, where prognostic means marginally associated with the outcome). For each sample size $n \in \{100, 500, 1500\}$ and for each scenario, we simulated 5000 datasets as follows. To generate datasets where covariates are prognostic, we draw $n$ pairs $(W, Y)$ randomly from the original dataset with replacement. This generates a dataset where the covariate prognostic strength is as observed in the real dataset. To simulate datasets where covariates are not prognostic, we first draw outcomes $Y$ at random with replacement from the original dataset, and then draw covariates $W$ at random with replacement and independent of the value $Y$ drawn.

For each scenario, a hypothetical treatment variable is assigned randomly for each patient with probability 0.5 independent of all other variables. This produces a data generating distribution with zero treatment effect. Next, a positive treatment effect is simulated for time-to-event outcomes by adding an independent random draw from a $\chi^2$ distribution four degrees of freedom to each patient's observed survival time in the treatment arm. This effect size translates to a difference in RMST of 1.04 and RD of 0.10, respectively. To simulate outcomes being missing completely at random, 5% of the patients are selected at random to be censored, and the censoring times are drawn from a uniform distribution between 1 and 14. A positive treatment effect is simulated for ordinal outcomes by subtracting from each patient's outcome in the treatment arm an independent random draw from a four-parameter Beta distribution with support $(0, 5)$ and parameters $(3, 15)$, rounded to the nearest nonnegative integer. This generates effect sizes for LOR of 0.60 and for MW of 0.46.

# 5   Simulation results

We evaluate several estimators. First, we evaluate unadjusted estimators based on substituting the empirical CDF for ordinal outcomes and the Kaplan-Meier estimator for time-to-event outcomes in the parameter definitions of Section 2. We then evaluate adjusted estimator $\tilde{S}_{\text{IE-TMLE}}(k, a)$ where the working models are:

LR: a fully adjusted estimator using logistic regression including all the variables listed in the previous section,
$\ell_1$-LR: $\ell_1$ regularization of the previous logistic regression,
RF: random forests,
MARS: multivariate adaptive regression splines, and
XGBoost: extreme gradient boosting tree ensembles.

For estimators RF, MARS, and XGBoost, we further evaluated cross-fitted versions of the working model. For all adjusted estimators the propensity score $\pi_A$ is estimated with an intercept-only model (A4), and the censoring mechanism $\pi_C$ is estimated using a Kaplan-Meier estimator fitted independently for each treatment arm (A5) (or equivalently for ordinal outcomes the proportion of missing outcomes within each treatment arm).

Confidence intervals and hypothesis tests are performed using Wald-type statistics, which use an estimate of the standard error. The standard error was estimated based on the asymptotic Gaussian approximation described in Theorem 1. We compare the performance

of the estimators in terms of the probability of type-1 error, power, the absolute bias, the variance, and the mean squared error.

We compute the relative efficiency RE of each estimator compared to the unadjusted estimator as a ratio of the mean squared errors. This relative efficiency can be interpreted as the ratio of sample sizes required by the estimators to achieve the same power at local alternatives, asymptotically (van der Vaart 1998). Equivalently, one minus the relative efficiency is the relative reduction (due to covariate adjustment) in the required sample size to achieve a desired power, asymptotically; e.g., a relative efficiency of 0.8 is approximately equivalent to needing 20% smaller sample size when using covariate adjustment.

In the presentation of the results, we append the prefix CF to cross-fitted estimators. For example, CF-RF will denote cross-fitted random forests.

Tables containing the comprehensive results of the simulations are presented in the supplementary materials. In the remainder of this section we present a summary of the results. First, we note that the use of random forests without cross-fitting exhibits very poor performance, failing to appropriately control type-1 error when the effect is null, and introducing significant bias when the effect is positive. We observed this poor performance across all simulations. Thus, in what follows we omit a discussion of this estimator.

Results for the LOR in Tables 3 and 11 show that covariate adjusted estimators have better performance than the unadjusted estimator at small sample sizes, even when the covariates are not prognostic. In these cases, the unadjusted estimator is unstable with large variance due to near-empty outcome categories in some simulated datasets, which causes division by near-zero numbers in the unadjusted LOR estimator. Some covariate adjusted estimators fix this problem by extrapolating model probabilities to obtain better estimates of the probabilities in the near-empty cells.

Tables 1-4 (in the web supplementary materials) display the results for the difference in RMST, RD, LOR, and MW estimands when covariates are prognostic and there is a positive effect size. At sample size $n = 1500$ all adjusted estimators yield efficiency gains, with CF-RF offering the best RE ranging from 0.51 to 0.67 compared to an unadjusted estimator, while appropriately controlling type-1 error. In contrast, the RE of $\ell_1$-LR at $n = 1500$ ranged from 0.79 to 0.89.

At sample size $n = 500$, $\ell_1$-LR, CF-RF, and XGBoost offer comparable efficiency gains, ranging from 0.29 to 0.99. As the sample size decreases to $n = 100$ most adjusted estimators yield efficiency losses and the only estimator that retains efficiency gains is $\ell_1$-LR, with RE from 0.86 to 0.92. (An exception is in estimation of the LOR, where the RE of $\ell_1$-LR was 0.1 due to the issue discussed above.)

Efficiency gains for $\ell_1$-LR did not always translate into power gains of a Wald-type hypothesis test compared to other estimators (e.g. LR at $n = 100$), possibly due to biased variance estimation and/or a poor Gaussian approximation of the distribution of the test statistic. At small sample size $n = 100$ power was uniformly better for a Wald-type test based on LR compared to $\ell_1$-LR. At sample size $n = 500$ a Wald-type test based on $\ell_1$-LR seemed to dominate all other algorithms, whereas at $n = 1500$ all algorithms had comparable power very close to one.

Results when the true treatment effect is zero and covariates are prognostic are presented in Tables 5-8 (in the web supplementary materials). At sample size $n = 1500$, CF-RF generally provides large efficiency gains with relative efficiencies ranging from 0.66 to 0.77.

15

For comparison, $\ell_1$-LR has RE ranging from 0.88 to 0.92. As the sample size decreases to $n = 500$, $\ell_1$-LR and CF-RF both offer the most efficiency gains while retaining type-1 error control, with RE ranging from 0.74 to 0.88. At small sample sizes $n = 100$, $\ell_1$-LR consistently leverages efficiency gains from covariate adjustment (RE ranging from 0.73 to 0.95) but its type-1 error (ranging from 0.07 to 0.09) is slightly larger than that of the unadjusted estimator. For estimation of LOR and MW, XGBoost has similar results at sample size $n = 100$.

Tables 9-12 (in the web supplementary materials) show results for scenarios where the covariates are not prognostic of the outcome but there is a positive effect. This case is interesting because it is well known that adjusted estimators can induce efficiency losses (i.e., RE > 1) by adding randomness to the estimator when there is nothing to be gained from covariate adjustment. We found that $\ell_1$-LR uniformly avoids efficiency losses associated with adjustment for independent covariates, with a maximum RE of 1.03. All other covariate adjustment methods had larger maximum RE. At sample size $n = 100$, the superior efficiency of the $\ell_1$-LR estimator did not always translate into better power (e.g., compared to LR) due to the use of a Wald-test which relies on an asymptotic approximation to the distribution of the estimator.

Results when the true treatment effect is zero and covariates are not prognostic are presented in Tables 13-16 (in the web supplementary materials). In this case, $\ell_1$-LR also avoids efficiency losses across all scenarios, while maintaining a type-1 error that is comparable to that of the unadjusted estimator.

Lastly, at large sample sizes all cross-fitted estimators along with logistic regression estimators yield correct type I error, illustrating the correctness of Wald-type tests proved in Theorem 1. Our simulation results also show that Wald-type hypothesis tests based on data-adaptive machine learning procedures fail to control type 1 error if the regressions procedures are not cross-fitted.

# 6  Recommendations and future directions

In our numerical studies we found that $\ell_1$-regularized logistic regression offers the best trade-off between type-I error control and efficiency gains across sample sizes, outcome types, and estimands. We found that this algorithm leverages efficiency gains when efficiency gains are feasible, while protecting the estimators from efficiency losses when efficiency gains are not feasible (e.g., adjusting for covariates with no prognostic power). A direction of future research is the evaluation of bootstrap estimators for the variance and confidence intervals of covariate-adjusted estimators, especially for cases where the Wald-type methods evaluated in this manuscript did not perform well (e.g., $\ell_1$-LR at $n = 100$).

We also found that logistic regression can result in large efficiency losses for small sample sizes, with relative efficiencies as large as 1.17 for the RMST estimand, and as large as 7.57 for the MW estimand. Covariate adjustment with $\ell_1$-regularized logistic regression solves this problem, maintaining efficiency when covariates are not prognostic for the outcome, even at small sample sizes. However, Wald-type hypothesis tests do not appropriately translate the efficiency gains of $\ell_1$-regularized logistic regression into more powerful tests. This requires the development of tests appropriate for small samples.

We recommend against using the LOR parameter since it is difficult to interpret and the corresponding estimators (even unadjusted ones) can be unstable at small sample sizes. Covariate adjustment with $\ell_1$-LR, CF-MARS, CF-RF, or CF-XGBoost can aid to improve efficiency in estimation of the LOR parameter over the unadjusted estimator when there are near-empty cells at small sample sizes. This improvement in efficiency did not translate into an improvement in power when using Wald-type hypothesis tests, due to poor small-sample Gaussian approximations or poor variance estimators.

We discourage the use of non-cross-fitted versions of the machine learning methods evaluated (i.e., RF, XGBoost, MARS) for covariate adjustment. Specifically, we found in simulations that non-cross-fitted random forests can lead to overly biased estimators in the case of a positive effect, and to anti-conservative Wald-type hypothesis tests in the case of a null treatment effect. We found that cross-fitting the random forests alleviated this problem and was able to produce small bias and acceptable type-1 error at all sample sizes. This is supported at large sample sizes by our main theoretical result (Theorem 1) which establishes asymptotic correctness of cross-fitted procedures under regularity conditions. In fact, we found that random forests with cross-fitting provided the most efficiency gains at large sample sizes.

Based on the results of our simulation studies, we recommend that cross-fitting with data-adaptive estimators such as random forests and extreme gradient boosting be considered for covariate selection in trials with large sample sizes ($n = 1500$ in our simulations). In large sample sizes, it is also possible to consider an ensemble approach such as Super Learning (van der Laan et al. 2007) that allows one to select the predictor that yields the most efficiency gains. Traditional model selection with statistical learning is focused on the goal of prediction, and an adaptation of those tools to the goal of maximizing efficiency in estimating the marginal treatment effect is the subject of future research.

The conditions for asymptotic normality and consistent variance estimation of $\tilde{S}_{\text{TMLE}}(k, a)$ established in Theorem 1 may be restrictive if censoring is informative. In that case, consistency of the $\tilde{S}_{\text{TMLE}}(k, a)$ and $\tilde{S}_{\text{IE-TMLE}}(k, a)$ estimators requires that censoring at random holds (i.e., $C \perp\!\!\!\perp Y \mid W, A$), and that either the outcome regression or censoring mechanism is consistently estimated. Thus, it is recommended to also estimate the censoring mechanism with machine learning methods that allow for flexible regression. Standard asymptotic normality results for the $\tilde{S}_{\text{TMLE}}(k, a)$ and $\tilde{S}_{\text{IE-TMLE}}(k, a)$ require consistent estimation of both the censoring mechanism and the outcome mechanism at certain rates (e.g., both estimated at a $n^{1/4}$ rate is sufficient). The development of estimators that remain asymptotically normal under the weaker condition that at least one of these regressions is consistently estimated has been the subject of recent research (e.g., Díaz and van der Laan 2017; Benkeser et al. 2017; Díaz 2019).

# References

Ibrahim A. Ahmad. A class of Mann—Whitney—Wilcoxon type statistics. *The American Statistician*, 50(4):324–327, 1996.

Peter C. Austin, Andrea Manca, Merrick Zwarenstein, David N. Juurlink, and Matthew B.

Stanbrook. A substantial and confusing variation exists in handling of baseline covariates in randomized controlled trials: a review of trials published in leading medical journals. *Journal of Clinical Epidemiology*, 63(2):142–153, 2010.

Lindsey R. Baden, Hana M. El Sahly, Brandon Essink, Karen Kotloff, Sharon Frey, Rick Novak, David Diemert, Stephen A. Spector, Nadine Rouphael, C. Buddy Creech, et al. Efficacy and safety of the mRNA-1273 SARS-CoV-2 vaccine. *New England Journal of Medicine*, 384(5):403–416, 2021.

David Benkeser, Marco Carone, M. J. Van Der Laan, and P. B. Gilbert. Doubly robust nonparametric inference on the average treatment effect. *Biometrika*, 104(4):863–880, 2017.

David Benkeser, Marco Carone, and Peter B. Gilbert. Improved estimation of the cumulative incidence of rare outcomes. *Statistics in Medicine*, 37(2):280–293, 2018.

David Benkeser, Iván Díaz, Alex Luedtke, Jodi Segal, Daniel Scharfstein, and Michael Rosenblum. Improving precision and power in randomized trials for COVID-19 treatments using covariate adjustment, for binary, ordinal, and time-to-event outcomes. *Biometrics*, 2020.

Adam Bloniarz, Hanzhong Liu, Cun-Hui Zhang, Jasjeet S. Sekhon, and Bin Yu. Lasso adjustments of treatment effect estimates in randomized experiments. *Proceedings of the National Academy of Sciences*, 113(27):7383–7390, 2016.

Leo Breiman. Random forests. *Machine Learning*, 45(1):5–32, 2001.

Jordan C. Brooks, Mark J. van der Laan, Daniel E. Singer, and Alan S. Go. Targeted minimum loss-based estimation of causal effects in right-censored survival data with time-dependent covariates: Warfarin, stroke, and death in atrial fibrillation. *Journal of Causal Inference*, 1(2):235–254, 2013. doi: 10.1515/jci-2013-0001.

Pei-Yun Chen and Anastasios A. Tsiatis. Causal inference on the difference of the restricted mean lifetime between two groups. *Biometrics*, 57(4):1030–1038, 2001.

Tianqi Chen, Tong He, Michael Benesty, Vadim Khotilovich, Yuan Tang, Hyunsu Cho, Kailong Chen, Rory Mitchell, Ignacio Cano, Tianyi Zhou, Mu Li, Junyuan Xie, Min Lin, Yifeng Geng, and Yutian Li. *xgboost: Extreme Gradient Boosting*, 2021. URL https://CRAN.R-project.org/package=xgboost. R package version 1.4.1.1.

Victor Chernozhukov, Denis Chetverikov, Mert Demirer, Esther Duflo, Christian Hansen, Whitney Newey, and James Robins. Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal*, 21(1):C1–C68, 2018. doi: 10.1111/ectj. 12097.

Francis S. Collins and Paul Stoffels. Accelerating COVID-19 therapeutic interventions and vaccines (activ): an unprecedented partnership for unprecedented times. *JAMA*, 323(24): 2455–2457, 2020.

Jeremy Coyle and Nima Hejazi. *origami: Generalized Framework for Cross-Validation*, 2020. URL https://CRAN.R-project.org/package=origami. R package version 1.0.3.

Iván Díaz. Statistical inference for data-adaptive doubly robust estimators with survival outcomes. *Statistics in Medicine*, 38(15):2735–2748, 2019.

Iván Díaz and Mark J. van der Laan. Doubly robust inference for targeted minimum loss– based estimation in randomized trials with missing outcome data. *Statistics in Medicine*, 36(24):3807–3819, 2017.

Iván Díaz, Elizabeth Colantuoni, and Michael Rosenblum. Enhanced precision in the analysis of randomized trials with ordinal outcomes. *Biometrics*, 72(2):422–431, 2016.

Iván Díaz and Nicholas Williams. *adjrct: Efficient Estimators for Survival and Ordinal Outcomes in RCTs Without Proportional Hazards and Odds Assumptions*, 2021. URL https://github.com/nt-williams/adjrct. R package version 0.1.0.

Iván Díaz, Elizabeth Colantuoni, Daniel F. Hanley, and Michael Rosenblum. Improved precision in the analysis of randomized trials with survival outcomes, without assuming proportional hazards. *Lifetime Data Analysis*, 25(3):439–468, 2019.

FDA and EMA. E9 statistical principles for clinical trials. *U.S. Food and Drug Administration: CDER/CBER. European Medicines Agency: CPMP/ICH/363/96*, 1998.

Jerome H. Friedman. Multivariate adaptive regression splines. *The Annals of Statistics*, 19 (1):1–67, 1991.

Jerome H. Friedman. Greedy function approximation: a gradient boosting machine. *The Annals of Statistics*, 29(5):1189–1232, 2001.

Parag Goyal, Justin J. Choi, Laura C. Pinheiro, Edward J. Schenck, Ruijun Chen, Assem Jabri, Michael J. Satlin, Thomas R. Campion, Musarrat Nahid, Joanna B. Ringel, Katherine L. Hoffman, Mark N. Alshak, Han A. Li, Graham T. Wehmeyer, Mangala Rajan, Evgeniya Reshetnyak, Nathaniel Hupert, Evelyn M. Horn, Fernando J. Martinez, Roy M. Gulick, and Monika M. Safford. Clinical characteristics of Covid-19 in new york city. *New England Journal of Medicine*, 382(24):2372–2374, 2020. doi: 10.1056/NEJMc2010419.

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.

Wei-jie Guan, Zheng-yi Ni, Yu Hu, Wen-hua Liang, Chun-quan Ou, Jian-xing He, Lei Liu, Hong Shan, Chun-liang Lei, David S.C. Hui, Bin Du, Lan-juan Li, Guang Zeng, Kwok-Yung Yuen, Ru-chong Chen, Chun-li Tang, Tao Wang, Ping-yan Chen, Jie Xiang, Shi-yue Li, Jin-lin Wang, Zi-jing Liang, Yi-xiang Peng, Li Wei, Yong Liu, Ya-hua Hu, Peng Peng, Jian-ming Wang, Ji-yang Liu, Zhong Chen, Gang Li, Zhi-jian Zheng, Shao-qin Qiu, Jie Luo, Chang-jiang Ye, Shao-yong Zhu, and Nan-shan Zhong. Clinical characteristics of Coronavirus Disease 2019 in China. *New England Journal of Medicine*, 382(18):1708–1720, 2020.

Rishi K. Gupta, Michael Marks, Thomas H.A. Samuels, Akish Luintel, Tommy Rampling, Humayra Chowdhury, Matteo Quartagno, Arjun Nair, Marc Lipman, Ibrahim Abubakar, Maarten van Smeden, Wai Keong Wong, Bryan Williams, and Mahdad Noursadeghi. Systematic evaluation and external validation of 22 prognostic models among hospitalised adults with COVID-19: an observational cohort study. *European Respiratory Journal*, 56 (6), 2020. doi: 10.1183/13993003.03498-2020.

Brennan C. Kahan, Vipul Jairath, Caroline J. Doré, and Tim P. Morris. The risks and rewards of covariate adjustment in randomized trials: an assessment of 12 outcomes from 8 studies. *Trials*, 15(1):139, 2014.

Chris A. J. Klaassen. Consistent estimation of the influence function of locally asymptotically linear estimators. *The Annals of Statistics*, 15(4):1548–1562, 1987.

Kai Kupferschmidt and Jon Cohen. Race to find COVID-19 treatments accelerates. *Science*, 367(6485):1412–1413, 2020.

Xiaomin Lu and Anastasios A. Tsiatis. Semiparametric estimation of treatment effect with time-lagged response in the presence of informative censoring. *Lifetime Data Analysis*, 17 (4):566–593, 2011.

J. C. Marshall, S. Murthy, J Diaz, N. K. Adhikari, D. C. Angus, Y. M. Arabi, et al. A minimal common outcome measure set for COVID-19 clinical research. *The Lancet Infectious Diseases*, 20:e192–e197, 2020.

Stephen Milborrow. *earth: Multivariate Adaptive Regression Splines*, 2020. URL https://CRAN.R-project.org/package=earth. R package version 5.3.0.

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, 2009a.

Kelly L. Moore and Mark J. van der Laan. Increasing power in randomized trials with right censored outcomes through covariate adjustment. *Journal of Biopharmaceutical Statistics*, 19(6):1099–1131, 2009b.

Layla Parast, Lu Tian, and Tianxi Cai. Landmark estimation of survival and treatment effect in a randomized clinical trial. *Journal of the American Statistical Association*, 109 (505):384–394, 2014.

Mee Young Park and Trevor Hastie. L1-regularization path algorithm for generalized linear models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 69(4): 659–677, 2007.

Fernando P. Polack, Stephen J. Thomas, Nicholas Kitchin, Judith Absalon, Alejandra Gurtman, Stephen Lockhart, John L. Perez, Gonzalo Pérez Marc, Edson D. Moreira, Cristiano Zerbini, et al. Safety and efficacy of the bnt162b2 mrna covid-19 vaccine. *New England Journal of Medicine*, 383(27):2603–2615, 2020.

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.

Patrick Royston and Mahesh K. B. Parmar. The use of restricted mean survival time to estimate the treatment effect in randomized clinical trials when the proportional hazards assumption is in doubt. *Statistics in Medicine*, 30(19):2409–2421, 2011.

Daniel B. Rubin and Mark J. van der Laan. Empirical efficiency maximization: Improved locally efficient covariate adjustment in randomized experiments and survival analysis. *The International Journal of Biostatistics*, 4(1), 2008.

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.

The RECOVERY Collaborative Group. Dexamethasone in hospitalized patients with Covid-19. *New England Journal of Medicine*, 384(8):693–704, 2021.

Lu Tian, Tianxi Cai, Lihui Zhao, and Lee-Jen Wei. On the covariate-adjusted estimation for an overall treatment difference with data from a randomized comparative clinical trial. *Biostatistics*, 13(2):256–273, 2012.

Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996. doi: 10.1111/j.2517-6161.1996.tb02080.x.

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.

U.S. Food and Drug Administration. Adjusting for covariates in randomized clinical trials for drugs and biological products: Guidance for industry. *U.S. Food and Drug Administration: CDER/CBER.*, 2021. URL https://www.fda.gov/regulatory-information/search-fda-guidance-documents/adjusting-cov

Mark J. van der Laan and Sherri Rose. *Targeted learning: causal inference for observational and experimental data*. Springer Science & Business Media, 2011.

Mark J. van der Laan, Eric C. Polley, and Alan E. Hubbard. Super learner. *Statistical Applications in Genetics and Molecular Biology*, 6(1), 2007.

A. W. van der Vaart. *Asymptotic Statistics*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 1998. doi: 10.1017/CBO9780511802256.

A. W. van der Vaart. *Asymptotic Statistics*. Cambridge University Press, 1998.

Stefan Wager, Wenfei Du, Jonathan Taylor, and Robert J. Tibshirani. High-dimensional regression adjustments in randomized experiments. *Proceedings of the National Academy of Sciences*, 113(45):12673–12678, 2016.

WHO Solidarity Trial Consortium. Repurposed antiviral drugs for Covid-19 — interim WHO SOLIDARITY trial results. *New England Journal of Medicine*, 384(6):497–511, 2021.

World Health Organization. Covid-19 weekly epidemiological update. 2021. Accessed: 2021-03-25.

Marvin N. Wright and Andreas Ziegler. ranger: A fast implementation of random forests for high dimensional data in C++ and R. *Journal of Statistical Software*, 77(1):1–17, 2017. doi: 10.18637/jss.v077.i01.

Edward Wu and Johann A. Gagnon-Bartsch. The LOOP estimator: Adjusting for covariates in randomized experiments. *Evaluation Review*, 42(4):458–488, 2018. doi: 10.1177/0193841X18808003.

Li Yang and Anastasios A. Tsiatis. Efficiency study of estimators for a treatment effect in a pretest–posttest trial. *The American Statistician*, 55(4):314–321, 2001.

Min Zhang. Robust methods to improve efficiency and reduce bias in estimating survival curves in randomized clinical trials. *Lifetime Data Analysis*, 21(1):119–137, 2014. doi: 10.1007/s10985-014-9291-y.

Min Zhang and Peter B. Gilbert. Increasing the efficiency of prevention trials by incorporating baseline covariates. *Statistical Communications in Infectious Diseases*, 2(1), 2010. doi: 10.2202/1948-4690.1002.

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.

Wenjing Zheng and Mark J. van der Laan. Cross-validated targeted minimum-loss-based estimation. In *Targeted Learning*, pages 459–474. 2011.

*Supplementary Materials for*

# Optimizing Precision and Power by Machine Learning in Randomized Trials, with an Application to COVID-19.

Nicholas Williams[1], Michael Rosenblum[2], and Iván Díaz[*1]

[1]Division of Biostatistics, Department of Population Health Sciences, Weill Cornell Medicine.
[2]Department of Biostatistics, Johns Hopkins Bloomberg School of Public Health.

September 10, 2021

## A   Auxiliary covariates for estimation algorithm

Denote the survival function for $Y$ at time $k \in \{1, \ldots, K\}$ conditioned on study arm $a$ and baseline variables $w$ by

$$S(k, a, w) = P(Y > k \mid A = a, W = w). \tag{2}$$

Similarly, define the following function of the censoring distribution:

$$G(k, a, w) = P(C \geq k \mid A = a, W = w). \tag{3}$$

Under the assumption $C \perp\!\!\!\perp (Y, W) \mid A = a$ for each treatment arm $a$, we have $Y \perp\!\!\!\perp C \mid A, W$ and therefore $S(k, a, w)$ and $G(k, a, w)$ have the following product formula representations:

$$S(k, a, w) = \prod_{u=1}^{k} \{1 - m(u, a, w)\}; \qquad \Pi_C(k, a, w) = \prod_{u=0}^{k-1} \{1 - \pi_C(u, a, w)\}. \tag{4}$$

At each iteration of the estimation algorithm, the auxiliary covariates fr $\tilde{S}_{\text{TMLE}}$ and $\tilde{S}_{\text{IE-TMLE}}$ are constructed as follows:

$$H_{Y,k,u} = -\frac{\mathbb{1}\{A = a\}}{\widehat{\pi}_A(a, W)\widehat{\Pi}_C(u, a, W)} \frac{\widehat{S}(k, a, W)}{\widehat{S}(u, a, W)}$$

$$H_A = \frac{S(k, a, W)}{\pi_A(a, W)},$$

$$H_{C,k,u} = -\frac{\mathbb{1}\{A = a\}}{\widehat{\pi}_A(a, W)} \frac{\widehat{S}(k, a, W)}{\widehat{S}(u, a, W)} \frac{1}{\widehat{\Pi}_C(u + 1, a, W)},$$

where $\widehat{\pi}_A$, $\widehat{S}$, and $\widehat{\Pi}_C$ are the estimates in the current step of the iteration.

---

*corresponding author: ild2005@med.cornell.edu

# B Efficiency theory

Before proving Theorem 1 given in the paper, we introduce some notation and efficiency theory for estimation of $S(k, a)$. We will use the notation of Díaz et al. (2019). First, we encode a single participant's data vector $O = (W, A, \Delta = \mathbb{1}\{Y \leq C\}, \widetilde{Y} = \min(C, Y))$ using the following longitudinal data structure:

$$O = (W, A, R_0, L_1, R_1, L_2 \ldots, R_{K-1}, L_K), \tag{5}$$

where $R_u = \mathbb{1}\{\widetilde{Y} = u, \Delta = 0\}$ and $L_u = \mathbb{1}\{\widetilde{Y} = u, \Delta = 1\}$, for $u \in \{0, \ldots, K\}$. The sequence $R_0, L_1, R_1, L_2 \ldots, R_{K-1}, L_K$ in the above display consists of all 0's until the first time that either the event is observed or censoring occurs, i.e., time $u = \widetilde{Y}$. In the former case $L_u = 1$; otherwise $R_u = 1$. For a random variable $X$, we denote its history through time $u$ as $\bar{X}_u = (X_0, \ldots, X_u)$. For a given scalar $x$, the expression $\bar{X}_u = x$ denotes element-wise equality. The corresponding vector (5) for participant $i$ is denoted by $(W_i, A_i, R_{0,i}, L_{1,i}, R_{1,i}, L_{2,i} \ldots, R_{K-1,i}, L_{K,i})$.

Define the following indicator variables for each $u \geq 1$:

$$I_u = \mathbb{1}\{\bar{R}_{u-1} = 0, \bar{L}_{u-1} = 0\}, \qquad J_u = \mathbb{1}\{\bar{R}_{u-1} = 0, \bar{L}_u = 0\}.$$

The variable $I_u$ is the indicator based on the data through time $u - 1$ that a participant is at risk of the event being observed at time $u$; in other words, $I_u = 1$ means that all the variables $R_0, L_1, R_1, L_2 ..., L_{u-1}, R_{u-1}$ in the data vector (5) equal 0, which makes it possible that $L_u = 1$. Analogously, $J_u$ is the indicator based on the outcome data through time $u$ and censoring data before time $u$ that a participant is at risk of censoring at time $u$. By convention we let $J_0 = 1$.

The efficient influence function for estimation of $S(k, a)$ (see Moore and van der Laan 2009) is equal to:

$$D_\eta(O) = \sum_{u=1}^k I_u \times H_Y(u, A, W) \{L_u - m(u, a, W)\} + S(k, a, W) - S(k, a), \tag{6}$$

where we have explicitly added the dependence of the auxiliary covariate $H_Y$ on $(u, A, W)$ to the notation, and have denoted the nuisance parameters with $\eta = (m, \pi_A, \pi_C)$. In what follows we will use $\theta = S(k, a)$, and will use $\theta(\eta_1)$ to refer to the target parameter evaluated at a specific distribution implied by $\eta_1$. We will denote $\mathsf{P}f = \int f(o)\mathrm{d}\mathsf{P}(o)$, and $\mathsf{P}h(t, a, W) = \int h(t, a, w)\mathrm{d}\mathsf{P}(w)$ for functions $f$ and $h$. The efficient influence function has important implications for estimation of $S(k, a)$. First, the variance of $D_\eta(O)$ is the non-parametric efficiency bound, meaning that it is the smallest possible variance achievable by any regular estimator (Bickel et al. 1997). Second, the efficient influence function characterizes the first order bias of a plug-in estimator based on data-adaptive regression. Correction for this first order bias will allow us to establish normality of the estimators. Specifically, for any estimate $\widehat{\eta}$ we have the following first order expansion around the true parameter value $\theta(\eta)$, proved in Lemma 1 in the Supplementary materials of Díaz et al. (2018):

$$\theta(\widehat{\eta}) - \theta(\eta) = -\mathsf{P}D_{\widehat{\eta}} + \mathsf{Rem}_1(\widehat{\eta}), \tag{7}$$

where $\mathsf{Rem}_1$ is a second order remainder term given by

$$\mathsf{Rem}_1(\widehat{\eta}) = -\sum_{u=1}^k \int \frac{\widehat{S}(k, a, w)}{\widehat{S}(u, a, w)} S(u-1, a, w)\{m(u, a, w) - \widehat{m}(u, a, w)\} \left\{ \frac{\pi_A(a, w)\Pi_C(u, a, w)}{\widehat{\pi}_A(a, w)\widehat{\Pi}_C(u, a, w)} - 1 \right\} \mathrm{d}\mathsf{P}(w),$$

and $\theta(\widehat{\eta})$ is the substitution estimator

$$\theta(\widehat{\eta}) = \frac{1}{n}\sum_{i=1}^n \prod_{u=1}^k \{1 - \widehat{m}(u, a, W_i)\}.$$

2

The following proposition establishing the robustness of $D_\eta$ to misspecification of the model $m$ will be useful to prove consistency of the estimator.

**Proposition 1.** *Let $\eta_1 = (m_1, \pi_{A,1}, \pi_{C,1})$ be such that either $m_1 = m$ or $(\pi_{A,1}, \pi_{C,1}) = (\pi_A, \pi_C)$. Then $PD_{\eta_1} = 0$.*

Recall the cross-fitting procedure described in the main document as follows. Let $\mathcal{V}_1, \ldots, \mathcal{V}_J$ denote a random partition of the index set $\{1, \ldots, n\}$ into $J$ prediction sets of approximately the same size. That is, $\mathcal{V}_j \subset \{1, \ldots, n\}$; $\bigcup_{j=1}^{J} \mathcal{V}_j = \{1, \ldots, n\}$; and $\mathcal{V}_j \cap \mathcal{V}_{j'} = \emptyset$. In addition, for each $j$, the associated training sample is given by $\mathcal{T}_j = \{1, \ldots, n\} \setminus \mathcal{V}_j$. Let $\widehat{m}_j$ denote the prediction algorithm trained in $\mathcal{T}_j$. Letting $j(i)$ denote the index of the validation set which contains observation $i$, cross-fitting entails using only observations in $\mathcal{T}_{j(i)}$ for fitting models when making predictions about observation $i$. That is, the outcome predictions for each subject $i$ are given by $\widehat{m}_{j(i)}(u, a, W_i)$. Since only $\widehat{m}$ and not $(\widehat{\pi}_A, \widehat{\pi}_C)$ is cross-fitted, we let $\widehat{\eta}_{j(i)} = (\widehat{m}_{j(i)}, \widehat{\pi}_A, \widehat{\pi}_C)$ and $\widetilde{\eta}_{j(i)} = (\widetilde{m}_{j(i)}, \widehat{\pi}_A, \widehat{\pi}_C)$.

# C  Proof of Theorem 1

In what follows we let $P_{n,j}$ denote the empirical distribution of the prediction set $\mathcal{V}_j$, and let $G_{n,j}$ denote the associated empirical process $\sqrt{n/J}(P_{n,j} - P)$. Let $G_n$ denote the empirical process $\sqrt{n}(P_n - P)$. We use $E(g(O_1, \ldots, O_n))$ to denote expectation with respect to the joint distribution of $(O_1, \ldots, O_n)$, and use $a_n \lesssim b_n$ to mean $a_n \le cb_n$ for universal constants $c$. The following lemmas will be useful in the proof of the theorem.

**Lemma 1.** *Assume A1, A4, and A5 . Then we have $\Pi_C(k, a, w)$ does not depend on $w$, and $\pi_A(a, w)$ does not depend on $w$. Furthermore, we have*

$$\sqrt{n}\{\widehat{\Pi}_C(k, a) - \Pi_C(k, a)\} = G_n \Delta_{k,a} + o_P(1),$$
$$\sqrt{n}\{\widehat{\pi}_A(a) - \pi_A(a)\} = G_n \Lambda_a + o_P(1),$$

*for mean-zero functions $\Delta_{k,a}(O_i)$ and $\Lambda_a(O_i)$ of $(k, a)$ and $O_i$ that do not depend on $W_i$.*

*Proof.* This lemma follows by application of the Delta method to the non-parametric maximum likelihood estimators $\widehat{\pi}_A$ and $\widehat{\Pi}_C$. $\square$

**Lemma 2.** *For two sequences $a_1, \ldots, a_m$ and $b_1, \ldots, b_m$ we have*

$$\prod_{t=1}^{m}(1 - a_t) - \prod_{t=1}^{m}(1 - b_t) = \sum_{t=1}^{m}\left\{\left[\prod_{k=1}^{t-1}(1 - a_k)\right](b_t - a_t)\left[\prod_{k=t+1}^{m}(1 - b_k)\right]\right\}.$$

*Proof.* Replace $(b_t - a_t)$ by $(1 - a_t) - (1 - b_t)$ in the right hand side and expand the sum to notice it is a telescoping sum. $\square$

The proof of Theorem 1 proceeds as follows.

*Proof.* Since censoring is completely at random by A1, we have $\theta = S(k, a) = \int S(k, a, w)dP(w)$. Let $\widetilde{\theta} = \widetilde{S}_{\text{TMLE}}(k, a)$. Define $\sigma^2 = \text{Var}[D_{\eta_1}(O)]$, where $\eta_1 = (m_1, \pi_A, \pi_C)$, and let

$$\widetilde{\Theta}_n = \sqrt{n}(\widetilde{\theta} - \theta)/\widetilde{\sigma}$$
$$\check{\Theta}_n = \sqrt{n}(\widetilde{\theta} - \theta)/\sigma$$
$$\Theta_n = G_n D_{\eta_1}/\sigma.$$

3

First, note that $\Theta_n \rightsquigarrow N(0,1)$ by the central limit theorem. We will now show that $|\widetilde{\Theta}_n - \Theta_n| = o_P(1)$, which would yield the result in the theorem. First, note that

$$
\begin{aligned}
|\widetilde{\Theta}_n - \Theta_n| &= |(\check{\Theta}_n - \Theta_n)(\sigma/\widetilde{\sigma}) + \Theta_n(\sigma - \widetilde{\sigma})/\widetilde{\sigma}| \\
&\leq |\check{\Theta}_n - \Theta_n|\,|\sigma/\widetilde{\sigma}| + |\Theta_n|\,|\sigma/\widetilde{\sigma} - 1| \\
&\lesssim |\check{\Theta}_n - \Theta_n| + o_P(1),
\end{aligned}
$$

where the last inequality follows because $|\sigma/\widetilde{\sigma} - 1| = o_P(1)$ (which follows by Lemma 1 and A3) and because $|\Theta_n| = O_P(1)$ by the central limit theorem. We will now show that $|\check{\Theta}_n - \Theta_n| = o_P(1)$.

An application of (7) with $\widehat{\eta} = \widetilde{\eta}$ yields

$$
\begin{aligned}
\sqrt{n}(\widetilde{\theta} - \theta) &= -\sqrt{n}\mathsf{P}D_{\widetilde{\eta}} + \sqrt{n}\mathsf{Rem}_1(\widetilde{\eta}) \\
&= \sqrt{n}(\mathsf{P}_n - \mathsf{P})D_{\widetilde{\eta}} + \sqrt{n}\mathsf{Rem}_1(\widetilde{\eta}) \\
&= \mathsf{G}_n D_{\eta_1} + \mathsf{G}_n(D_{\widetilde{\eta}} - D_{\eta_1}) + \sqrt{n}\mathsf{Rem}_1(\widetilde{\eta}),
\end{aligned}
$$

where the second equality follows because $\mathsf{P}_n D_{\widetilde{\eta}} = 0$ by definition of $\widetilde{\eta}$ (see Díaz et al. 2019). This implies

$$
\check{\Theta}_n - \Theta_n = B_{n,2} + B_{n,1},
$$

where $B_{n,2} = \mathsf{G}_n(D_{\widetilde{\eta}} - D_{\eta_1})$ and $B_{n,1} = \sqrt{n}\mathsf{Rem}_1(\widetilde{\eta})$.

We first tackle the case of A6.2, where the estimators for $m$ are cross-fitted. Note that

$$
B_{n,2} = \frac{1}{\sqrt{J}}\sum_{j=1}^{J}\mathsf{G}_{n,j}(D_{\widetilde{\eta}_j} - D_{\eta_1}),
$$

and that $D_{\widetilde{\eta}_j}$ depends on the full sample through the estimate of the parameter $\varepsilon$ of the logistic tilting model. To make this dependence explicit, we introduce the notation $D_{\widehat{\eta}_j, \widehat{\varepsilon}} = D_{\widetilde{\eta}_j}$. Let $\varepsilon_1$ denote the probability limit of $\widehat{\varepsilon}$, which exists and is finite by Assumption A7. We can find a deterministic sequence $\delta_n \to 0$ satisfying $P(|\widehat{\varepsilon} - \varepsilon_1| < \delta_n) \to 1$. Let $\mathcal{F}_n^j = \{D_{\widehat{\eta}_j, \varepsilon} - D_{\eta_1} : |\varepsilon - \varepsilon_1| < \delta_n\}$. Because the function $\widehat{\eta}_j$ is fixed given the training data, we can apply Theorem 2.14.2 of van der Vaart and Wellner (1996) to obtain

$$
E\left\{\sup_{f \in \mathcal{F}_n^j}|\mathsf{G}_{n,j}f|\,\Big|\,\mathcal{T}_j\right\} \lesssim \|F_n^j\|\int_0^1\sqrt{1 + N_{[\,]}(\alpha\|F_n^j\|, \mathcal{F}_n^j, L_2(\mathsf{P}))}\,\mathrm{d}\alpha, \tag{8}
$$

where $N_{[\,]}(\alpha\|F_n^j\|, \mathcal{F}_n^j, L_2(\mathsf{P}))$ is the bracketing number and we take $F_n^j = \sup_{\varepsilon:|\varepsilon - \varepsilon_1| < \delta_n}|D_{\widehat{\eta}_j, \varepsilon} - D_{\eta_1}|$ as an envelope for the class $\mathcal{F}_n^j$. Theorem 2.7.2 of van der Vaart and Wellner (1996) shows

$$
\log N_{[\,]}(\alpha\|F_n^j\|, \mathcal{F}_n^j, L_2(\mathsf{P})) \lesssim \frac{1}{\alpha\|F_n^j\|}.
$$

This shows

$$
\begin{aligned}
\|F_n^j\|\int_0^1\sqrt{1 + N_{[\,]}(\alpha\|F_n^j\|, \mathcal{F}_n^j, L_2(\mathsf{P}))}\,\mathrm{d}\alpha &\lesssim \int_0^1\sqrt{\|F_n^j\|^2 + \frac{\|F_n^j\|}{\alpha}}\,\mathrm{d}\alpha \\
&\leq \|F_n^j\| + \|F_n^j\|^{1/2}\int_0^1\frac{1}{\alpha^{1/2}}\,\mathrm{d}\alpha \\
&\leq \|F_n^j\| + 2\|F_n^j\|^{1/2}.
\end{aligned}
$$

Since $D_{\widehat{\eta}_j,\widehat{\varepsilon}} \to D_{\eta 1}$ and $\delta_n \to 0$, $\|F_n^j\| = o_P(1)$. The above argument shows that $\sup_{f \in \mathcal{F}_n^j} |\mathsf{G}_{n,j} f| = o_P(1)$ for each $j$, conditional on $\mathcal{T}_j$. Thus $|B_{n,2}| = o_P(1)$.

In the case of A6.1, where the estimators for $m$ are *not* cross-fitted but belong in a parametric family, standard empirical process theory such as Example 19.7 of van der Vaart (1998) shows that $D_{\widetilde{\eta}}$ takes values in a Donsker class. Therefore, an application of Theorem 19.24 of van der Vaart (1998) yields $|B_{n,2}| = o_P(1)$.

We now show that $|B_{n,1}| = o_P(1)$. First, Lemma 1 along with the Delta method show that

$$\sqrt{n} \left\{ \frac{\pi_A(a,w)\Pi_C(u,a,w)}{\widetilde{\pi}_A(a,w)\widetilde{\Pi}_C(u,a,w)} - 1 \right\} = \mathsf{G}_n \Gamma_{k,a} + o_P(1),$$

for some function $\Gamma_{k,a}(O)$ not depending on $W$. Thus

$$
\begin{aligned}
B_{n,1} &= -\sum_{u=1}^{k} \int \frac{\widetilde{S}(k,a,w)}{\widetilde{S}(u,a,w)} S(u-1,a,w)\{m(u,a,w) - \widetilde{m}(u,a,w)\} \{\mathsf{G}_n \Gamma_{k,a} + o_P(1)\} \, \mathrm{d}P(w) \\
&= -\mathsf{G}_n \Gamma_{k,a} \sum_{u=1}^{k} \int \frac{\widetilde{S}(k,a,w)}{\widetilde{S}(u,a,w)} S(u-1,a,w)\{m(u,a,w) - \widetilde{m}(u,a,w)\} \mathrm{d}P(w) + o_P(1) \\
&= \mathsf{G}_n \Gamma_{k,a} \int \{S(k,a,w) - \widetilde{S}(k,a,w)\} \mathrm{d}P(w) + o_P(1),
\end{aligned}
$$

where the last equality follows from Lemma 2. Expression (7) together with the assumptions of the theorem and Proposition 1 show that the estimator $\widetilde{\theta}$ is consistent, and thus

$$\int \{S(k,a,w) - \widetilde{S}(k,a,w)\} \mathrm{d}P(w) = o_P(1).$$

The central limit theorem shows that $\mathsf{G}_n \Gamma_{k,a} = O_P(1)$, which yields $|B_{n,1}| = o_P(1)$, concluding the proof of the theorem. $\qquad\square$

# D Tables with simulation results

## D.1 Results for simulations with a positive effect and where the covariates are prognostic of the outcome

Table 1: Simulation results for the RMST of time to intubation or death at day 14 under a positive effect and covariates with prognostic power.

| Estimator | $n$ | Effect size | P(Reject $H_0$) | $n \times$ MSE | $n \times$ Var | \|Bias\| | Rel. eff. |
|---|---|---|---|---|---|---|---|
| CF-MARS | 100 | 1.04 | 0.26 | 83.22 | 83.07 | 0.04 | 1.29 |
| CF-RF | 100 | 1.04 | 0.24 | 71.13 | 71.13 | 0.00 | 1.10 |
| CF-XGBoost | 100 | 1.04 | 0.23 | 70.33 | 70.33 | 0.00 | 1.09 |
| $\ell_1$-LR | 100 | 1.04 | 0.28 | 58.49 | 58.47 | 0.02 | 0.90 |
| LR | 100 | 1.04 | 0.34 | 66.61 | 66.61 | 0.00 | 1.03 |
| MARS | 100 | 1.04 | 0.27 | 70.14 | 69.26 | 0.09 | 1.09 |
| RF | 100 | 1.04 | 0.62 | 51.07 | 44.94 | 0.25 | 0.79 |
| Unadjusted | 100 | 1.04 | 0.25 | 64.64 | 64.64 | 0.01 | 1.00 |
| XGBoost | 100 | 1.04 | 0.27 | 64.07 | 64.06 | 0.01 | 0.99 |
| CF-MARS | 500 | 1.04 | 0.85 | 61.36 | 61.35 | 0.00 | 0.93 |
| CF-RF | 500 | 1.04 | 0.85 | 60.10 | 60.10 | 0.00 | 0.91 |
| CF-XGBoost | 500 | 1.04 | 0.82 | 64.70 | 64.69 | 0.01 | 0.98 |
| $\ell_1$-LR | 500 | 1.04 | 0.87 | 58.00 | 58.00 | 0.00 | 0.88 |
| LR | 500 | 1.04 | 0.87 | 60.49 | 60.49 | 0.00 | 0.92 |
| MARS | 500 | 1.04 | 0.86 | 58.01 | 57.95 | 0.01 | 0.88 |
| RF | 500 | 1.04 | 0.97 | 66.02 | 55.76 | 0.14 | 1.00 |
| Unadjusted | 500 | 1.04 | 0.82 | 65.85 | 65.85 | 0.00 | 1.00 |
| XGBoost | 500 | 1.04 | 0.88 | 58.13 | 58.08 | 0.01 | 0.88 |
| CF-MARS | 1500 | 1.04 | 1.00 | 60.51 | 60.50 | 0.00 | 0.93 |
| CF-RF | 1500 | 1.04 | 1.00 | 47.88 | 47.88 | 0.00 | 0.74 |
| CF-XGBoost | 1500 | 1.04 | 1.00 | 61.51 | 61.50 | 0.00 | 0.95 |
| $\ell_1$-LR | 1500 | 1.04 | 1.00 | 58.15 | 58.14 | 0.00 | 0.89 |
| LR | 1500 | 1.04 | 1.00 | 57.15 | 57.15 | 0.00 | 0.88 |
| MARS | 1500 | 1.04 | 1.00 | 60.15 | 60.15 | 0.00 | 0.92 |
| RF | 1500 | 1.04 | 1.00 | 61.29 | 50.61 | 0.08 | 0.94 |
| Unadjusted | 1500 | 1.04 | 1.00 | 65.07 | 65.06 | 0.00 | 1.00 |
| XGBoost | 1500 | 1.04 | 1.00 | 56.24 | 55.91 | 0.01 | 0.86 |

Table 2: Simulation results for the RD of time to intubation or death at day 7 under a positive effect and covariates with prognostic power.

| Estimator | $n$ | Effect size | P(Reject $H_0$) | $n \times$ MSE | $n \times$ Var | \|Bias\| | Rel. eff. |
|---|---|---|---|---|---|---|---|
| CF-MARS | 100 | 0.10 | 0.27 | 0.74 | 0.73 | 0.00 | 1.23 |
| CF-RF | 100 | 0.10 | 0.25 | 0.64 | 0.64 | 0.00 | 1.08 |
| CF-XGBoost | 100 | 0.10 | 0.25 | 0.64 | 0.64 | 0.00 | 1.07 |
| $\ell_1$-LR | 100 | 0.10 | 0.28 | 0.55 | 0.55 | 0.00 | 0.92 |
| LR | 100 | 0.10 | 0.32 | 0.60 | 0.60 | 0.00 | 1.01 |
| MARS | 100 | 0.10 | 0.28 | 0.64 | 0.63 | 0.01 | 1.07 |
| RF | 100 | 0.10 | 0.50 | 0.56 | 0.52 | 0.02 | 0.94 |
| Unadjusted | 100 | 0.10 | 0.26 | 0.60 | 0.60 | 0.00 | 1.00 |
| XGBoost | 100 | 0.10 | 0.28 | 0.60 | 0.60 | 0.00 | 1.01 |
| CF-MARS | 500 | 0.10 | 0.85 | 0.57 | 0.57 | 0.00 | 0.92 |
| CF-RF | 500 | 0.10 | 0.85 | 0.56 | 0.56 | 0.00 | 0.91 |
| CF-XGBoost | 500 | 0.10 | 0.84 | 0.59 | 0.59 | 0.00 | 0.96 |
| $\ell_1$-LR | 500 | 0.10 | 0.87 | 0.54 | 0.54 | 0.00 | 0.89 |
| LR | 500 | 0.10 | 0.86 | 0.57 | 0.57 | 0.00 | 0.92 |
| MARS | 500 | 0.10 | 0.85 | 0.56 | 0.55 | 0.00 | 0.90 |
| RF | 500 | 0.10 | 0.94 | 0.66 | 0.58 | 0.01 | 1.07 |
| Unadjusted | 500 | 0.10 | 0.83 | 0.61 | 0.61 | 0.00 | 1.00 |
| XGBoost | 500 | 0.10 | 0.87 | 0.55 | 0.55 | 0.00 | 0.90 |
| CF-MARS | 1500 | 0.10 | 1.00 | 0.58 | 0.58 | 0.00 | 0.96 |
| CF-RF | 1500 | 0.10 | 1.00 | 0.47 | 0.47 | 0.00 | 0.77 |
| CF-XGBoost | 1500 | 0.10 | 1.00 | 0.56 | 0.56 | 0.00 | 0.94 |
| $\ell_1$-LR | 1500 | 0.10 | 1.00 | 0.56 | 0.56 | 0.00 | 0.93 |
| LR | 1500 | 0.10 | 1.00 | 0.55 | 0.55 | 0.00 | 0.92 |
| MARS | 1500 | 0.10 | 1.00 | 0.57 | 0.57 | 0.00 | 0.94 |
| RF | 1500 | 0.10 | 1.00 | 0.61 | 0.54 | 0.01 | 1.01 |
| Unadjusted | 1500 | 0.10 | 1.00 | 0.60 | 0.60 | 0.00 | 1.00 |
| XGBoost | 1500 | 0.10 | 1.00 | 0.54 | 0.53 | 0.00 | 0.89 |

Table 3: Simulation results for the LOR of the modified WHO scale at day 14 under a positive effect and covariates with prognostic power.

| Estimator | $n$ | Effect size | P(Reject $H_0$) | $n \times$ MSE | $n \times$ Var | \|Bias\| | Rel. eff. |
|---|---|---|---|---|---|---|---|
| CF-MARS | 100 | 0.60 | 0.35 | 73.00 | 70.25 | 0.17 | 0.13 |
| CF-RF | 100 | 0.60 | 0.36 | 37.97 | 36.75 | 0.11 | 0.07 |
| CF-XGBoost | 100 | 0.60 | 0.36 | 34.80 | 34.30 | 0.07 | 0.06 |
| $\ell_1$-LR | 100 | 0.60 | 0.45 | 57.23 | 54.57 | 0.16 | 0.10 |
| LR | 100 | 0.60 | 0.48 | 472.16 | 392.98 | 0.89 | 0.81 |
| MARS | 100 | 0.60 | 0.45 | 192.05 | 176.11 | 0.40 | 0.33 |
| RF | 100 | 0.60 | 0.48 | 32.02 | 32.01 | 0.01 | 0.05 |
| Unadjusted | 100 | 0.60 | 0.42 | 582.27 | 441.29 | 1.19 | 1.00 |
| XGBoost | 100 | 0.60 | 0.41 | 49.81 | 49.47 | 0.06 | 0.09 |
| CF-MARS | 500 | 0.60 | 0.90 | 30.48 | 30.27 | 0.02 | 0.54 |
| CF-RF | 500 | 0.60 | 0.93 | 16.15 | 16.13 | 0.01 | 0.29 |
| CF-XGBoost | 500 | 0.60 | 0.92 | 18.88 | 18.86 | 0.01 | 0.34 |
| $\ell_1$-LR | 500 | 0.60 | 0.93 | 18.74 | 18.68 | 0.01 | 0.33 |
| LR | 500 | 0.60 | 0.92 | 32.86 | 32.53 | 0.03 | 0.59 |
| MARS | 500 | 0.60 | 0.93 | 32.10 | 32.07 | 0.01 | 0.57 |
| RF | 500 | 0.60 | 0.99 | 15.84 | 15.10 | 0.04 | 0.28 |
| Unadjusted | 500 | 0.60 | 0.86 | 56.01 | 55.76 | 0.02 | 1.00 |
| XGBoost | 500 | 0.60 | 0.95 | 16.29 | 15.42 | 0.04 | 0.29 |
| CF-MARS | 1500 | 0.60 | 1.00 | 17.71 | 17.71 | 0.00 | 0.88 |
| CF-RF | 1500 | 0.60 | 1.00 | 12.28 | 12.24 | 0.01 | 0.61 |
| CF-XGBoost | 1500 | 0.60 | 1.00 | 15.28 | 15.27 | 0.00 | 0.76 |
| $\ell_1$-LR | 1500 | 0.60 | 1.00 | 17.90 | 17.89 | 0.00 | 0.89 |
| LR | 1500 | 0.60 | 1.00 | 17.12 | 17.04 | 0.01 | 0.85 |
| MARS | 1500 | 0.60 | 1.00 | 16.72 | 16.71 | 0.00 | 0.83 |
| RF | 1500 | 0.60 | 1.00 | 9.93 | 9.21 | 0.02 | 0.49 |
| Unadjusted | 1500 | 0.60 | 1.00 | 20.20 | 20.20 | 0.00 | 1.00 |
| XGBoost | 1500 | 0.60 | 1.00 | 11.21 | 10.81 | 0.02 | 0.55 |

Table 4: Simulation results for the MW estimand of the modified WHO scale at day 14 under a positive effect and covariates with prognostic power.

| Estimator | $n$ | Effect size | P(Reject $H_0$) | $n \times$ MSE | $n \times$ Var | \|Bias\| | Rel. eff. |
|---|---|---|---|---|---|---|---|
| CF-MARS | 100 | 0.46 | 0.24 | 0.66 | 0.66 | 0.00 | 2.52 |
| CF-RF | 100 | 0.46 | 0.17 | 0.24 | 0.24 | 0.00 | 0.90 |
| CF-XGBoost | 100 | 0.46 | 0.15 | 0.24 | 0.24 | 0.00 | 0.92 |
| $\ell_1$-LR | 100 | 0.46 | 0.18 | 0.23 | 0.23 | 0.00 | 0.86 |
| LR | 100 | 0.46 | 0.45 | 1.68 | 1.66 | 0.01 | 6.40 |
| MARS | 100 | 0.46 | 0.23 | 0.49 | 0.49 | 0.00 | 1.85 |
| RF | 100 | 0.46 | 0.22 | 0.20 | 0.19 | 0.01 | 0.76 |
| Unadjusted | 100 | 0.46 | 0.14 | 0.26 | 0.26 | 0.00 | 1.00 |
| XGBoost | 100 | 0.46 | 0.19 | 0.21 | 0.21 | 0.00 | 0.81 |
| CF-MARS | 500 | 0.46 | 0.51 | 0.52 | 0.51 | 0.01 | 2.01 |
| CF-RF | 500 | 0.46 | 0.56 | 0.20 | 0.20 | 0.00 | 0.79 |
| CF-XGBoost | 500 | 0.46 | 0.53 | 0.23 | 0.22 | 0.00 | 0.87 |
| $\ell_1$-LR | 500 | 0.46 | 0.53 | 0.21 | 0.21 | 0.00 | 0.82 |
| LR | 500 | 0.46 | 0.54 | 0.24 | 0.24 | 0.00 | 0.92 |
| MARS | 500 | 0.46 | 0.51 | 0.26 | 0.26 | 0.00 | 1.02 |
| RF | 500 | 0.46 | 0.76 | 0.14 | 0.14 | 0.00 | 0.55 |
| Unadjusted | 500 | 0.46 | 0.44 | 0.26 | 0.26 | 0.00 | 1.00 |
| XGBoost | 500 | 0.46 | 0.59 | 0.20 | 0.20 | 0.00 | 0.76 |
| CF-MARS | 1500 | 0.46 | 0.93 | 0.23 | 0.22 | 0.00 | 0.86 |
| CF-RF | 1500 | 0.46 | 0.99 | 0.15 | 0.14 | 0.00 | 0.57 |
| CF-XGBoost | 1500 | 0.46 | 0.97 | 0.19 | 0.18 | 0.00 | 0.73 |
| $\ell_1$-LR | 1500 | 0.46 | 0.93 | 0.23 | 0.22 | 0.00 | 0.86 |
| LR | 1500 | 0.46 | 0.94 | 0.23 | 0.21 | 0.00 | 0.85 |
| MARS | 1500 | 0.46 | 0.94 | 0.21 | 0.21 | 0.00 | 0.80 |
| RF | 1500 | 0.46 | 1.00 | 0.11 | 0.11 | 0.00 | 0.40 |
| Unadjusted | 1500 | 0.46 | 0.88 | 0.27 | 0.26 | 0.00 | 1.00 |
| XGBoost | 1500 | 0.46 | 0.99 | 0.14 | 0.13 | 0.00 | 0.52 |

## D.2 Results for simulations with null treatment effect and where the covariates are prognostic of the outcome

Table 5: Simulation results for the RMST of time to intubation or death at day 14 under null treatment effect and covariates with prognostic power.

| Estimator | $n$ | Effect size | P(Reject $H_0$) | $n \times$ MSE | $n \times$ Var | \|Bias\| | Rel. eff. |
|---|---|---|---|---|---|---|---|
| CF-MARS | 100 | 0.00 | 0.10 | 126.72 | 121.56 | 0.23 | 1.44 |
| CF-RF | 100 | 0.00 | 0.05 | 88.87 | 88.81 | 0.02 | 1.01 |
| CF-XGBoost | 100 | 0.00 | 0.05 | 90.87 | 90.84 | 0.02 | 1.03 |
| $\ell_1$-LR | 100 | 0.00 | 0.07 | 82.16 | 82.14 | 0.01 | 0.93 |
| LR | 100 | 0.00 | 0.09 | 86.83 | 86.82 | 0.01 | 0.99 |
| MARS | 100 | 0.00 | 0.09 | 88.35 | 87.88 | 0.07 | 1.00 |
| RF | 100 | 0.00 | 0.30 | 51.79 | 51.78 | 0.01 | 0.59 |
| Unadjusted | 100 | 0.00 | 0.06 | 87.92 | 87.85 | 0.03 | 1.00 |
| XGBoost | 100 | 0.00 | 0.05 | 78.09 | 78.09 | 0.01 | 0.89 |
| CF-MARS | 500 | 0.00 | 0.05 | 80.81 | 80.76 | 0.01 | 0.95 |
| CF-RF | 500 | 0.00 | 0.05 | 75.68 | 75.68 | 0.00 | 0.89 |
| CF-XGBoost | 500 | 0.00 | 0.05 | 83.71 | 83.70 | 0.00 | 0.98 |
| $\ell_1$-LR | 500 | 0.00 | 0.05 | 73.05 | 73.03 | 0.01 | 0.86 |
| LR | 500 | 0.00 | 0.06 | 79.64 | 79.63 | 0.00 | 0.94 |
| MARS | 500 | 0.00 | 0.05 | 78.44 | 78.44 | 0.00 | 0.92 |
| RF | 500 | 0.00 | 0.29 | 48.37 | 48.37 | 0.00 | 0.57 |
| Unadjusted | 500 | 0.00 | 0.05 | 85.12 | 85.10 | 0.01 | 1.00 |
| XGBoost | 500 | 0.00 | 0.07 | 72.98 | 72.98 | 0.00 | 0.86 |
| CF-MARS | 1500 | 0.00 | 0.05 | 73.18 | 73.18 | 0.00 | 0.85 |
| CF-RF | 1500 | 0.00 | 0.05 | 56.36 | 56.35 | 0.00 | 0.66 |
| CF-XGBoost | 1500 | 0.00 | 0.05 | 80.36 | 80.36 | 0.00 | 0.94 |
| $\ell_1$-LR | 1500 | 0.00 | 0.05 | 75.52 | 75.51 | 0.00 | 0.88 |
| LR | 1500 | 0.00 | 0.06 | 75.39 | 75.38 | 0.00 | 0.88 |
| MARS | 1500 | 0.00 | 0.05 | 76.55 | 76.51 | 0.00 | 0.89 |
| RF | 1500 | 0.00 | 0.30 | 27.65 | 27.64 | 0.00 | 0.32 |
| Unadjusted | 1500 | 0.00 | 0.05 | 85.87 | 85.83 | 0.00 | 1.00 |
| XGBoost | 1500 | 0.00 | 0.07 | 64.40 | 64.40 | 0.00 | 0.75 |

Table 6: Simulation results for the RD of time to intubation or death at day 7 under null treatment effect and covariates with prognostic power.

| Estimator | $n$ | Effect size | P(Reject $H_0$) | $n \times$ MSE | $n \times$ Var | \|Bias\| | Rel. eff. |
|---|---|---|---|---|---|---|---|
| CF-MARS | 100 | 0.00 | 0.09 | 1.00 | 0.96 | 0.02 | 1.36 |
| CF-RF | 100 | 0.00 | 0.05 | 0.75 | 0.75 | 0.00 | 1.02 |
| CF-XGBoost | 100 | 0.00 | 0.05 | 0.75 | 0.75 | 0.00 | 1.03 |
| $\ell_1$-LR | 100 | 0.00 | 0.07 | 0.69 | 0.69 | 0.00 | 0.95 |
| LR | 100 | 0.00 | 0.09 | 0.72 | 0.72 | 0.00 | 0.98 |
| MARS | 100 | 0.00 | 0.08 | 0.72 | 0.72 | 0.01 | 0.99 |
| RF | 100 | 0.00 | 0.26 | 0.51 | 0.51 | 0.00 | 0.70 |
| Unadjusted | 100 | 0.00 | 0.06 | 0.73 | 0.73 | 0.00 | 1.00 |
| XGBoost | 100 | 0.00 | 0.05 | 0.66 | 0.66 | 0.00 | 0.90 |
| CF-MARS | 500 | 0.00 | 0.05 | 0.69 | 0.68 | 0.00 | 0.95 |
| CF-RF | 500 | 0.00 | 0.04 | 0.63 | 0.63 | 0.00 | 0.88 |
| CF-XGBoost | 500 | 0.00 | 0.04 | 0.70 | 0.70 | 0.00 | 0.98 |
| $\ell_1$-LR | 500 | 0.00 | 0.05 | 0.63 | 0.63 | 0.00 | 0.88 |
| LR | 500 | 0.00 | 0.06 | 0.69 | 0.69 | 0.00 | 0.95 |
| MARS | 500 | 0.00 | 0.06 | 0.68 | 0.68 | 0.00 | 0.94 |
| RF | 500 | 0.00 | 0.25 | 0.45 | 0.45 | 0.00 | 0.62 |
| Unadjusted | 500 | 0.00 | 0.05 | 0.72 | 0.72 | 0.00 | 1.00 |
| XGBoost | 500 | 0.00 | 0.06 | 0.62 | 0.62 | 0.00 | 0.86 |
| CF-MARS | 1500 | 0.00 | 0.04 | 0.62 | 0.62 | 0.00 | 0.86 |
| CF-RF | 1500 | 0.00 | 0.05 | 0.49 | 0.49 | 0.00 | 0.67 |
| CF-XGBoost | 1500 | 0.00 | 0.05 | 0.68 | 0.68 | 0.00 | 0.94 |
| $\ell_1$-LR | 1500 | 0.00 | 0.05 | 0.64 | 0.64 | 0.00 | 0.89 |
| LR | 1500 | 0.00 | 0.05 | 0.65 | 0.65 | 0.00 | 0.90 |
| MARS | 1500 | 0.00 | 0.05 | 0.65 | 0.65 | 0.00 | 0.90 |
| RF | 1500 | 0.00 | 0.25 | 0.26 | 0.26 | 0.00 | 0.36 |
| Unadjusted | 1500 | 0.00 | 0.05 | 0.72 | 0.72 | 0.00 | 1.00 |
| XGBoost | 1500 | 0.00 | 0.07 | 0.56 | 0.56 | 0.00 | 0.78 |

Table 7: Simulation results for the LOR of the modified WHO scale at day 14 under null treatment effect and covariates with prognostic power.

| Estimator | $n$ | Effect size | P(Reject $H_0$) | $n \times$ MSE | $n \times$ Var | \|Bias\| | Rel. eff. |
|-----------|-----|-------------|-----------------|----------------|----------------|----------|-----------|
| CF-MARS | 100 | 0.00 | 0.07 | 21.27 | 21.27 | 0.00 | 0.55 |
| CF-RF | 100 | 0.00 | 0.06 | 16.93 | 16.93 | 0.01 | 0.44 |
| CF-XGBoost | 100 | 0.00 | 0.07 | 18.12 | 18.12 | 0.00 | 0.47 |
| $\ell_1$-LR | 100 | 0.00 | 0.09 | 28.01 | 28.01 | 0.00 | 0.73 |
| LR | 100 | 0.00 | 0.11 | 111.68 | 111.61 | 0.03 | 2.89 |
| MARS | 100 | 0.00 | 0.09 | 60.09 | 60.08 | 0.01 | 1.56 |
| RF | 100 | 0.00 | 0.13 | 21.91 | 21.90 | 0.01 | 0.57 |
| Unadjusted | 100 | 0.00 | 0.06 | 38.60 | 38.60 | 0.00 | 1.00 |
| XGBoost | 100 | 0.00 | 0.10 | 18.97 | 18.96 | 0.01 | 0.49 |
| CF-MARS | 500 | 0.00 | 0.05 | 14.70 | 14.70 | 0.00 | 0.88 |
| CF-RF | 500 | 0.00 | 0.05 | 12.57 | 12.57 | 0.00 | 0.75 |
| CF-XGBoost | 500 | 0.00 | 0.06 | 14.07 | 14.06 | 0.00 | 0.84 |
| $\ell_1$-LR | 500 | 0.00 | 0.06 | 13.51 | 13.51 | 0.00 | 0.81 |
| LR | 500 | 0.00 | 0.06 | 13.76 | 13.75 | 0.00 | 0.82 |
| MARS | 500 | 0.00 | 0.05 | 12.75 | 12.75 | 0.00 | 0.76 |
| RF | 500 | 0.00 | 0.14 | 7.34 | 7.34 | 0.00 | 0.44 |
| Unadjusted | 500 | 0.00 | 0.05 | 16.74 | 16.74 | 0.00 | 1.00 |
| XGBoost | 500 | 0.00 | 0.08 | 11.07 | 11.07 | 0.00 | 0.66 |
| CF-MARS | 1500 | 0.00 | 0.05 | 13.35 | 13.35 | 0.00 | 0.82 |
| CF-RF | 1500 | 0.00 | 0.05 | 8.36 | 8.36 | 0.00 | 0.51 |
| CF-XGBoost | 1500 | 0.00 | 0.05 | 11.03 | 11.03 | 0.00 | 0.68 |
| $\ell_1$-LR | 1500 | 0.00 | 0.05 | 12.87 | 12.87 | 0.00 | 0.79 |
| LR | 1500 | 0.00 | 0.05 | 13.27 | 13.27 | 0.00 | 0.81 |
| MARS | 1500 | 0.00 | 0.05 | 12.52 | 12.52 | 0.00 | 0.77 |
| RF | 1500 | 0.00 | 0.15 | 4.56 | 4.56 | 0.00 | 0.28 |
| Unadjusted | 1500 | 0.00 | 0.05 | 16.30 | 16.29 | 0.00 | 1.00 |
| XGBoost | 1500 | 0.00 | 0.16 | 6.37 | 6.37 | 0.00 | 0.39 |

Table 8: Simulation results for the MW estimand of the modified WHO scale at day 14 under null treatment effect and covariates with prognostic power.

| Estimator | $n$ | Effect size | P(Reject $H_0$) | $n \times$ MSE | $n \times$ Var | \|Bias\| | Rel. eff. |
|---|---|---|---|---|---|---|---|
| CF-MARS | 100 | 0.50 | 0.11 | 0.54 | 0.54 | 0.00 | 1.95 |
| CF-RF | 100 | 0.50 | 0.07 | 0.23 | 0.23 | 0.00 | 0.84 |
| CF-XGBoost | 100 | 0.50 | 0.07 | 0.25 | 0.25 | 0.00 | 0.92 |
| $\ell_1$-LR | 100 | 0.50 | 0.07 | 0.22 | 0.22 | 0.00 | 0.80 |
| LR | 100 | 0.50 | 0.48 | 2.14 | 2.14 | 0.00 | 7.77 |
| MARS | 100 | 0.50 | 0.12 | 0.43 | 0.43 | 0.00 | 1.58 |
| RF | 100 | 0.50 | 0.13 | 0.16 | 0.16 | 0.00 | 0.57 |
| Unadjusted | 100 | 0.50 | 0.06 | 0.28 | 0.28 | 0.00 | 1.00 |
| XGBoost | 100 | 0.50 | 0.09 | 0.21 | 0.21 | 0.00 | 0.75 |
| CF-MARS | 500 | 0.50 | 0.06 | 0.28 | 0.28 | 0.00 | 1.07 |
| CF-RF | 500 | 0.50 | 0.05 | 0.20 | 0.20 | 0.00 | 0.74 |
| CF-XGBoost | 500 | 0.50 | 0.06 | 0.22 | 0.22 | 0.00 | 0.85 |
| $\ell_1$-LR | 500 | 0.50 | 0.05 | 0.21 | 0.21 | 0.00 | 0.80 |
| LR | 500 | 0.50 | 0.08 | 0.79 | 0.79 | 0.00 | 2.98 |
| MARS | 500 | 0.50 | 0.06 | 0.27 | 0.27 | 0.00 | 1.01 |
| RF | 500 | 0.50 | 0.14 | 0.13 | 0.13 | 0.00 | 0.48 |
| Unadjusted | 500 | 0.50 | 0.05 | 0.27 | 0.27 | 0.00 | 1.00 |
| XGBoost | 500 | 0.50 | 0.08 | 0.18 | 0.18 | 0.00 | 0.69 |
| CF-MARS | 1500 | 0.50 | 0.05 | 0.22 | 0.22 | 0.00 | 0.84 |
| CF-RF | 1500 | 0.50 | 0.05 | 0.13 | 0.13 | 0.00 | 0.52 |
| CF-XGBoost | 1500 | 0.50 | 0.06 | 0.18 | 0.18 | 0.00 | 0.71 |
| $\ell_1$-LR | 1500 | 0.50 | 0.05 | 0.21 | 0.21 | 0.00 | 0.81 |
| LR | 1500 | 0.50 | 0.06 | 0.22 | 0.22 | 0.00 | 0.84 |
| MARS | 1500 | 0.50 | 0.05 | 0.22 | 0.22 | 0.00 | 0.84 |
| RF | 1500 | 0.50 | 0.16 | 0.20 | 0.20 | 0.00 | 0.78 |
| Unadjusted | 1500 | 0.50 | 0.05 | 0.26 | 0.26 | 0.00 | 1.00 |
| XGBoost | 1500 | 0.50 | 0.15 | 0.11 | 0.11 | 0.00 | 0.43 |

## D.3 Results for simulations with a positive effect and where the covariates are not prognostic of the outcome

Table 9: Simulation results for the RMST of the time to intubation or death at day 14 under a positive effect and covariates with no prognostic power.

| Estimator | $n$ | Effect size | P(Reject $H_0$) | $n \times$ MSE | $n \times$ Var | |Bias| | Rel. eff. |
|---|---|---|---|---|---|---|---|
| CF-MARS | 100 | 1.04 | 0.25 | 94.58 | 94.56 | 0.01 | 1.46 |
| CF-RF | 100 | 1.04 | 0.21 | 77.13 | 77.13 | 0.00 | 1.19 |
| CF-XGBoost | 100 | 1.04 | 0.23 | 77.43 | 77.41 | 0.01 | 1.20 |
| $\ell_1$-LR | 100 | 1.04 | 0.24 | 64.21 | 64.21 | 0.01 | 0.99 |
| LR | 100 | 1.04 | 0.30 | 75.60 | 75.55 | 0.02 | 1.17 |
| MARS | 100 | 1.04 | 0.25 | 73.38 | 72.61 | 0.09 | 1.14 |
| RF | 100 | 1.04 | 0.58 | 68.20 | 64.12 | 0.20 | 1.06 |
| Unadjusted | 100 | 1.04 | 0.25 | 64.64 | 64.64 | 0.01 | 1.00 |
| XGBoost | 100 | 1.04 | 0.26 | 66.74 | 66.74 | 0.01 | 1.03 |
| CF-MARS | 500 | 1.04 | 0.81 | 67.36 | 67.36 | 0.00 | 1.02 |
| CF-RF | 500 | 1.04 | 0.77 | 74.08 | 74.06 | 0.01 | 1.12 |
| CF-XGBoost | 500 | 1.04 | 0.77 | 73.22 | 73.21 | 0.00 | 1.11 |
| $\ell_1$-LR | 500 | 1.04 | 0.83 | 66.61 | 66.61 | 0.00 | 1.01 |
| LR | 500 | 1.04 | 0.83 | 68.06 | 68.04 | 0.01 | 1.03 |
| MARS | 500 | 1.04 | 0.82 | 64.56 | 64.54 | 0.01 | 0.98 |
| RF | 500 | 1.04 | 0.91 | 81.56 | 79.70 | 0.06 | 1.24 |
| Unadjusted | 500 | 1.04 | 0.82 | 65.85 | 65.85 | 0.00 | 1.00 |
| XGBoost | 500 | 1.04 | 0.83 | 63.75 | 63.70 | 0.01 | 0.97 |
| CF-MARS | 1500 | 1.04 | 1.00 | 65.62 | 65.61 | 0.00 | 1.01 |
| CF-RF | 1500 | 1.04 | 1.00 | 73.92 | 73.91 | 0.00 | 1.14 |
| CF-XGBoost | 1500 | 1.04 | 1.00 | 70.56 | 70.55 | 0.00 | 1.08 |
| $\ell_1$-LR | 1500 | 1.04 | 1.00 | 65.55 | 65.54 | 0.00 | 1.01 |
| LR | 1500 | 1.04 | 1.00 | 65.55 | 65.55 | 0.00 | 1.01 |
| MARS | 1500 | 1.04 | 1.00 | 65.68 | 65.67 | 0.00 | 1.01 |
| RF | 1500 | 1.04 | 1.00 | 74.73 | 74.72 | 0.00 | 1.15 |
| Unadjusted | 1500 | 1.04 | 1.00 | 65.07 | 65.06 | 0.00 | 1.00 |
| XGBoost | 1500 | 1.04 | 1.00 | 66.71 | 66.69 | 0.00 | 1.03 |

Table 10: Simulation results for the RD of the time to intubation or death at day 7 under a positive effect and covariates with no prognostic power.

| Estimator | $n$ | Effect size | P(Reject $H_0$) | $n \times$ MSE | $n \times$ Var | \|Bias\| | Rel. eff. |
|---|---|---|---|---|---|---|---|
| CF-MARS | 100 | 0.10 | 0.26 | 0.83 | 0.83 | 0.00 | 1.39 |
| CF-RF | 100 | 0.10 | 0.22 | 0.68 | 0.68 | 0.00 | 1.14 |
| CF-XGBoost | 100 | 0.10 | 0.26 | 0.70 | 0.70 | 0.00 | 1.17 |
| $\ell_1$-LR | 100 | 0.10 | 0.26 | 0.60 | 0.60 | 0.00 | 1.00 |
| LR | 100 | 0.10 | 0.30 | 0.67 | 0.67 | 0.00 | 1.13 |
| MARS | 100 | 0.10 | 0.27 | 0.67 | 0.67 | 0.01 | 1.13 |
| RF | 100 | 0.10 | 0.51 | 0.71 | 0.68 | 0.02 | 1.18 |
| Unadjusted | 100 | 0.10 | 0.26 | 0.60 | 0.60 | 0.00 | 1.00 |
| XGBoost | 100 | 0.10 | 0.27 | 0.62 | 0.62 | 0.00 | 1.03 |
| CF-MARS | 500 | 0.10 | 0.82 | 0.61 | 0.61 | 0.00 | 1.00 |
| CF-RF | 500 | 0.10 | 0.79 | 0.66 | 0.66 | 0.00 | 1.08 |
| CF-XGBoost | 500 | 0.10 | 0.79 | 0.66 | 0.66 | 0.00 | 1.07 |
| $\ell_1$-LR | 500 | 0.10 | 0.83 | 0.61 | 0.61 | 0.00 | 0.99 |
| LR | 500 | 0.10 | 0.83 | 0.62 | 0.62 | 0.00 | 1.02 |
| MARS | 500 | 0.10 | 0.83 | 0.59 | 0.59 | 0.00 | 0.96 |
| RF | 500 | 0.10 | 0.89 | 0.79 | 0.77 | 0.01 | 1.28 |
| Unadjusted | 500 | 0.10 | 0.83 | 0.61 | 0.61 | 0.00 | 1.00 |
| XGBoost | 500 | 0.10 | 0.83 | 0.60 | 0.59 | 0.00 | 0.97 |
| CF-MARS | 1500 | 0.10 | 1.00 | 0.61 | 0.61 | 0.00 | 1.01 |
| CF-RF | 1500 | 0.10 | 1.00 | 0.67 | 0.67 | 0.00 | 1.11 |
| CF-XGBoost | 1500 | 0.10 | 1.00 | 0.63 | 0.63 | 0.00 | 1.05 |
| $\ell_1$-LR | 1500 | 0.10 | 1.00 | 0.60 | 0.60 | 0.00 | 0.99 |
| LR | 1500 | 0.10 | 1.00 | 0.60 | 0.60 | 0.00 | 0.99 |
| MARS | 1500 | 0.10 | 1.00 | 0.60 | 0.60 | 0.00 | 0.99 |
| RF | 1500 | 0.10 | 1.00 | 0.69 | 0.69 | 0.00 | 1.14 |
| Unadjusted | 1500 | 0.10 | 1.00 | 0.60 | 0.60 | 0.00 | 1.00 |
| XGBoost | 1500 | 0.10 | 1.00 | 0.61 | 0.61 | 0.00 | 1.01 |

Table 11: Simulation results for the LOR of the modified WHO scale at day 14 under a positive effect and covariates with no prognostic power.

| Estimator | $n$ | Effect size | P(Reject $H_0$) | $n \times$ MSE | $n \times$ Var | \|Bias\| | Rel. eff. |
|---|---|---|---|---|---|---|---|
| CF-MARS | 100 | 0.60 | 0.35 | 67.15 | 64.59 | 0.16 | 0.12 |
| CF-RF | 100 | 0.60 | 0.37 | 62.46 | 58.98 | 0.19 | 0.11 |
| CF-XGBoost | 100 | 0.60 | 0.40 | 76.06 | 72.37 | 0.19 | 0.13 |
| $\ell_1$-LR | 100 | 0.60 | 0.44 | 95.99 | 90.49 | 0.23 | 0.16 |
| LR | 100 | 0.60 | 0.41 | 392.58 | 330.71 | 0.79 | 0.67 |
| MARS | 100 | 0.60 | 0.41 | 187.85 | 167.51 | 0.45 | 0.32 |
| RF | 100 | 0.60 | 0.45 | 61.04 | 59.51 | 0.12 | 0.10 |
| Unadjusted | 100 | 0.60 | 0.42 | 582.27 | 441.29 | 1.19 | 1.00 |
| XGBoost | 100 | 0.60 | 0.43 | 132.91 | 124.69 | 0.29 | 0.23 |
| CF-MARS | 500 | 0.60 | 0.87 | 31.28 | 31.16 | 0.02 | 0.56 |
| CF-RF | 500 | 0.60 | 0.87 | 23.06 | 23.05 | 0.00 | 0.41 |
| CF-XGBoost | 500 | 0.60 | 0.88 | 25.83 | 25.62 | 0.02 | 0.46 |
| $\ell_1$-LR | 500 | 0.60 | 0.87 | 24.01 | 23.98 | 0.01 | 0.43 |
| LR | 500 | 0.60 | 0.86 | 53.78 | 53.45 | 0.03 | 0.96 |
| MARS | 500 | 0.60 | 0.87 | 39.20 | 39.04 | 0.02 | 0.70 |
| RF | 500 | 0.60 | 0.91 | 33.77 | 33.77 | 0.00 | 0.60 |
| Unadjusted | 500 | 0.60 | 0.86 | 56.01 | 55.76 | 0.02 | 1.00 |
| XGBoost | 500 | 0.60 | 0.89 | 22.54 | 22.54 | 0.00 | 0.40 |
| CF-MARS | 1500 | 0.60 | 1.00 | 20.98 | 20.98 | 0.00 | 1.04 |
| CF-RF | 1500 | 0.60 | 1.00 | 20.93 | 20.92 | 0.00 | 1.04 |
| CF-XGBoost | 1500 | 0.60 | 1.00 | 22.14 | 22.09 | 0.01 | 1.10 |
| $\ell_1$-LR | 1500 | 0.60 | 1.00 | 20.81 | 20.80 | 0.00 | 1.03 |
| LR | 1500 | 0.60 | 1.00 | 20.27 | 20.26 | 0.00 | 1.00 |
| MARS | 1500 | 0.60 | 1.00 | 20.16 | 20.15 | 0.00 | 1.00 |
| RF | 1500 | 0.60 | 1.00 | 20.51 | 20.48 | 0.00 | 1.02 |
| Unadjusted | 1500 | 0.60 | 1.00 | 20.20 | 20.20 | 0.00 | 1.00 |
| XGBoost | 1500 | 0.60 | 1.00 | 20.48 | 20.44 | 0.01 | 1.01 |

Table 12: Simulation results for the MW of the modified WHO scale at day 14 under a positive effect and covariates with no prognostic power.

| Estimator | $n$ | Effect size | P(Reject $H_0$) | $n \times$ MSE | $n \times$ Var | \|Bias\| | Rel. eff. |
|-----------|-----|-------------|-----------------|----------------|----------------|----------|-----------|
| CF-MARS | 100 | 0.46 | 0.18 | 0.54 | 0.54 | 0.00 | 2.05 |
| CF-RF | 100 | 0.46 | 0.15 | 0.29 | 0.29 | 0.00 | 1.11 |
| CF-XGBoost | 100 | 0.46 | 0.15 | 0.29 | 0.29 | 0.00 | 1.10 |
| $\ell_1$-LR | 100 | 0.46 | 0.16 | 0.26 | 0.26 | 0.00 | 1.00 |
| LR | 100 | 0.46 | 0.44 | 1.69 | 1.67 | 0.01 | 6.43 |
| MARS | 100 | 0.46 | 0.18 | 0.43 | 0.43 | 0.00 | 1.63 |
| RF | 100 | 0.46 | 0.19 | 0.24 | 0.23 | 0.01 | 0.91 |
| Unadjusted | 100 | 0.46 | 0.14 | 0.26 | 0.26 | 0.00 | 1.00 |
| XGBoost | 100 | 0.46 | 0.18 | 0.26 | 0.26 | 0.00 | 1.00 |
| CF-MARS | 500 | 0.46 | 0.44 | 0.50 | 0.49 | 0.00 | 1.93 |
| CF-RF | 500 | 0.46 | 0.44 | 0.28 | 0.28 | 0.00 | 1.09 |
| CF-XGBoost | 500 | 0.46 | 0.46 | 0.27 | 0.27 | 0.00 | 1.04 |
| $\ell_1$-LR | 500 | 0.46 | 0.46 | 0.26 | 0.26 | 0.00 | 1.02 |
| LR | 500 | 0.46 | 0.46 | 0.29 | 0.28 | 0.00 | 1.10 |
| MARS | 500 | 0.46 | 0.44 | 0.30 | 0.30 | 0.00 | 1.16 |
| RF | 500 | 0.46 | 0.51 | 0.24 | 0.24 | 0.00 | 0.93 |
| Unadjusted | 500 | 0.46 | 0.44 | 0.26 | 0.26 | 0.00 | 1.00 |
| XGBoost | 500 | 0.46 | 0.47 | 0.26 | 0.25 | 0.00 | 0.98 |
| CF-MARS | 1500 | 0.46 | 0.87 | 0.28 | 0.27 | 0.00 | 1.04 |
| CF-RF | 1500 | 0.46 | 0.88 | 0.27 | 0.27 | 0.00 | 1.03 |
| CF-XGBoost | 1500 | 0.46 | 0.88 | 0.28 | 0.27 | 0.00 | 1.06 |
| $\ell_1$-LR | 1500 | 0.46 | 0.88 | 0.27 | 0.26 | 0.00 | 1.00 |
| LR | 1500 | 0.46 | 0.89 | 0.26 | 0.25 | 0.00 | 0.98 |
| MARS | 1500 | 0.46 | 0.88 | 0.27 | 0.26 | 0.00 | 1.03 |
| RF | 1500 | 0.46 | 0.91 | 0.26 | 0.26 | 0.00 | 0.98 |
| Unadjusted | 1500 | 0.46 | 0.88 | 0.27 | 0.26 | 0.00 | 1.00 |
| XGBoost | 1500 | 0.46 | 0.89 | 0.27 | 0.26 | 0.00 | 1.01 |

## D.4 Results for simulations with null treatment effect and where the covariates are not prognostic of the outcome

Table 13: Simulation results for the RMST of the time to intubation or death at day 14 under null treatment effect and covariates with no prognostic power.

| Estimator | $n$ | Effect size | P(Reject $H_0$) | $n \times$ MSE | $n \times$ Var | \|Bias\| | Rel. eff. |
|---|---|---|---|---|---|---|---|
| CF-MARS | 100 | 0.00 | 0.11 | 142.00 | 135.82 | 0.25 | 1.62 |
| CF-RF | 100 | 0.00 | 0.05 | 103.14 | 103.11 | 0.02 | 1.17 |
| CF-XGBoost | 100 | 0.00 | 0.05 | 95.76 | 95.76 | 0.00 | 1.09 |
| $\ell_1$-LR | 100 | 0.00 | 0.06 | 87.21 | 87.13 | 0.03 | 0.99 |
| LR | 100 | 0.00 | 0.09 | 99.98 | 99.95 | 0.02 | 1.14 |
| MARS | 100 | 0.00 | 0.08 | 97.67 | 97.24 | 0.07 | 1.11 |
| RF | 100 | 0.00 | 0.31 | 70.72 | 70.68 | 0.02 | 0.80 |
| Unadjusted | 100 | 0.00 | 0.06 | 87.92 | 87.85 | 0.03 | 1.00 |
| XGBoost | 100 | 0.00 | 0.06 | 87.83 | 87.82 | 0.01 | 1.00 |
| CF-MARS | 500 | 0.00 | 0.06 | 91.82 | 91.53 | 0.02 | 1.08 |
| CF-RF | 500 | 0.00 | 0.05 | 98.39 | 98.38 | 0.00 | 1.16 |
| CF-XGBoost | 500 | 0.00 | 0.05 | 97.91 | 97.91 | 0.00 | 1.15 |
| $\ell_1$-LR | 500 | 0.00 | 0.05 | 87.33 | 87.32 | 0.01 | 1.03 |
| LR | 500 | 0.00 | 0.06 | 92.74 | 92.74 | 0.00 | 1.09 |
| MARS | 500 | 0.00 | 0.05 | 88.10 | 88.10 | 0.00 | 1.04 |
| RF | 500 | 0.00 | 0.17 | 86.34 | 86.34 | 0.00 | 1.01 |
| Unadjusted | 500 | 0.00 | 0.05 | 85.12 | 85.10 | 0.01 | 1.00 |
| XGBoost | 500 | 0.00 | 0.06 | 86.89 | 86.89 | 0.00 | 1.02 |
| CF-MARS | 1500 | 0.00 | 0.05 | 85.92 | 85.89 | 0.00 | 1.00 |
| CF-RF | 1500 | 0.00 | 0.05 | 94.58 | 94.57 | 0.00 | 1.10 |
| CF-XGBoost | 1500 | 0.00 | 0.04 | 93.27 | 93.26 | 0.00 | 1.09 |
| $\ell_1$-LR | 1500 | 0.00 | 0.05 | 85.62 | 85.61 | 0.00 | 1.00 |
| LR | 1500 | 0.00 | 0.06 | 90.04 | 90.01 | 0.00 | 1.05 |
| MARS | 1500 | 0.00 | 0.05 | 89.32 | 89.32 | 0.00 | 1.04 |
| RF | 1500 | 0.00 | 0.11 | 89.12 | 89.12 | 0.00 | 1.04 |
| Unadjusted | 1500 | 0.00 | 0.05 | 85.87 | 85.83 | 0.00 | 1.00 |
| XGBoost | 1500 | 0.00 | 0.07 | 89.61 | 89.59 | 0.00 | 1.04 |

Table 14: Simulation results for the RD of the time to intubation or death at day 7 under null treatment effect and covariates with no prognostic power.

| Estimator | $n$ | Effect size | P(Reject $H_0$) | $n \times$ MSE | $n \times$ Var | \|Bias\| | Rel. eff. |
|---|---|---|---|---|---|---|---|
| CF-MARS | 100 | 0.00 | 0.10 | 1.11 | 1.07 | 0.02 | 1.52 |
| CF-RF | 100 | 0.00 | 0.05 | 0.85 | 0.85 | 0.00 | 1.17 |
| CF-XGBoost | 100 | 0.00 | 0.05 | 0.80 | 0.80 | 0.00 | 1.09 |
| $\ell_1$-LR | 100 | 0.00 | 0.06 | 0.72 | 0.72 | 0.00 | 0.99 |
| LR | 100 | 0.00 | 0.08 | 0.82 | 0.82 | 0.00 | 1.12 |
| MARS | 100 | 0.00 | 0.07 | 0.79 | 0.79 | 0.00 | 1.09 |
| RF | 100 | 0.00 | 0.27 | 0.66 | 0.66 | 0.00 | 0.90 |
| Unadjusted | 100 | 0.00 | 0.06 | 0.73 | 0.73 | 0.00 | 1.00 |
| XGBoost | 100 | 0.00 | 0.06 | 0.73 | 0.73 | 0.00 | 1.00 |
| CF-MARS | 500 | 0.00 | 0.05 | 0.76 | 0.76 | 0.00 | 1.06 |
| CF-RF | 500 | 0.00 | 0.05 | 0.82 | 0.82 | 0.00 | 1.13 |
| CF-XGBoost | 500 | 0.00 | 0.05 | 0.80 | 0.80 | 0.00 | 1.11 |
| $\ell_1$-LR | 500 | 0.00 | 0.05 | 0.72 | 0.72 | 0.00 | 1.00 |
| LR | 500 | 0.00 | 0.06 | 0.78 | 0.78 | 0.00 | 1.08 |
| MARS | 500 | 0.00 | 0.05 | 0.74 | 0.74 | 0.00 | 1.03 |
| RF | 500 | 0.00 | 0.16 | 0.73 | 0.73 | 0.00 | 1.02 |
| Unadjusted | 500 | 0.00 | 0.05 | 0.72 | 0.72 | 0.00 | 1.00 |
| XGBoost | 500 | 0.00 | 0.06 | 0.72 | 0.72 | 0.00 | 1.00 |
| CF-MARS | 1500 | 0.00 | 0.05 | 0.71 | 0.71 | 0.00 | 0.99 |
| CF-RF | 1500 | 0.00 | 0.05 | 0.79 | 0.79 | 0.00 | 1.09 |
| CF-XGBoost | 1500 | 0.00 | 0.04 | 0.77 | 0.77 | 0.00 | 1.07 |
| $\ell_1$-LR | 1500 | 0.00 | 0.05 | 0.72 | 0.72 | 0.00 | 0.99 |
| LR | 1500 | 0.00 | 0.06 | 0.74 | 0.74 | 0.00 | 1.03 |
| MARS | 1500 | 0.00 | 0.05 | 0.74 | 0.74 | 0.00 | 1.03 |
| RF | 1500 | 0.00 | 0.10 | 0.74 | 0.74 | 0.00 | 1.03 |
| Unadjusted | 1500 | 0.00 | 0.05 | 0.72 | 0.72 | 0.00 | 1.00 |
| XGBoost | 1500 | 0.00 | 0.06 | 0.75 | 0.75 | 0.00 | 1.05 |

Table 15: Simulation results for the LOR of the modified WHO scale at day 14 under null treatment effect and covariates with no prognostic power.

| Estimator | $n$ | Effect size | P(Reject $H_0$) | $n \times$ MSE | $n \times$ Var | \|Bias\| | Rel. eff. |
|---|---|---|---|---|---|---|---|
| CF-MARS | 100 | 0.00 | 0.07 | 38.32 | 38.30 | 0.01 | 0.99 |
| CF-RF | 100 | 0.00 | 0.06 | 20.08 | 20.07 | 0.01 | 0.52 |
| CF-XGBoost | 100 | 0.00 | 0.07 | 21.21 | 21.21 | 0.00 | 0.55 |
| $\ell_1$-LR | 100 | 0.00 | 0.09 | 22.15 | 22.15 | 0.01 | 0.57 |
| LR | 100 | 0.00 | 0.08 | 96.02 | 96.02 | 0.00 | 2.49 |
| MARS | 100 | 0.00 | 0.07 | 28.91 | 28.90 | 0.01 | 0.75 |
| RF | 100 | 0.00 | 0.11 | 16.61 | 16.60 | 0.01 | 0.43 |
| Unadjusted | 100 | 0.00 | 0.06 | 38.60 | 38.60 | 0.00 | 1.00 |
| XGBoost | 100 | 0.00 | 0.10 | 34.50 | 34.48 | 0.01 | 0.89 |
| CF-MARS | 500 | 0.00 | 0.05 | 16.92 | 16.91 | 0.00 | 1.01 |
| CF-RF | 500 | 0.00 | 0.05 | 17.27 | 17.27 | 0.00 | 1.03 |
| CF-XGBoost | 500 | 0.00 | 0.06 | 17.75 | 17.74 | 0.00 | 1.06 |
| $\ell_1$-LR | 500 | 0.00 | 0.06 | 16.45 | 16.45 | 0.00 | 0.98 |
| LR | 500 | 0.00 | 0.06 | 17.28 | 17.28 | 0.00 | 1.03 |
| MARS | 500 | 0.00 | 0.06 | 17.26 | 17.26 | 0.00 | 1.03 |
| RF | 500 | 0.00 | 0.09 | 15.20 | 15.20 | 0.00 | 0.91 |
| Unadjusted | 500 | 0.00 | 0.05 | 16.74 | 16.74 | 0.00 | 1.00 |
| XGBoost | 500 | 0.00 | 0.07 | 16.54 | 16.54 | 0.00 | 0.99 |
| CF-MARS | 1500 | 0.00 | 0.05 | 16.60 | 16.60 | 0.00 | 1.02 |
| CF-RF | 1500 | 0.00 | 0.06 | 17.37 | 17.37 | 0.00 | 1.07 |
| CF-XGBoost | 1500 | 0.00 | 0.06 | 17.11 | 17.11 | 0.00 | 1.05 |
| $\ell_1$-LR | 1500 | 0.00 | 0.05 | 16.33 | 16.28 | 0.01 | 1.00 |
| LR | 1500 | 0.00 | 0.05 | 16.61 | 16.61 | 0.00 | 1.02 |
| MARS | 1500 | 0.00 | 0.05 | 16.25 | 16.24 | 0.00 | 1.00 |
| RF | 1500 | 0.00 | 0.08 | 16.21 | 16.21 | 0.00 | 0.99 |
| Unadjusted | 1500 | 0.00 | 0.05 | 16.30 | 16.29 | 0.00 | 1.00 |
| XGBoost | 1500 | 0.00 | 0.06 | 16.30 | 16.29 | 0.00 | 1.00 |

Table 16: Simulation results for the MW of the modified WHO scale at day 14 under null treatment effect and covariates with no prognostic power.

| Estimator | $n$ | Effect size | P(Reject $H_0$) | $n \times$ MSE | $n \times$ Var | \|Bias\| | Rel. eff. |
|---|---|---|---|---|---|---|---|
| CF-MARS | 100 | 0.50 | 0.08 | 0.42 | 0.42 | 0.00 | 1.54 |
| CF-RF | 100 | 0.50 | 0.06 | 0.28 | 0.28 | 0.00 | 1.02 |
| CF-XGBoost | 100 | 0.50 | 0.07 | 0.29 | 0.29 | 0.00 | 1.07 |
| $\ell_1$-LR | 100 | 0.50 | 0.07 | 0.27 | 0.27 | 0.00 | 0.99 |
| LR | 100 | 0.50 | 0.47 | 2.09 | 2.09 | 0.00 | 7.57 |
| MARS | 100 | 0.50 | 0.08 | 0.36 | 0.36 | 0.00 | 1.32 |
| RF | 100 | 0.50 | 0.11 | 0.22 | 0.22 | 0.00 | 0.79 |
| Unadjusted | 100 | 0.50 | 0.06 | 0.28 | 0.28 | 0.00 | 1.00 |
| XGBoost | 100 | 0.50 | 0.09 | 0.27 | 0.27 | 0.00 | 0.97 |
| CF-MARS | 500 | 0.50 | 0.05 | 0.28 | 0.28 | 0.00 | 1.04 |
| CF-RF | 500 | 0.50 | 0.06 | 0.28 | 0.28 | 0.00 | 1.04 |
| CF-XGBoost | 500 | 0.50 | 0.06 | 0.29 | 0.29 | 0.00 | 1.07 |
| $\ell_1$-LR | 500 | 0.50 | 0.05 | 0.26 | 0.26 | 0.00 | 0.99 |
| LR | 500 | 0.50 | 0.08 | 0.72 | 0.72 | 0.00 | 2.69 |
| MARS | 500 | 0.50 | 0.06 | 0.28 | 0.28 | 0.00 | 1.05 |
| RF | 500 | 0.50 | 0.09 | 0.24 | 0.24 | 0.00 | 0.91 |
| Unadjusted | 500 | 0.50 | 0.05 | 0.27 | 0.27 | 0.00 | 1.00 |
| XGBoost | 500 | 0.50 | 0.07 | 0.26 | 0.26 | 0.00 | 0.99 |
| CF-MARS | 1500 | 0.50 | 0.05 | 0.27 | 0.27 | 0.00 | 1.03 |
| CF-RF | 1500 | 0.50 | 0.06 | 0.28 | 0.28 | 0.00 | 1.08 |
| CF-XGBoost | 1500 | 0.50 | 0.06 | 0.28 | 0.28 | 0.00 | 1.06 |
| $\ell_1$-LR | 1500 | 0.50 | 0.05 | 0.26 | 0.26 | 0.00 | 1.02 |
| LR | 1500 | 0.50 | 0.05 | 0.27 | 0.27 | 0.00 | 1.03 |
| MARS | 1500 | 0.50 | 0.05 | 0.26 | 0.26 | 0.00 | 1.00 |
| RF | 1500 | 0.50 | 0.08 | 0.26 | 0.26 | 0.00 | 1.00 |
| Unadjusted | 1500 | 0.50 | 0.05 | 0.26 | 0.26 | 0.00 | 1.00 |
| XGBoost | 1500 | 0.50 | 0.06 | 0.27 | 0.27 | 0.00 | 1.02 |

# References

P.J. Bickel, C.A.J. Klaassen, Y. Ritov, and J. Wellner. *Efficient and Adaptive Estimation for Semiparametric Models*. Springer-Verlag, 1997.

Iván Díaz, Oleksandr Savenkov, and Karla Ballman. Targeted learning ensembles for optimal individualized treatment rules with time-to-event outcomes. *Biometrika*, 105(3):723–738, 2018.

Iván Díaz, Elizabeth Colantuoni, Daniel F. Hanley, and Michael Rosenblum. Improved precision in the analysis of randomized trials with survival outcomes, without assuming proportional hazards. *Lifetime Data Analysis*, 25(3):439–468, 2019.

Kelly L. Moore and Mark J. van der Laan. Increasing power in randomized trials with right censored outcomes through covariate adjustment. *Journal of Biopharmaceutical Statistics*, 19(6): 1099–1131, 2009.

A. W. van der Vaart. *Asymptotic Statistics.* Cambridge University Press, 1998.

Aad W. van der Vaart and Jon A. Wellner. *Weak Convergence and Emprical Processes.* Springer-Verlag New York, 1996.