

ASHLEY I. NAIMI, PHD

DOUBLE ROBUST ESTI- MATION

Outline

- Double Robust Estimation and Machine Learning
- Observed Data and Target Parameter
- Parametric Estimation
 - Estimation via Parametric Models
 - Estimation via Parametric Exposure Models
 - Parametric Doubly Robust Estimation
- Nonparametric Singly Robust Estimation: The Curse of Dimensionality
- Nonparametric Doubly Robust Estimation
- Simulation Study
- The AIPW Package

Both machine learning methods and doubly robust estimators are becoming increasingly popular, yet the critical relation between them remains poorly understood. Machine learning methods consist of a wide range of analytic techniques that do not require hard to verify modeling assumptions. Because of this, they are often assumed to be less biased than their standard parametric counterparts. This perceived property has motivated many to either recommended or use machine learning methods to estimate statistical parameters that correspond to causal quantities of interest (Lee, Lessler, and Stuart 2010; Westreich, Lessler, and Funk 2010; Snowden, Rose, and Mortimer 2011; Oulhote et al. 2019) However, it is generally not recognized that machine learning methods are subject to problems that arise from the curse of dimensionality, a term first coined by Bellman (1957) to refer to a set of problems encountered when estimating models with many variables (Wasserman 2006).

Doubly robust estimators are so named because these methods allow two chances for adjustment (J. Robins and Rotnitzky 1995, 2001; Bang and Robins 2005) In the case of confounding adjustment, these chances arise because the analyst must fit two models: a model for the outcome conditional on the exposure and all confounders (outcome model); and a model for the exposure conditional all confounders (the propensity score model). These are then combined to estimate the effect of interest (Rotnitzky and Vansteelandt 2014).

The benefits of doubly robust methods have been explained by pointing out that if a confounding variable is left out of either the exposure or the outcome model (but not both), unbiased estimates can still be obtained (Jonsson-Funk et al. 2011). While true, analysts would not typically leave confounding variables out of either the exposure or outcome model. Such justifications ignore a critically important benefit conferred by doubly robust estimators: under relatively mild conditions, they remain unbiased, with asymptotically nominal confidence interval coverage, even when machine learning methods are used to fit the exposure and outcome models (Mark J. van der Laan and Rubin 2006; Edward H. Kennedy and Balakrishnan 2017). In effect, doubly robust methods can mitigate or resolve

problems caused by the curse of dimensionality.

This little recognized relation between machine learning and doubly robust estimators has important implications for applied researchers, particularly those interested in using machine learning methods to estimate causal effects. Here, we examine these implications using simple Monte Carlo simulations (Metropolis and Ulam 1949). Our intent is to clarify that machine learning methods should be used with doubly robust methods; they should not generally be used to estimate causal effects with singly robust techniques, such as model-based standardization (i.e., the parametric g formula, or g computation), or inverse probability weighting.

Observed Data & Target Parameter

We consider a simple setting with a single binary exposure (X), a set of continuous confounders ($\mathbf{C} = \{C_1, C_2, C_3, C_4\}$) measured at baseline, and a single continuous outcome (Y) measured at the end of follow-up. In an observational cohort study to estimate the effect of X on Y , \mathbf{C} might be assumed a minimally sufficient adjustment set (Greenland, Pearl, and Robins 1999), and the outcome and exposure would be assumed generated according to some unknown models, for example:

$$E(Y \mid X, \mathbf{C}) = g(X, \mathbf{C}), \quad (\text{Model 1})$$

$$P(X = 1 \mid \mathbf{C}) = f(\mathbf{C}). \quad (\text{Model 2})$$

In the above equations, we use $g(\bullet)$ and $f(\bullet)$ to emphasize that the expected outcome conditional on X and \mathbf{C} , and the probability of the exposure given \mathbf{C} need not be considered standard linear or logistic regression functions. Rather, $g(\bullet)$ and $f(\bullet)$ represent arbitrary functions relating the exposure and confounders to the outcome, and the confounders to the exposure. Importantly, in an observational cohort study assuming a correct confounder adjustment set, these arbitrary functions usually represent the extent of what is known about the exposure and outcome models (J. Robins 2001). That is, while these models may typically be assumed to be in the family of generalized

linear models (Nelder and Wedderburn 1972), we note below why this may not often be ideal.

We focus here on the average treatment effect:

$$\psi = E(Y^{x=1} - Y^{x=0})$$

where Y^x is the outcome that would be observed if X were set to x . This estimand is (point) identified under positivity, consistency, and exchangeability (James M. Robins and Hernán 2009; A. I. Naimi, Cole, and Kennedy 2017). If these assumptions hold, ψ can be estimated using a number of approaches. In the equations that follow, we let i index sample observations which range from 1 to N , $\hat{g}_i(X = x, \mathbf{C})$ and $\hat{f}_i(\mathbf{C})$ are individual sample predictions for $E(Y | X = x, \mathbf{C})$ and $P(X = 1 | \mathbf{C})$, respectively.

With predictions from Model 1, ψ can be estimated via model-based standardization (henceforth g computation) (A. I. Naimi, Cole, and Kennedy 2017):

$$\hat{\psi}_{gComp} = \frac{1}{N} \sum_{i=1}^N \{ \hat{g}_i(X = 1, \mathbf{C}) - \hat{g}_i(X = 0, \mathbf{C}) \}. \quad (1)$$

With predictions from Model 2, ψ can be estimated via inverse probability weighting (Hernán and Robins 2006) as:

$$\hat{\psi}_{ipw} = \frac{1}{N} \sum_{i=1}^N \left\{ \left[\frac{X_i Y_i}{\hat{f}_i(\mathbf{C})} \right] - \left[\frac{(1 - X_i) Y_i}{1 - \hat{f}_i(\mathbf{C})} \right] \right\}. \quad (2)$$

Both approaches 1 and 2 are “singly robust” in that they typically rely entirely on the correct specification of the appropriate single regression model. If these models are misspecified, the estimators will not generally converge to the true value.

Alternatively, one may employ a “doubly robust” technique where predictions from both the exposure and outcome models are combined into a single estimator to quantify the effect of interest. For example, using predictions from both Models Model 2 and Model 1,

ψ can be estimated as:

$$\hat{\psi}_{aipw} = \frac{1}{N} \sum_{i=1}^N \left\{ \frac{(2X_i - 1)[Y_i - \hat{g}_i(X, \mathbf{C})]}{(2X_i - 1)\hat{f}_i(\mathbf{C}) + (1 - X_i)} + \hat{g}_i(X = 1, \mathbf{C}) - \hat{g}_i(X = 0, \mathbf{C}) \right\}. \quad (3)$$

Equation 3 is an augmented inverse probability weighted estimator, and will converge to the true value as the sample size grows if either $f(\mathbf{C})$ or $g(X, \mathbf{C})$, but not necessarily both, are consistently estimated. The estimator 3 can be viewed as either a bias-corrected version of the g computation estimator (where the correction is the term incorporating the propensity score defined in Model 2), or an efficiency enhanced version of the IPW estimator (where the enhancement is the term incorporating the outcome model defined in Model 1) (Daniel 2018).

There is a recently developed R package that can be used to implement AIPW (Zhong et al. 2021). Among other things, this package enables the use of the Super Learner to estimate the exposure and outcome models, and implements the sample splitting procedure described below, which is required for optimal performance of any DR estimator using machine learning. Details on the package and its implementation are available in the citation above, as well as here: <https://yqzhong7.github.io/AIPW/>.

Here is an example set of code one can use to implement the AIPW estimator in the package:

```
library(AIPW)
library(SuperLearner)
set.seed(1234)
#load simulated dataset (RCT)
data(eager_sim_rct)
#Specify SuperLearner libraries
sl.lib = c("SL.gam", "SL.earth", "SL.ranger", "SL.xgboost")
#Create a vector of covariates
Cov = c("loss_num", "age", "time_try_pregnant", "BMI", "meanAP")
#create a new AIPW object called AIPW_SL
AIPW_SL <- AIPW$new(Y = eager_sim_rct$sim_Y,
                    A = eager_sim_rct$sim_Tx,
```

```

W.g = eager_sim_rct$eligibility,
W.Q = subset(eager_sim_rct,select=Cov), #covariates
Q.SL.library = sl.lib, #outcome model
g.SL.library = sl.lib, #exposure model
k_split = 10, #num of folds for cross-fitting
verbose=TRUE)

#fit the data stored in the AIPW_SL object
AIPW_SL$fit()

#summarise the results using truncated propensity scores
AIPW_SL$summary(g.bound = 0.025)

```

Alternatively, Model 2 can be used to “update” Model 1 via targeted minimum loss-based estimation: (Rose and Laan 2011, (p72–3))

$$\hat{\psi}_{tmle} = \frac{1}{N} \sum_{i=1}^N \{ \hat{g}_i^u(X=1, \mathbf{C}) - \hat{g}_i^u(X=0, \mathbf{C}) \}, \quad (4)$$

where $\hat{g}_i^u(X=x, \mathbf{C})$ are predictions from an “updated” outcome model. For the average treatment effect, this outcome model is updated by first generating a modified inverse probability weight, defined as:

$$H(X, \mathbf{C}) = \begin{cases} \frac{1}{\hat{f}_i(\mathbf{C})} & \text{if } X = 1 \\ -\frac{1}{1-\hat{f}_i(\mathbf{C})} & \text{otherwise} \end{cases}$$

and then including this inverse probability weight in a no-intercept logistic regression model for the outcome that includes the previous outcome predictions $\hat{g}_i(X, \mathbf{C})$ as an offset. The $\hat{g}_i^u(X=x, \mathbf{C})$ predictions are then generated from this model by setting X to 1 and then to 0 for all individuals in the sample. TMLE is asymptotically equivalent to equation 3 but can have better finite-sample performance (Gruber and Laan 2012).

Parametric Estimation

For continuous Y and binary X , it is customary to specify models Model 1 and Model 2 parametrically using linear and logistic regression, respectively. Doing so effectively states that we know enough

about the form of $g(X, \mathbf{C})$ and $f(\mathbf{C})$ to define them as:

$$g(X, \mathbf{C}) = E(Y \mid X, \mathbf{C}) = \beta_0 + \beta_1 X + \beta_2 C_1 + \beta_3 C_2 + \beta_4 C_3 + \beta_5 C_4, \quad (5)$$

$$Y \mid X, \mathbf{C} \sim \mathcal{N}(E(Y \mid X, \mathbf{C}), \sigma^2)$$

$$f(\mathbf{C}) = P(X = 1 \mid \mathbf{C}) = \text{expit}(\alpha_0 + \alpha_1 C_1 + \alpha_2 C_2 + \alpha_3 C_3 + \alpha_4 C_4), \quad (6)$$

$$\text{expit}(\bullet) = 1 / (1 + \exp[-\bullet])$$

Imposing these forms on $g(X, \mathbf{C})$ and $f(\mathbf{C})$ permits use of maximum likelihood for estimation and inference (Cole, Chu, and Greenland 2013).

Estimation via Parametric Models

Equation 5 imposes several parametric constraints on the form of $g(X, \mathbf{C})$: (i) Y follows a conditional normal distribution with constant variance not depending on X or \mathbf{C} ; and (ii) the conditional mean of Y is related to the covariates X and \mathbf{C} additively, as defined in equation 5. If these constraints on $g(X, \mathbf{C})$ are true, and other identification and regularity conditions hold (Longford 2008, ch2), the maximum likelihood estimates of β are asymptotically efficient (Rencher 2000, (p144)). Relatedly, under the model constraints and identification and regularity conditions, as the sample size increases, the estimates of $g(X, \mathbf{C})$ and/or $f(\mathbf{C})$ will converge to the true values at an optimal (i.e., \sqrt{N}) rate, and their distribution will be such that confidence intervals can be easily derived.

If constraint (i) is violated, the maximum likelihood estimator is no longer the most efficient, but can still be used to estimate ψ consistently. If constraint (ii) is violated, then the maximum likelihood estimator is no longer consistent. Depending on the severity to which constraint (ii) is violated, the bias may be substantial. Unfortunately, in an observational study the true form of equation 5 is almost never known. This means that such maximum likelihood estimates are almost always biased, with the degree of bias depending on the (unknown) extent to which the model is mis-specified (Box 1976).

Estimation via Parametric Exposure Model

One way to avoid relying on correct outcome model specification is to use a parametric approach for Model 2, and estimate ψ via $\hat{\psi}_{ipw}$. Specifically, with IP-weighting, one need not model the interactions between the exposure and any covariates (Hernán, Brumback, and Robins 2001). Such an estimator is not as efficient as $\hat{\psi}_{gComp}$, and can be subject to important finite-sample biases when weights are very large, or when there are no observations to weight in certain exposure-confounder strata. But as the sample size increases, the inverse probability weighted estimator converges at the same standard \sqrt{N} rate as the g computation estimator (Westreich et al. 2012). Unfortunately, as with the outcome model, the true form of Model 2 will almost never be known in an observational study. Mis-specification of equation 6 will also lead to biased estimation of ψ , again with the degree of bias depending on the unknown extent of model mis-specification.

Parametric Doubly Robust Estimation

To mitigate against mis-specification of the exposure or outcome models, numerous authors have advocated for the use of estimators such as equations 3 or 4. These doubly robust estimators remain consistent even if either the exposure model or the outcome model is mis-specified, but not both. However, if it is unlikely that either equations 5 or 6 is correct, then the doubly robust estimator will also likely be biased, and not much better than the singly robust estimators (Kang and L. 2007; Edward H. Kennedy and Balakrishnan 2017).

Nonparametric Singly Robust Estimation: The Curse of Dimensionality

Nonparametric methods are an alternative to parametric models. For example, nonparametric maximum likelihood estimation (NPMLE) for Model 2 or Model 1 would entail fitting equations 5 or 6, but

with a parameter for each unique combination of values defined by the cross-classification of all covariates (i.e., saturating the model). However, the NPMLE will be undefined in any finite sample with a continuous confounder, since there will be no covariate patterns containing both treated and untreated subjects.

Alternatively, one can use nonparametric “machine learning” methods like kernel regression, splines, random forests, boosting, etc., which exploit smoothness across covariate patterns to estimate the regression function. However, for any nonparametric approach there is an explicit bias-variance trade-off that arises in the choice of tuning parameters; less smoothing yields smaller bias but larger variance, while more smoothing yields smaller variance but larger bias (parametric models can be viewed as an extreme form of smoothing). This tradeoff has important consequences. In particular, it is generally impossible to estimate regression functions nonparametrically at the standard \sqrt{N} rates attained by correctly specified parametric estimators (Vaart 2000). These slow rates generally require sample sizes that are exponentially larger than those required for (fast converging) parametric methods to maintain the same degree of accuracy.

Convergence rates for nonparametric estimators become slower with more flexibility and more covariates. For example, a standard rate for estimating smooth regression functions is $N^{-\beta/(2\beta+d)}$, where β represents the number of derivatives of the true regression function, and d represents the dimension of, or number of covariates in, the true regression function. This issue is known as the curse of dimensionality (Györfi et al. 2002; J. M. Robins and Ritov 1997; Wasserman 2006). Sometimes this is viewed as a disadvantage of nonparametric methods; however, it is just the cost of making weaker assumptions: if a parametric model is misspecified, it will converge very quickly to the wrong answer.

In addition to slower convergence rates, confidence intervals are harder to obtain. Specifically, even in the rare case where one can derive asymptotic distributions for nonparametric estimators, it is typically not possible to construct confidence intervals (even via the bootstrap, as it requires certain convergence rate conditions to hold)

without impractically undersmoothing the regression function (i.e., overfitting the data) (Hahn 1998).

These complications (slow rates and lack of valid confidence intervals) are generally inherited by the singly robust estimators 2 and 1 (apart from a few special cases which require simple estimators, such as kernel methods with strong smoothness assumptions and careful tuning parameter choices that are suboptimal for estimating f or g). For general nonparametric estimators \hat{f} and \hat{g} , the estimators 2 and 1 will converge at slow rates, and honest confidence intervals (defined as confidence intervals that are at least nominal over a large nonparametric class of regression functions) (Li 1989) will not be computable.

Nonparametric Doubly Robust Estimation

Fortunately, doubly robust estimators that rely on nonparametric estimates of f and g do not suffer from the same limitations as the nonparametric versions of the singly robust estimators. In particular the doubly robust estimators 3 and 4 can be \sqrt{N} -consistent, asymptotically normal, and optimally efficient even if the estimators \hat{f} and \hat{g} are converging at slower nonparametric rates. In other words, the doubly robust estimator is less susceptible to the curse of dimensionality. This is because the singly robust estimators are combined in a way that their combined convergence rates are as fast or faster than the convergence rate of each estimator separately. In particular, if \hat{f} and \hat{g} are converging to their targets at least faster than $n^{-1/4}$ rates (technically, in L_2 norm), the doubly robust estimator will behave (asymptotically) just as if both f and g were estimated with correct parametric models. Importantly, $n^{-1/4}$ rates can be attained nonparametrically under relatively weak (smoothness, sparsity, or other structural) assumptions [Györfi et al. (2002); Wasserman2006]. This improved performance of nonparametric methods when used with doubly robust techniques has important implications for applied researchers.

Simulation Study

Data Generating Mechanism: Correct Specification

To explore these implications, we carried out a simulation study of singly and doubly robust estimators with parametric and non-parametric methods. We simulated 100 Monte Carlo samples, with sample sizes of {200, 1200, 5000} using data generating mechanisms that would lead to both simple and challenging conditions for estimation and inference. Specifically, we generated four independent standard normal confounders, denoted C . Both the exposure and outcome models included each of these confounders. The exposure was generated from a logistic model with:

$$P(X = 1 | C) = \text{expit} \{-1 + \log(1.75)C_1 + \log(1.75)C_2 + \log(1.75)C_3 + \log(1.75)C_4\}, \quad (7)$$

A continuous outcome was generated as:

$$Y = 120 + 6X + 3C_1 + 3C_2 + 3C_3 + 3C_4 + \epsilon, \quad (8)$$

where the true average treatment effect $\psi = 6$, with ϵ drawn from a normal distribution with mean $\mu = 0$ and standard deviation $\sigma = 6$.

Data Generating Mechanism: Model Misspecification

To induce model misspecification, we followed previous research (Kang and L. 2007) and transformed each of the continuous confounders as follows:

$$\begin{aligned} Z_1 &= \exp(C_1/2) \\ Z_2 &= C_2 / (1 + \exp(C_1)) + 10 \\ Z_3 &= (C_1 C_3 / 25 + 0.6)^3 \\ Z_4 &= (C_2 + C_4 + 20)^2 \end{aligned}$$

Thus, while the true models generating the exposure and outcome variables included only the untransformed variables C , analyses conducted under parametric model misspecification included only

the transformed variables Z .

Simulation Analysis

In each Monte Carlo sample, we estimated the average treatment effect $\psi = E(Y^1 - Y^0) = 6$ using g computation, inverse probability weighting, augmented inverse probability weighting, and targeted minimum loss-based estimation under two settings: (i) only the simple confounder data C were available and adjusted for in all estimators (parametric and nonparametric), and (ii) only the transformed confounder data Z were available adjusted for in all estimators (parametric and nonparametric).

Parametric estimation was accomplished via generalized linear models, with a binomial distribution and logistic link for the exposure, and a Gaussian distribution and identity link for the outcome. As described above, these parametric models are correctly specified when the simple confounders are used, but highly misspecified when the transformed confounders are used.

Nonparametric estimation was accomplished via a stacking algorithm (Super Learner) (Mark J. van der Laan, Polley, and Hubbard 2007). To explore the importance of the selected algorithm, we implemented a wide variety of different stacking algorithms that included different sets of base algorithms. Full details on all variations of the stacking algorithms explored are available in the GitHub Repository provided below. Here, we present the results based on a stacked generalizations that included:

version 1) (i) random forests with 500 trees, random subspace selection value of two, and a minimum node size of 30 and 60; (ii) the extreme gradient boosting algorithm with 500 trees, a maximum tree depth of 4, shrinkage parameter of 0.1, and minimum node size of 30 and 60.

version 2) Both random forests and extreme gradient boosting included in version 1, as well as (iii) generalized additive models with univariate smoothing splines with effective degrees of freedom between 3 and 8.

We also explored estimating the average treatment effects of interest with the stacking algorithms in version 2 that included 2-way interactions between all four confounders in the adjustment set. For all stacking algorithms, cross validation was used to compute the learner weights with fold sizes of $K = 10, 5$, and 5 for the sample sizes 200 , 1200 , and 5000 , respectively (A. Naimi, Platt, and JC 2018). For each machine learning based doubly robust estimator, we also explored the impact of sample splitting (Rinaldo et al. 2018; Zivich and Breskin 2020). This procedure involves splitting the sample into K equal size folds, fitting models for $f(\mathbf{C})$ and $g(X, \mathbf{C})$ in one fold, using these models to predict exposure and outcome values in all remaining folds, and then repeating the process with the folds switched. We note that sample splitting is distinct from cross-validation of the super learner algorithm. The final effect estimate is computed over the entire sample as usual. The sample splitting procedure used here is equivalent to the CV-TMLE approach such as is implemented in the `tlverse` R package (Coyle, Hejazi, and van der Laan 2020). However, different variations exist (Rinaldo et al. 2018; Zivich and Breskin 2020).

Standard errors for g computation were obtained from the standard deviation of 100 bootstrap resamples using the normal interval approximation (i.e., Wald method). However, for computational reasons, we were only able to apply the bootstrap to the nonparametric g computation estimator in select scenarios. Standard errors for the inverse probability weighted approach were obtained using the robust variance estimator. Standard errors for both doubly robust approaches were obtained using the variance of the efficient influence function. All confidence intervals were computed via the normal interval (i.e., Wald) equation. For each estimator in each scenario, we computed the bias: $B(\hat{\psi}) = E(\hat{\psi}) - \psi$, and 95% confidence interval coverage, defined as the proportion of 95% confidence intervals that included the true value over all 200 Monte Carlo runs. Simulations were done in R version 4.0.3 (The R Foundation, Vienna, Austria). Code to reproduce our results and additional details are available on GitHub: <https://github.com/amishler/nonparametricDoublyRobust>.

Simulation Results

Figure 1 shows the estimated absolute bias across all sample sizes for all scenarios with the stacking algorithm that included random forests and extreme gradient boosting, and which did not use sample splitting. As expected, when using the correct parametric models, all methods are unbiased. In contrast, when the transformed confounders are used with parametric models (and thus parametric models are all mis-specified), all four estimators are subject to considerable bias which does not improve as the sample size increases (Figure 1).

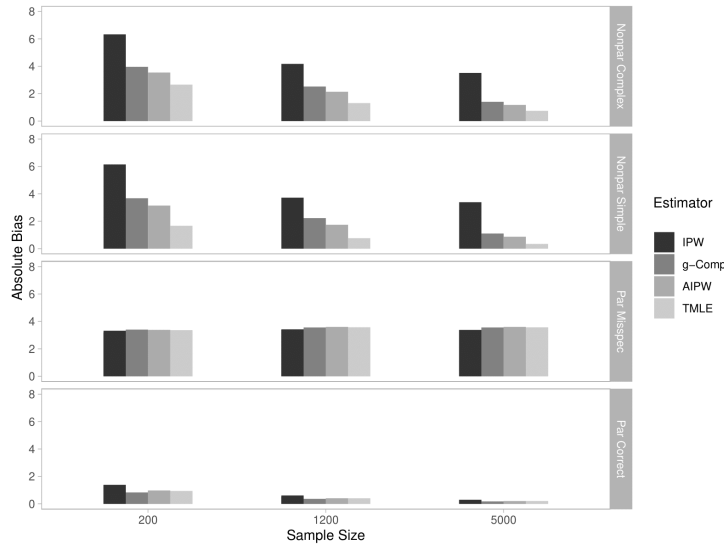


Figure 1: Absolute bias of inverse probability weighted, g-computation, and doubly robust estimators for sample sizes of $N=200$, $N=1200$, and $N=5000$. Bar color intensity, from black to light gray, represent IPW, g Computation, AIPW, and TMLE estimators, respectively. Plot panels: A) nonparametric regression with transformed covariates; B) nonparametric regression with untransformed covariates; C) parametric regression with transformed covariates; D) nonparametric regression with untransformed covariates (correctly specified parametric regression). Parametric regression included logistic regression for the exposure model, and linear regression for the outcome model. Nonparametric method consisted of a stacked generalization with random forests and extreme gradient boosting algorithms, and no sample splitting.

When models are fit nonparametrically using the simple confounders, IP-weighting displays considerable bias. G computation is also biased, but less than IP-weighting. In the nonparametric simple and complex settings (with transformed confounders), the bias decreases when doubly robust estimators are used (Figure 1). Generally, these results demonstrate what is expected from theory: the bias of singly robust estimators is larger than the bias of doubly robust estimators. Notably, in our simulation scenario under select sample sizes, the bias of the IP-weighted estimator under a nonparametric model with simple and transformed confounders is comparable to the bias of the misspecified parametric models (Figure 1).

Table 1 shows the 95% confidence interval coverage for each scenario. When correct parametric models were used, CI coverage was nominal, except for the robust variance estimator used for IP-weighting, which is known to be conservative (Hernán, Brumback, and Robins 2001) When parametric models were fit with the transformed covariates (Parametric Misspecified), coverage dropped to 46% or lower.

Table 1. Confidence interval coverage^a for sample sizes of $N = 200$, $N = 1200$, and $N = 5000$ obtained from parametric and nonparametric^b models under simple and complex confounding scenarios without sample splitting.

N	IPW	g-Comp	AIPW	TMLE	IPW	g-Comp	AIPW	TMLE
	Parametric	True			Parametric	Mispecified		
200	0.96	0.95	0.95	0.94	0.46	0.23	0.28	0.24
1200	0.98	0.93	0.94	0.94	0.01	0.00	0.00	0.00
5000	0.97	0.92	0.92	0.92	0.00	0.00	0.00	0.00
	Nonparametric Simple				Nonparametric Complex			
200	0.01	NA	0.02	0.22	0.00	NA	0.00	0.07
1200	0.02	NA	0.00	0.24	0.01	NA	0.00	0.05
5000	0.00	NA	0.02	0.29	0.00	NA	0.00	0.03

Abbreviations: IPW, inverse-probability weighting; g-Comp, g Computation; AIPW, augmented inverse-probability weighting; TMLE, targeted minimum loss-based estimation.

^a Confidence interval coverage, defined as the proportion of 95% confidence intervals that included the true value.

^b Nonparametric estimation was based on a stacked generalization with random forests and extreme gradient boosting algorithms.

The machine learning results presented in Table 1 represent version 1 of the stacked generalization when sample splitting was not used. When fit with machine learning algorithms, coverage for all estimators was well below the nominal threshold of 95%. This was true for both singly and doubly robust approaches in both simple and transformed confounder settings (Table 1).

The poor performance of machine learning methods observed in Table 1 improved under the additional strategies explored. These results are presented in Figures ?? to ??, which includes confidence

interval coverage from scenarios in which: sample splitting, generalized additive models, and confounder interactions were used with the stacking algorithms and estimators. Indeed, the highest observed coverage was 29% for TMLE in the simple confounder setting. In contrast, the lowest coverage in the simple confounder setting was 44% for TMLE with sample splitting. When sample splitting was used, AIPW and TMLE almost reached nominal coverage rates in the simple confounder setting. Coverage improved in the transformed confounder setting with sample splitting, but did not reach nominal rates.

When GAMs were combined with sample splitting, nominal coverage was attained in the simple confounder setting, but was still quite low for the transformed confounders. Coverage in the transformed confounder setting only attained nominal rates for AIPW and TMLE when sample splitting was combined with GAMs, and all confounder-confounder interactions were included in the models (Figure ??).

Discussion

Both machine learning and doubly robust estimation are becoming increasingly popular, however the relation between them remains poorly understood. Here, we have shown how machine learning methods are biased when used with singly robust estimators such as inverse probability weighting or g computation (also known as marginal standardization). Performance, however, is greatly improved when used with doubly robust approaches, particularly with sample splitting and flexible regression methods.

Doubly robust estimators can enable use of machine learning algorithms to estimate causal effects, and thus offer some protection against model misspecification. A misspecified model form can occur if the analyst fails to correctly account for the manner in which exposure and confounders relate to the outcome. For a generalized linear model, this would occur if chosen link function is not compatible with how the data were actually generated (Weisberg and

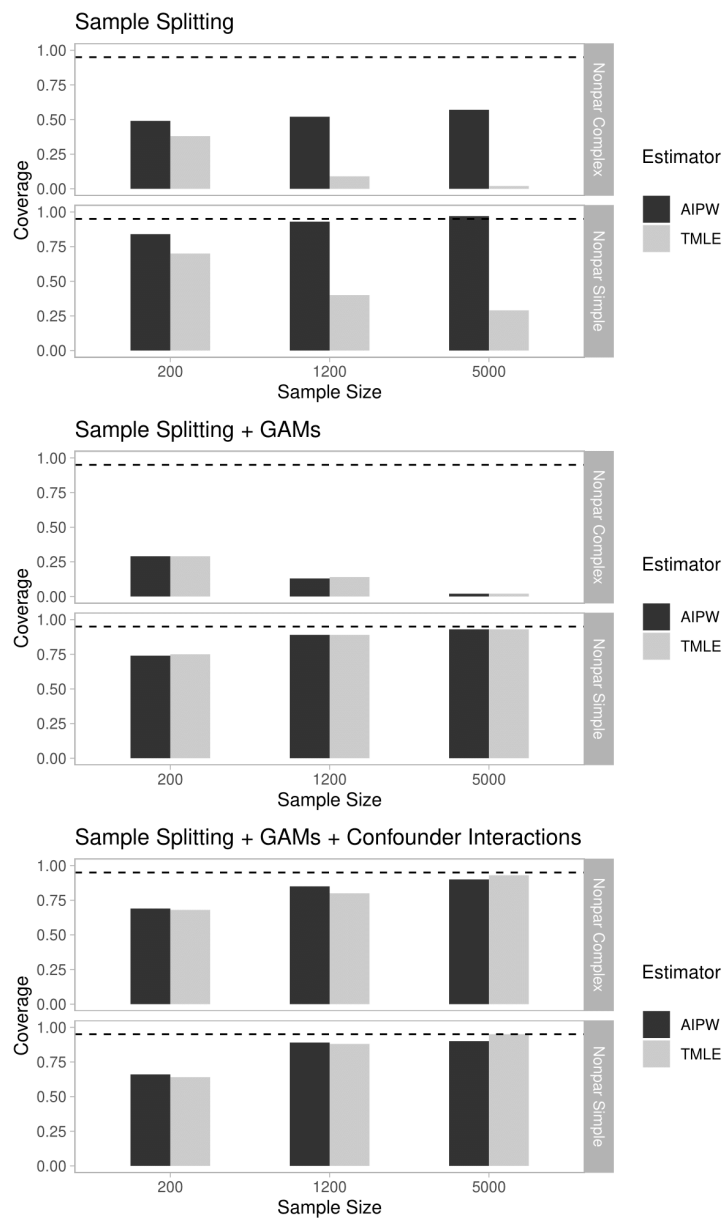


Figure 2: Coverage of doubly robust estimators for sample sizes of $N = 200$, $N = 1200$, and $N = 5000$ when models for each estimator are specified non-parametrically in the simple confounder and complex (transformed) confounder settings. Bar color black and light gray represent AIPW and TMLE estimators, respectively. Nonparametric method consisted of: a stacked generalization with random forests and extreme gradient boosting algorithms with sample splitting.

Welsh 1994). if the analyst fails to account for curvilinear relations between the covariates and the outcome, or fails to include important exposure-confounder or confounder-confounder interactions. Unfortunately, in an observational study the true nature of these relations is typically not known, which is one reason underlying the increasing popularity of machine learning methods. However, misspecification resulting in an incomplete confounder adjustment set, or incorrectly adjusting for a mediator, cannot be fixed with doubly robust machine learning methods (Keil et al. 2018).

The problems that can be encountered when using machine learning algorithms to estimate causal effects are typically attributed to the curse of dimensionality. Generally, the curse of dimensionality describes a situation where, for a given estimator, as the number of variables in a model increases, the sample size needed to maintain the same level of accuracy (expressed in terms of, e.g., bias, MSE, or coverage) increases exponentially. As we have shown, such problems will affect nonparametric (i.e., machine learning based) more profoundly, unless double robust methods are used. Indeed, Under our chosen data generating mechanisms, implementing each estimator using correct parametric models resulted in unbiased estimation. However, when implemented nonparametrically using the correct set of confounders, both g computation and inverse probability weighting were biased, while both doubly robust approaches were less biased. These results align with other work on the use of machine learning methods with double robust estimators (Zivich and Breskin 2020; Chernozhukov et al. 2018; Edward H. Kennedy 2016) and suggest that researchers should carefully weigh all considerations when using machine learning methods to estimate causal effects.

More specifically, our results suggest that when machine learning is used to quantify average treatment effects, researchers should employ the following practices to maximize the performance of the estimation approach:

1. Use doubly robust estimation methods, such as augmented inverse probability weighted or targeted minimum loss-based estimation.

2. Use sample splitting, also referred to as cross-fitting, double cross-fitting, which improves estimation of standard errors and confidence interval coverage.
3. Use a richly specified library of flexible regression, tree-based, gradient based, and other algorithms, that maximize the diversity of a given stacking algorithm.
4. Include first and higher order interactions between selected adjustment variables in a given stacking algorithm. Additionally, one may include other transformations (e.g., log, non-product interactions, or polynomial terms), as well as consider the use of screening algorithms that remove potentially unnecessary variable transformations.

While our recommendations are general enough to be considered any time researchers seek to use machine learning methods when estimating causal effects, certain limitations of our simulation study should be taken into consideration. First, we relied on only one hundred Monte Carlo samples, which is small. However, our intent was not to provide an in depth evaluation of the performance of doubly and singly robust estimators with and without machine learning methods, which has been done extensively in more technical areas (Chernozhukov et al. 2018; Tan 2010; Rose and Laan 2011). Rather, we sought to demonstrate properties of machine learning methods that are well-known in some fields, but seem to not be well appreciated among applied epidemiologists. Second, we did not focus our simulations on evaluating the relative performance of AIPW versus TMLE. Though our results might suggest that one or the other estimator performs better in certain settings, we would recommend against making such interpretations without a more in-depth exploration. Third, we only explored average treatment effect estimation for a binary point treatment and continuous confounders, but doubly robust-type methods have been developed for a wide variety of settings, including continuous (Munoz and Laan 2012; Edward H. Kennedy et al. 2017) and time-varying exposures (Edward H. Kennedy 2019) instrumental variables (Ogburn, Rot-

nitzky, and Robins 2015), mediation (Tchetgen Tchetgen and Shpitser 2012), and missing data Sun and Tchetgen Tchetgen (2018). However, we do expect our findings would apply more generally (Edward H. Kennedy 2016). Finally, our simulations were very limited in that they explored two relatively unrealistic data generating mechanisms: one simple (with untransformed confounders), and one complex (with confounders transformed via complex nonlinear functions). Nevertheless, even under our simple data generating mechanism, we were able to achieve low bias and nominal coverage only when sample splitting and flexible regression methods were used (for the simple confounder scenario), or when sample splitting, flexible regression, and confounder interactions were used (for the transformed confounder setting). We believe these findings should inform future analyses using machine learning methods with double robust estimators.

References

- Bang, Heejung, and James M. Robins. 2005. "Doubly Robust Estimation in Missing Data and Causal Inference Models." *Biometrics* 61 (4): 962–973.
- Bellman, R. 1957. *Dynamic Programming*. Rand Corporation Research Study. Princeton University Press. <https://books.google.it/books?id=wdtoPwAACAAJ>.
- Box, G. E. P. 1976. "Science and Statistics." *JASA* 71 (356): 791–99.
- Chernozhukov, Victor, Denis Chetverikov, Mert Demirer, Esther Duflo, Christian Hansen, Whitney Newey, and James Robins. 2018. "Double/Debiased Machine Learning for Treatment and Structural Parameters." *The Econometrics Journal* 21 (1): C1–68.
- Cole, Stephen R, Haitao Chu, and Sander Greenland. 2013. "Maximum Likelihood, Profile Likelihood, and Penalized Likelihood: A Primer." *Am J Epidemiol* 179 (2): 252–60.
- Coyle, Jeremey, Nima Hejazi, and Mark van der Laan. 2020. *Tlverse: R Packages for Targeted Learning*. <https://github.com/tlverse>.
- Daniel, Rhian M. 2018. "Double Robustness." In *Wiley StatsRef: Statistics Reference Online*. John Wiley & Sons, Ltd.
- Greenland, Sander, Judea Pearl, and JM Robins. 1999. "Causal Diagrams for Epidemiological Research." *Epidemiol* 10 (1): 37–48.
- Gruber, Susan, and Mark J van der Laan. 2012. "Tmle: An r Package for Targeted Maximum Likelihood Estimation." *Journal of Statistical Software* 51 (13): 1–35.
- Györfi, L., M. Kohler, A. Krzyzak, and H. Walk. 2002. *A Distribution-Free Theory of Nonparametric Regression*. New York, NY: Springer.
- Hahn, Jinyong. 1998. "On the Role of the Propensity Score in Efficient Semiparametric Estimation of Average Treatment Effects."

Econometrica 66 (2): 315–31.

Hernán, Miguel A, Babette Brumback, and James M Robins. 2001.

“Marginal Structural Models to Estimate the Joint Causal Effect of Nonrandomized Treatments.” *J Am Stat Assoc* 96 (454): 440–48.

Hernán, Miguel A, and James M Robins. 2006. “Estimating Causal Effects from Epidemiological Data.” *J Epidemiol Community Health* 60 (7): 578–86.

Jonsson-Funk, Michele, Daniel Westreich, Chris Wiesen, Til Stürmer, M. Alan Brookhart, and Marie Davidian. 2011. “Doubly Robust Estimation of Causal Effects.” *Am J Epidemiol* 173 (7): 761–67.

Kang, J. D., and Schafer J. L. 2007. “Demystifying Double Robustness: A Comparison of Alternative Strategies for Estimating a Population Mean from Incomplete Data.” *Stat Sci* 22 (4): 523–39.

Keil, Alexander P, Stephen J Mooney, Michele Jonsson Funk, Stephen R Cole, Jessie K Edwards, and Daniel Westreich. 2018. “RESOLVING AN APPARENT PARADOX IN DOUBLY ROBUST ESTIMATORS.” *Am J Epidemiol* 187 (4): 891–92.

Kennedy, Edward H. 2016. “Semiparametric Theory and Empirical Processes in Causal Inference.” In *Statistical Causal Inferences and Their Applications in Public Health Research*, edited by Hua He, Pan Wu, and Ding-Geng (Din) Chen. Switzerland: Springer International.

Kennedy, Edward H. 2019. “Nonparametric Causal Effects Based on Incremental Propensity Score Interventions.” *Journal of the American Statistical Association* 114 (526): 645–56.

Kennedy, Edward H., Zongming Ma, Matthew D. McHugh, and Dylan S. Small. 2017. “Non-Parametric Methods for Doubly Robust Estimation of Continuous Treatment Effects.” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 79 (4): 1229–45.

Kennedy, Edward H, and Sivaraman Balakrishnan. 2017. “Discussion of ‘Data-driven confounder selection via Markov and Bayesian networks’ by Jenny Häggström.” *Biometrics* In Press.

Laan, Mark J van der, Eric C Polley, and Alan E Hubbard. 2007. “Super Learner.” *Statistical Applications in Genetics and Molecular Biology* 6 (1): Article 25.

- Laan, Mark J. van der, and Daniel Rubin. 2006. "Targeted Maximum Likelihood Learning." *Int J Biostat* 2 (1): Article 11.
- Lee, Brian K., Justin Lessler, and Elizabeth A. Stuart. 2010. "Improving Propensity Score Weighting Using Machine Learning." *Stat Med* 29 (3): 337–46.
- Li, Ker-Chau. 1989. "Honest Confidence Regions for Nonparametric Regression." *The Annals of Statistics* 17 (3): 1001–8.
- Long, Qi, Chiu-Hsieh Hsu, and Yisheng Li. 2012. "Doubly Robust Nonparametric Multiple Imputation for Ignorable Missing Data." *Stat Sin* 22: 149–72.
- Longford, NT. 2008. *Studying Human Populations: An Advanced Course in Statistics*. New York: Springer.
- Metropolis, N, and S Ulam. 1949. "The Monte Carlo method." *J Am Stat Assoc* 44 (247): 335–41.
- Munoz, Ivan Diaz, and Mark van der Laan. 2012. "Population Intervention Causal Effects Based on Stochastic Interventions." *Biometrics* 68 (2): 541–49.
- Naimi, AI, RW Platt, and Larkin JC. 2018. "Machine Learning for Fetal Growth Prediction." *Epidemiol* 29 (2): 290–98.
- Naimi, Ashley I, Stephen R Cole, and Edward H Kennedy. 2017. "An Introduction to G Methods." *Int J Epidemiol* 46 (2): 756–62.
- Nelder, J. A., and R. W. M. Wedderburn. 1972. "Generalized Linear Models." *JRSS-A* 135 (3): 370–84.
- Ogburn, Elizabeth L., Andrea Rotnitzky, and James M. Robins. 2015. "Doubly Robust Estimation of the Local Average Treatment Effect Curve." *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 77 (2): 373–96.
- Oulhote, Youssef, Brent Coull, Marie-Abele Bind, Frodi Debes, Fleming Nielsen, Ibon Tamayo, Pal Weihe, and Philippe Grandjean. 2019. "Joint and Independent Neurotoxic Effects of Early Life Exposures to a Chemical Mixture: A Multi-Pollutant Approach Combining Ensemble Learning and g-Computation." *Environmental Epidemiology* 3 (5): e063.
- Rencher, Alvin C. 2000. *Linear Models in Statistics*. New York: Wiley.
- Rinaldo, Alessandro, Larry Wasserman, Max G'Sell, and Jing Lei.

2018. "Bootstrapping and Sample Splitting for High-Dimensional, Assumption-Free Inference." <https://arxiv.org/abs/1611.05401>.
- Robins, J M, and Y Ritov. 1997. "Toward a Curse of Dimensionality Appropriate (CODA) Asymptotic Theory for Semi-Parametric Models." *Stat Med* 16 (1-3): 285-319.
- Robins, James M, and Miguel Á Hernán. 2009. "Estimation of the Causal Effects of Time-Varying Exposures." In *Advances in Longitudinal Data Analysis*, edited by G Fitzmaurice, M. Davidian, G Verbeke, and G Molenberghs, 553-99. Boca Raton, FL: Chapman & Hall.
- Robins, JM. 2001. "Data, Design, and Background Knowledge in Etiologic Inference." *Epidemiol* 12 (3): 313-20.
- Robins, JM, and Andrea Rotnitzky. 1995. "Semiparametric Efficiency in Multivariate Regression Models with Missing Data." *JASA* 90 (429): 122-29.
- . 2001. "Comment: Inference for Semiparametric Models: Some Questions and an Answer." *Statistica Sinica* 11 (4): 920-36.
- Rose, Sherri, and Mark J van der Laan. 2011. *Targeted Learning: Causal Inference for Observational and Experimental Data*. New York, NY: Springer.
- Rotnitzky, Andrea, and Stijn Vansteelandt. 2014. "Double-Robust Methods." In *Handbook of Missing Data Methodology*, edited by Geert Molenberghs, Garrett Fitzmaurice, Michael G. Kenward, Anastasios Tsiatis, and Geert Verbeke, 185-209. CRC Press.
- Snowden, Jonathan M., Sherri Rose, and Kathleen M. Mortimer. 2011. "Implementation of g-Computation on a Simulated Data Set: Demonstration of a Causal Inference Technique." *Am J Epidemiol* 173 (7): 731-38.
- Sun, BaoLuo, and Eric J Tchetgen Tchetgen. 2018. "On Inverse Probability Weighting for Nonmonotone Missing at Random Data." *J Am Stat Assoc* 113 (521): 369-79.
- Tan, Zhiqiang. 2010. "Bounded, Efficient and Doubly Robust Estimation with Inverse Weighting." *Biometrika* 97 (3): 661-82.
- Tchetgen Tchetgen, Eric J., and Ilya Shpitser. 2012. "Semiparametric Theory for Causal Mediation Analysis: Efficiency Bounds, Mul-

- tiple Robustness and Sensitivity Analysis." *Annals of Statistics* 40 (3): 1816–45.
- Vaart, A. W. van der. 2000. *Asymptotic Statistics*. Cambridge: Cambridge University Press.
- Wasserman, Larry. 2006. *All of Nonparametric Statistics*. New York; London: Springer.
- Weisberg, S., and A. H. Welsh. 1994. "Adapting for the Missing Link." *The Annals of Statistics* 22 (4): 1674—1700.
- Westreich, Daniel, Stephen R. Cole, Enrique F. Schisterman, and Robert W. Platt. 2012. "A simulation study of finite-sample properties of marginal structural Cox proportional hazards models." *Stat Med* 31 (19): 2098–2109.
- Westreich, Daniel, Justin Lessler, and Michele Jonsson Funk. 2010. "Propensity Score Estimation: Neural Networks, Support Vector Machines, Decision Trees (CART), and Meta-Classifiers as Alternatives to Logistic Regression." *J Clin Epidemiol* 63 (8): 826–33.
- Zhong, Yongqi, Edward H Kennedy, Lisa M Bodnar, and Ashley I Naimi. 2021. "AIPW: An r Package for Augmented Inverse Probability Weighted Estimation of Average Causal Effects." *American Journal of Epidemiology* In Press.
- Zivich, Paul N, and Alexander Breskin. 2020. "Machine Learning for Causal Inference: On the Use of Cross-Fit Estimators." *arXiv:2004.10337*.