BAYESIAN REGRESSION TREE MODELS FOR CAUSAL INFERENCE: REGULARIZATION, CONFOUNDING, AND HETEROGENEOUS EFFECTS

By P. Richard Hahn, Jared Murray and Carlos M. Carvalho

Arizona State University and University of Texas

This paper develops a semi-parametric Bayesian regression model for estimating heterogeneous treatment effects from observational data. Standard nonlinear regression models, which may work quite well for prediction, can yield badly biased estimates of treatment effects when fit to data with strong confounding. Our Bayesian causal forest model avoids this problem by directly incorporating an estimate of the propensity function in the specification of the response model, implicitly inducing a covariate-dependent prior on the regression function. This new parametrization also allows treatment heterogeneity to be regularized separately from the prognostic effect of control variables, making it possible to informatively "shrink to homogeneity", in contrast to existing Bayesian non- and semi-parametric approaches. We illustrate the benefits of this approach via the reanalysis of an observational study assessing the causal effects of smoking on medical expenditures as well as extensive simulation studies.

1. Introduction. The success of modern predictive modeling is founded on the understanding that flexible predictive models must be carefully regularized in order to achieve good out-of-sample performance (low generalization error). In a causal inference setting, regularization is less straightforward: In the presence of confounding, regularized models originally designed for prediction can bias causal estimates towards some unknown function of high dimensional nuisance parameters (Hahn et al., 2016). That is, despite offering excellent predictive performance, the causal conclusions from a naively regularized nonlinear regression are likely to be substantially biased, leading to high estimation error of the target parameter. A key finding in this paper is that this effect will be especially pronounced in flexible models which allow for heterogeneous effects.

To mitigate these estimation problems we propose a flexible sum-of-regression-trees — a forest

Keywords and phrases: Bayesian; Causal inference; Heterogeneous treatment effects; Predictor-dependent priors; Machine learning; Regression trees; Regularization; Shrinkage

— to model a response variable as a function of a binary treatment indicator and a vector of control variables. Our work departs from existing contributions — primarily Hill (2011) and later extensions — in two important respects: First, we develop a novel prior for the response surface that depends explicitly on estimates of the propensity score as an important 1-dimensional transformation of the covariates (including the treatment assignment). Incorporating this transformation of the covariates is not strictly necessary in response surface modeling, but we show that it can substantially improve treatment effect estimation in the presence of moderate to strong confounding, especially when that confounding is driven by targeted selection — individuals selecting into treatment based on somewhat accurate predictions of the potential outcomes.

Second, we represent our regression as a sum of two functions: the first models the *prognostic* impact of the control variables (the component of the conditional mean of the response that is unrelated to the treatment effect), while the second represents the treatment effect directly, which itself is a nonlinear function of the observed attributes (capturing possibly heterogeneous effects). We represent each function as a forest. This approach allows the degree of shrinkage on the treatment effect to be modulated *directly* and *separately* of the prognostic effect. In particular, under this parametrization, standard regression tree priors shrink towards homogeneous effects.

In previous approaches, the prior distribution over treatment effects is induced indirectly, and is therefore difficult to understand and control. Our approach interpolates between two extemes: Modeling the conditional means of treated and control units entirely separately, or including treatment assignment as "just another covariate" (see also Künzel et al. (2017) for detailed discussion of the tradeoffs between these two approaches). The former precludes any borrowing or regularization entirely, while the second can be rather difficult to understand using flexible models. Parametrizing non- and semiparametric models this way is attractive regardless of the specific priors in use.

Comparisons on simulated data show that the new model — which we call the Bayesian causal forest model — performs at least as well as existing approaches for estimating heterogenous treatment effects across a range of plausible data generating processes. More importantly, it performs dramatically better in many cases, especially those with strong confounding, targeted selection, and relatively weak treatment effects, which we believe to be common in applied settings.

In section 7, we demonstrate how our flexible Bayesian model allows us to make rich inferences on heterogeneous treatment effects, including estimates of average and conditional average treatment effects at various levels, in a re-analysis of data from an observational study of the effect of smoking on medical expenditures.

As we have noted, the Bayesian causal forest model directly extends ideas from two earlier papers: Hill (2011) and Hahn et al. (2016). Specifically, this paper studies the "regularization-induced confounding" of Hahn et al. (2016) in the context of nonparametric Bayesian models as utilized by Hill (2011). In terms of implementation, this paper builds explicitly on the work of Chipman, George and McCulloch (2010); see also Gramacy and Lee (2008) and Murray (2017). Other notable work on Bayesian treatment effect estimation includes Gustafson and Greenland (2006), Zigler and Dominici (2014), Heckman, Lopes and Piatek (2014), Li and Tobias (2014), ? and Taddy et al. (2016). A more complete discussion of how the new method relates to this earlier literature, including non-Bayesian approaches, is deferred until Section 8.

2. Problem statement and notation. Let Y denote a scalar response variable and Z denote a binary treatment indicator variable. Capital Roman letters denote random variables, while realized values appear in lower case, that is, y and z. Let x denote a length d vector of observed control variables. Throughout, we will consider an observed sample of size n independent observations (Y_i, Z_i, x_i) , for $i = 1, \ldots n$. When Y or Z (respectively, y or z) are without a subscript, they denote length n column vectors; likewise, X will denote the $n \times d$ matrix of control variables.

We are interested in estimating various treatment effects. In particular, we are interested in conditional average treatment effects (CATE) — the amount by which the response Y_i would differ between hypothetical worlds in which the treatment was set to $Z_i = 1$ versus $Z_i = 0$, averaged across subpopulations defined by attributes x. This kind of counterfactual estimand can be formalized in the potential outcomes framework (Imbens and Rubin (2015), chapter 1) by using $Y_i(0)$ and $Y_i(1)$ to denote the outcomes we would have observed if treatment were set to zero or one, respectively. We make the stable unit treatment value assumption (SUTVA) throughout (excluding interference between units and multiple versions of treatment (Imbens and Rubin, 2015)). We observe the potential outcome that corresponds to the realized treatment: $Y_i = Z_i Y_i(1) + (1 - Z_i) Y_i(0)$.

Throughout the paper we will assume that strong ignorability holds, which stipulates that

$$(1) Y_i(0), Y_i(1) \perp Z_i \mid \mathbf{X}_i.$$

and also that

(2)
$$0 < \Pr(Z_i = 1 \mid x_i) < 1$$

for all i = 1, ..., n. The first condition assumes we have no unmeasured confounders, and the second condition (overlap) is necessary to estimate treatment effects everywhere in covariate space. Provided that these conditions hold, it follows that $E(Y_i(z) \mid x_i) = E(Y_i \mid x_i, Z_i = z)$ so our estimand may be expressed as

(3)
$$\tau(\mathbf{x}_i) := \mathbf{E}(Y_i \mid \mathbf{x}_i, Z_i = 1) - \mathbf{E}(Y_i \mid \mathbf{x}_i, Z_i = 0).$$

For simplicity, we restrict attention to mean-zero additive error representations

(4)
$$Y_i = f(\mathbf{x}_i, Z_i) + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2)$$

so that $E(Y_i \mid x_i, Z_i = z_i) = f(x_i, z_i)$. In this context, (1) can be expressed equivalently as $\epsilon_i \perp \!\!\! \perp Z_i \mid x_i$. The treatment effect of setting $z_i = 1$ versus $z_i = 0$ can therefore be expressed as

$$\tau(\mathbf{x}_i) := f(\mathbf{x}_i, 1) - f(\mathbf{x}_i, 0).$$

Our contribution in this paper is a careful study of prior specification for f. We propose new prior distributions that improve estimation of the parameter of interest, namely τ . Previous work (Hill, 2011) advocated using a Bayesian additive regression tree (BART) prior for $f(\mathbf{x}_i, z_i)$ directly. We instead recommend expressing the response surface as

(5)
$$E(Y_i \mid \mathbf{x}_i, Z_i = z_i) = \mu(\mathbf{x}_i, \hat{\pi}(\mathbf{x}_i)) + \tau(\mathbf{x}_i)z_i,$$

where the functions μ and τ are given independent BART priors and $\hat{\pi}(\mathbf{x}_i)$ is an estimate of the propensity score $\pi(\mathbf{x}_i) = \Pr(Z_i = 1 \mid \mathbf{x}_i)$. The following sections motivate this model specification and provide additional context; further modeling details are given in Section 5.

3. Bayesian additive regression trees for heterogeneous treatment effect estimation.

Hill (2011) observed that under strong ignorability, treatment effect estimation reduces to response surface estimation. That is, provided that a sufficiently rich collection of control variables are available (to ensure strong ignorability), treatment effect estimation can proceed "merely" by estimating the conditional expectations $E(Y \mid x, Z = 1)$ and $E(Y \mid x, Z = 0)$. Noting its strong performance in prediction tasks, Hill (2011) advocates the use of the Bayesian additive regression tree (BART) model of Chipman, George and McCulloch (2010) for estimating these conditional expectations.

BART is particularly well-suited to detecting interactions and discontinuities, can be made invariant to monotone transformations of the covariates, and typically requires little parameter tuning.

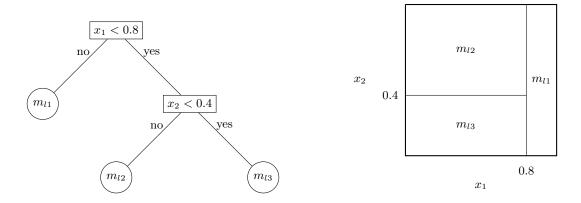


FIG 1. (Left) An example binary tree, with internal nodes labelled by their splitting rules and terminal nodes labelled with the corresponding parameters m_{lb} . (Right) The corresponding partition of the sample space and the step function.

Chipman, George and McCulloch (2010) provide extensive evidence of BART's excellent predictive performance. BART has also been used successfully for applications in causal inference, for example Green and Kern (2012), Hill et al. (2013), Kern et al. (2016), and Sivaganesan, Müller and Huang (2017). It has subsequently been demonstrated to successfully infer heterogeneous and average treatment effects in multiple independent simulation studies (Dorie et al., 2017; Wendling et al., 2018), frequently outperforming competitors (and never lagging far behind).

Despite its excellent performance in practice, there are limited theoretical results about BART. There have been recent developments on posterior consistency and rates of posterior concentration for Bayesian tree models in prediction contexts (Linero and Yang, 2017; Rockova and van der Pas, 2017). These results require significant modifications to the BART prior, however, which we do not further investigate here. To the extent that these results are informative or suggestive about BART's impressive performance in prediction, however, we expect those insights to carry over to the treatment effect models presented here.

3.1. Specifying the BART prior. The BART prior expresses an unknown function f(x) as a sum of many piecewise constant binary regression trees. (In this section, we suppress z in the notation; implicitly z may be considered as a coordinate of x.) Each tree T_l , $1 \le l \le L$, consists of a set of internal decision nodes which define a partition of the covariate space (say $A_1, \ldots, A_{B(l)}$), as well as a set of terminal nodes or leaves corresponding to each element of the partition. Further, each element of the partition A_b is associated a parameter value, m_{lb} . Taken together the partition and the leaf parameters define a piecewise constant function: $g_l(x) = m_{lb}$ if $x \in A_b$; see Figure 1.

Individual regression trees are then additively combined into a single regression forest: f(x) =

 $\sum_{l=1}^{L} g_l(\mathbf{x})$. Each of the functions g_l are constrained by their prior to be "weak learners" in the sense that the prior favors small trees and leaf parameters that are near zero. Each tree follows (independently) the prior described in Chipman, George and McCulloch (1998): the probability that a node at depth h splits is given by $\eta(1+h)^{-\beta}$, $\eta \in (0,1)$, $\beta \in [0,\infty)$.

A variable to split on, as well as a cut-point to split at, are then selected uniformly at random from the available splitting rules. Large, deep trees are given extremely low prior probability by taking $\eta = 0.95$ and $\beta = 2$ as in Chipman, George and McCulloch (2010). The leaf parameters are assigned independent priors $m_{lb} \sim N(0, \sigma_m^2)$ where $\sigma_m = \sigma_0/\sqrt{L}$. The induced marginal prior for $f(\mathbf{x})$ is centered at zero and puts approximately 95% of the prior mass within $\pm 2\sigma_0$ (pointwise), and σ_0 can be used to calibrate the plausible range of the regression function. Full details of the BART prior and its implementation are given by Chipman, George and McCulloch (2010).

In our context we are concerned with the impact that the prior over f(x, z) has on estimating $\tau(x) = f(x, 1) - f(x, 0)$. The choice of BART as a prior over f has particular implications for the induced prior on τ that are difficult to understand: In particular, the induced prior will vary with the dimension of x and the degree of dependence with z. In Section 5 we propose an alternative parameterization that mitigates this problem. But first, the next section develops a more general framework for investigating the influence of prior specification and regularization on treatment effect estimates.

- 4. The central role of the propensity score in regularized causal modeling. In this section we explore the joint impacts of regularization and confounding on estimation of heterogeneous treatment effects. We find that including an estimate of the propensity score as a covariate reduces the bias of regularized treatment effect estimates in finite samples. We recommend including an estimated propensity score as a covariate as routine practice regardless of the particular models or algorithms used to estimate treatment effects since regularization is necessary to estimate heterogeneous treatment effects non- or semiparamaterically or in high dimensions. To illustrate the potential for biased estimation and motivate our fix, we introduce two key concepts: Regularization induced confounding and targeted selection.
- 4.1. Regularization-induced confounding. Since treatment effects may be deduced from the conditional expectation function $f(\mathbf{x}_i, z_i)$, a likelihood perspective suggests that the conditional distribution of Y given x and Z is sufficient for estimating treatment effects. While this is true in terms

of identification of treatment effects, the question of estimation with finite samples is more nuanced. In particular, many functions in the support of the prior will yield approximately equivalent likelihood evaluations, but may imply substantially different treatment effects. This is particularly true in a strong confounding-modest treatment effect regime, where the conditional expectation of Y is largely determined by x rather than Z.

Accordingly, the posterior estimate of the treatment effect is apt to be substantially influenced by the prior distribution over f for realistic sample sizes. This issue was explored by Hahn et al. (2016) in the narrow context of linear regression with continuous treatment and homogenous treatment effect; they call this phenomenon "regularization-induced confounding" (RIC). In the linear regression setting an exact expression for the bias on the treatment effect under standard regularization priors is available in closed form.

Example: RIC in the linear model. Suppose the treatment effect is homogenous and response and treatment model are both linear:

(6)
$$Y_i = \tau Z_i + \beta^t \mathbf{x}_i + \varepsilon_i,$$
$$Z_i = \gamma^t \mathbf{x}_i + \nu_i;$$

where the error terms are mean zero Gaussian and a multivariate Gaussian prior is placed over all regression coefficients. The Bayes estimator under squared error loss is the posterior mean, so we examine the expression for the bias of $\hat{\tau}_{rr} \equiv E(\tau \mid Y, z, \mathbf{X})$. We begin from a standard expression for the bias of the ridge estimator, as given, for example, in Giles and Rayner (1979). Write $\theta = (\tau, \beta^t)^t$,

$$\tilde{\mathbf{X}} = \begin{pmatrix} \mathbf{z} & \mathbf{X} \end{pmatrix}$$

and let $\theta \sim N(0, \mathbf{M}^{-1})$. Then the bias of the Bayes estimator is

(7)
$$\operatorname{bias}(\hat{\theta}_{rr}) = -(\mathbf{M} + \tilde{\mathbf{X}}^t \tilde{\mathbf{X}})^{-1} \mathbf{M} \theta$$

where the bias expectation is taken over Y, conditional on X and all model parameters.

Consider $M = \begin{pmatrix} 0 & 0 \\ 0 & I_p \end{pmatrix}$, where I_p denotes a p-by-p identity matrix, which corresponds to a ridge prior (with ridge parameter $\lambda = 1$ for simplicity) on the control variables and a non-informative "flat" prior over the first element (τ , the treatment effect). Plugging this into the bias equation (7) and noting that

$$(\mathbf{M} + ilde{\mathbf{X}}^t ilde{\mathbf{X}})^{-1} = egin{pmatrix} \mathbf{z}^t \mathbf{z} & \mathbf{z}^t \mathbf{X} \ \mathbf{X}^t \mathbf{z} & \mathbf{X}^t \mathbf{X} + \mathbf{I}_p \end{pmatrix}^{-1}$$

we obtain

(8)
$$\operatorname{bias}(\hat{\tau}_{rr}) = -\left((\mathbf{z}^t \mathbf{z})^{-1} \mathbf{z}^t \mathbf{X} \right) (\mathbf{I} + \mathbf{X}^t (\mathbf{X} - \hat{\mathbf{X}}_{\mathbf{z}}))^{-1} \beta,$$

where $\hat{\mathbf{X}}_z = \mathbf{z}(\mathbf{z}^t\mathbf{z})^{-1}\mathbf{z}^t\mathbf{X}$. Notice that the leading term $((\mathbf{z}^t\mathbf{z})^{-1}\mathbf{z}^t\mathbf{X})$ is a vector of regression coefficients from p univariate regressions predicting X_j given z. With completely randomized treatment assignment these terms will tend to be near zero (and precisely zero in expectation over Z). This ensures that the ridge estimate of τ is nearly unbiased, despite the fact that the middle matrix is generally nonzero. However, in the presence of selection some of these regression coefficients will be non-zero due to the correlation between Z and the covariates in \mathbf{X} . As a result, the bias of $\hat{\tau}_{rr}$ will depend on the form of the design matrix and unknown nuisance parameters β .

The problem here is not simply that $\hat{\tau}_{rr}$ is biased — after all, the insight behind regularization is that some bias can actually improve our average estimation error. Rather, the problem is that the degree of bias is not under the analyst's control (as it depends on unknown nuisance parameters). The use of a naive regularization prior in the presence of counfounding can unwittingly induce extreme bias in estimation of the target parameter, even when all the confounders are measured and the parametric model is correctly specified.

In more complicated nonparametric regression models with heterogeneous treatment effects a closed-form expression of the bias is not generally available; see Yang, Cheng and Dunson (2015) and Chernozhukov et al. (2016) for related results in a partially linear model where effects are homogeneous but the β^t x term above is replaced by a nonlinear function. However, note that both of these theoretical results consider asymptotic bias in semi- and non-parametric Bayesian and frequentist inference; our attention here to the simple case of the linear model shows that the phenomenon occurs in finite samples even in a parametric model. That said, the RIC phenomenon can be reliably recreated in nonlinear, semiparametric settings. The easiest way to demonstrate this is by considering scenarios where selection into treatment is based on expected outcomes under no treatment, a situation we call targeted selection.

4.2. Targeted selection. Targeted selection refers to settings where treatment is assigned based on a prediction of the outcome in the absence of treatment, given measured covariates. That is, targeted selection asserts that treatment is being assigned, in part, based on an estimate of the expected potential outcome $\mu(x) := E(Y(0) \mid x)$ and that the probability of treatment is generally increasing or decreasing as a function of this estimate. We suspect this selection process is quite

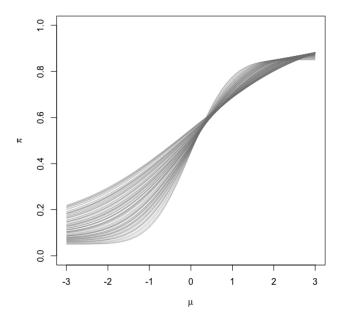


FIG 2. For any value of \tilde{x} , the propensity score $\pi(\mu, \tilde{x})$ is monotone in the prognostic function μ . Here, many realizations of this function are shown for different values of \tilde{x} .

common in practice; for example, in medical contexts where risk factors for adverse outcomes are well-understood physicians are more likely to assign treatment to patients with worse expected outcomes in its absence.

Targeted selection implies that there is a particular functional relationship between the propensity score π and the expected outcomes without treatment μ . In particular, suppose for simplicity that there exists a change of variables $\mathbf{x} \to (\mu(\mathbf{x}), \tilde{\mathbf{x}})$ that takes the prognostic function $\mu(\mathbf{x})$ to the first element of the covariate vector. Then targeted selection says that for every $\tilde{\mathbf{x}}$, the propensity function $\mathbf{E}(Z \mid \mathbf{x}) = \pi(\mu, \tilde{\mathbf{x}})$ is (approximately) monotone in μ ; see Figure 2 for a visual depiction. If the relationship is strictly monotone so that π is invertible in μ for any $\tilde{\mathbf{x}}$, this in turn implies that $\mu(\mathbf{x})$ is a function of $\pi(\mathbf{x})$.

Targeted selection and RIC in the linear model. To help understand how targeted selection leads to RIC, it is helpful to again consider the linear model. There, one can describe RIC in terms of three components: the coefficients defining the propensity function $E(Z \mid x) = \gamma x$, the coefficients defining the prognostic function, $E(Y \mid Z = 0, x = x)$, and the strength of the selection as measured

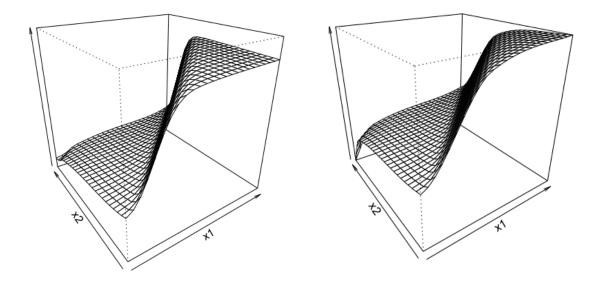


Fig 3. Left panel: The propensity function, π , shown for various values of $\tilde{\mathbf{x}}$. The "shelf" at the line $x_1 = x_2$ is a complex shape for many regression methods to represent. Right panel: the analogous plot for the prognostic function μ . Note the similar shapes due to targeted selection; the π function falls between 0 to 1, while the μ function ranges from -3 to 3.

by $Var(Z \mid x) = Var(\nu)$. Specifically, note the identity

(9)
$$E(Y \mid \mathbf{x}, Z) = (\tau + b)Z + (\beta - b\gamma)^t \mathbf{x} - b(Z - \gamma^t \mathbf{x}) = \hat{\tau}Z + \hat{\beta}^t \mathbf{x} - \hat{\epsilon},$$

which is true for any value of the scalar parameter b, the bias of $\hat{\tau}$. Intuitively, if neighborhoods of $\hat{\beta} = (\beta - b\gamma)$ have higher prior probability than β and $Var(\hat{\epsilon}) = b^2Var(\nu)$ is small on average relative to σ^2 , then the posterior distribution for τ is apt to be biased toward $\hat{\tau} = \tau + b$.

The bias will be large precisely when confounding is strong and the selection is targeted: For non-negligible bias the term $b^2 \text{Var}(\nu)$ is smallest when $\text{Var}(\nu)$ is small, that is, when selection (hence, confounding) is strong. For priors on β that are centered at zero —which is overwhelmingly the default — the $(\beta - b\gamma)$ term can be made most favorable with respect to the prior when the vector β and γ have the same direction, which corresponds to perfectly targeted selection.

Targeted selection and RIC in nonlinear models. To investigate RIC in more complex regression settings, we start with a simple 2-d example characterized by targeted selection:

Table 1

The standard BART prior exhibits substantial bias in estimating the treatment effect, poor coverage of 95% posterior (quantile-based) credible intervals, and high root mean squared error (rmse). A modified BART prior (denoted BCF) allows splits in an estimated propensity score; it performs markedly better on all three metrics.

Prior	bias	coverage	rmse
BART	0.27	65%	0.31
BCF	0.14	95%	0.21

Example 1: d = 2, n = 250, homogeneous effects. Consider the following simple data generating process:

$$Y_{i} = \mu(x_{1}, x_{2}) - \tau Z_{i} + \epsilon_{i},$$

$$E(Y_{i} \mid x_{i1}, x_{i2}, Z_{i} = 1) = \mu(x_{1}, x_{2}),$$

$$(10) \qquad E(Z_{i} \mid x_{i1}, x_{i2}) = \pi(\mu(x_{i1}, x_{i2}), x_{i1}, x_{i2}),$$

$$= 0.8\Phi\left(\frac{\mu(x_{i1}, x_{i2})}{0.1(2 - x_{i1} - x_{i2}) + 0.25}\right) + 0.025(x_{i1} + x_{i2}) + 0.05$$

$$\epsilon_{i} \stackrel{\text{iid}}{\sim} N(0, 1), \quad x_{i1}, x_{i2} \stackrel{\text{iid}}{\sim} \text{Uniform}(0, 1).$$

Suppose that in (10) Y is a continuous biometric measure of heart distress, Z is an indicator for having received a heart medication, and x_1 and x_2 are systolic and diastolic blood pressure (in standardized units), respectively. Suppose that it is known that the difference between these two measurements is prognostic of high distress levels, with positive levels of $x_1 - x_2$ being a critical threshold. At the same time, suppose that prescribers are targeting the drug towards patients with high levels of diagnostic markers, so the probability of receiving the drug is an increasing function in μ . Figure 3 shows π as a function of x_1 and x_2 ; figure 2 shows the relationship between μ and π for various values of $\tilde{\mathbf{x}} = x_1 + x_2$.

We simulated 200 datasets of size n=250 according to this data generating process with $\tau=-1$. With only a few covariates, low noise, and a relatively large sample size, we might expect most methods to perform well here. Table 1 shows that standard, unmodified BART exhibits high bias and root mean squared error (RMSE) as well as poor coverage of 95% credible intervals. Our proposed fix (detailed below) improves on both estimation error and coverage, primarily by including an estimate of π as a covariate.

What explains BART's relatively poor performance on this DGP? First, strong confounding and targeted selection implies that μ is approximately a monotone function of π alone (Figure 4). However, π (and hence μ) is difficult to learn via regression trees — it takes many axis-aligned splits to approximate the "shelf" across the diagonal (see Figure 5), and the BART prior specifically

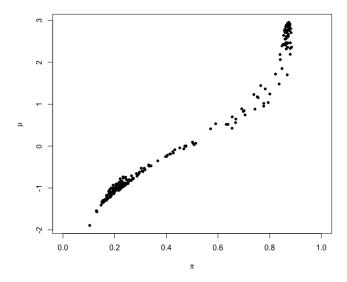


FIG 4. This scatterplot depicts $\mu(x) = E(Y \mid Z = 0, x)$ and $\pi(x) = E(Z \mid x)$ for a realization from the data generating process from the above example. It shows clear evidence of targeted selection. Such plots, based on estimates $(\hat{\mu}, \hat{\pi})$ can provide evidence of (strong) targeted selection in empirical data.

penalizes this kind of complexity. At the same time, due to the strong confounding in this example a single split in Z can stand in for many splits on x_1 and x_2 that would be required to approximate $\mu(\mathbf{x})$. These simpler structures are favored by the BART prior, leading to RIC.

Before discussing how we reduce RIC, we note that this example is somewhat stylized in that we designed it specifically to be difficult to learn for tree-based models. Other models might suffer less from RIC on this particular example. However, any informative, sparse, or nonparametric prior distribution – any method that imposes meaningful regularization – is susceptible to similar effects, as they prioritize some data-generating processes at the expense of others. Absent prior knowledge of the form of the treatment assignment and outcome models, it is impossible to know a prior whether RIC will be an issue. Fortunately it is relatively straightforward to minimize the risk of bias due to RIC.

4.3. Mitigating RIC with covariate-dependent priors. Finally, we arrive at the role of the propensity score in a regularized regression context. The potential for RIC is strongest when $\mu(x)$ is exactly or approximately a function of $\pi(x)$ and when the composition of the two has relatively low prior support. This can lead the model to misattribute the variability of μ , in the direction of π , to Z. A natural solution to this problem would be to include $\pi(x)$ as a covariate, so that it is penalized

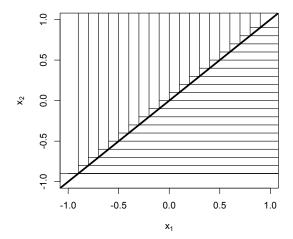


FIG 5. Many axis-aligned splits are required to approximate a step function (or near-step function) along the diagonal in the outcome model, as in Fig. 3 (right panel). Since these two regions correspond also to disparate rates of treatment, tree-based regularized regression is apt to overstate the treatment effect.

equitably with changes in the treatment variable Z. That is, when evaluating candidate functions for our estimate of $E(Y \mid x, z)$ we want representations involving $\pi(x)$ to be regularized/penalized the same as representations involving z. Of course π is unknown and must be estimated, but this is a straightforward regression problem.

Mitigating RIC in the linear model. Given an estimate of the propensity function $\hat{z}_i \approx \gamma^t \mathbf{x}_i$, we consider the over-complete regression that includes as regressors both z and \hat{z} . Our design matrix becomes

$$\tilde{\mathbf{X}} = \begin{pmatrix} \mathbf{z} & \hat{\mathbf{z}} & \mathbf{X} \end{pmatrix}$$
.

This covariate matrix is degenerate because \hat{z} is in the column span of **X** by construction. In a regularized regression proglem this degeneracy is no obstacle. Applying the expression for the bias from above, with a flat prior over the coefficient associated with \hat{z} , yields

$$\operatorname{bias}(\hat{\tau}_{rr}) = -\left\{ (\tilde{\mathbf{z}}^t \tilde{\mathbf{z}})^{-1} \tilde{\mathbf{z}}^t \mathbf{X} \right\}_1 (\mathbf{I} + \mathbf{X}^t (\mathbf{X} - \hat{\mathbf{X}}_{\mathbf{z}}))^{-1} \beta = 0,$$

where $\tilde{\mathbf{z}} = (\mathbf{z} \ \hat{\mathbf{z}})$ and $\{(\tilde{\mathbf{z}}^t\tilde{\mathbf{z}})^{-1}\tilde{\mathbf{z}}^t\mathbf{X}\}_1$ denotes the top row of $\{(\tilde{\mathbf{z}}^t\tilde{\mathbf{z}})^{-1}\tilde{\mathbf{z}}^t\mathbf{X}\}$, which corresponds to the regression coefficient associated with z in the two variable regression predicting X_j given \tilde{z} . Because \hat{z} captures the observed association between z and x, z is conditionally independent of x given \hat{z} , from which we conclude that these regression coefficients will be zero. See Yang, Cheng

and Dunson (2015) for a similar de-biasing strategy in a partially linear semiparametric context.

Mitigating RIC in nonlinear models. The same strategy also proves effective in the nonlinear setting — simply by including an estimate of the propensity score as a covariate in the BART model, the RIC effect is dramatically mitigated, as can be seen in the second row of Table 1. From a Bayesian perspective, this is simply a judicious variable transformation since our regression model is specified conditional on both Z and x — we are not obliged to consider uncertainty in our estimate of π to obtain valid posterior inference. We obtain another example of a covariate dependent prior, similar to Zellner's g-prior (albeit motivated by very different considerations). See section 8 for additional discussion of this point.

To summarize, although it has long been known that the propensity score is a sufficient dimension reduction for estimation of the ATE – and that combining estimates of the response surface and propensity score can improve estimation of average treatment effects (Bang and Robins, 2005), we find that incorporating an estimate of the propensity score into estimation of the response surface can improve estimation of average treatment effects in finite samples. As we will demonstrate in Section 6, these benefits also accrue when estimating (heterogeneous) conditional average treatment effects. Estimating heterogenous effects also calls for careful consideration of regularization applied to the treatment effect function, which we consider in the next section.

5. Regularization for heterogeneous treatment effects: Bayesian causal forests. In much the same way that a direct BART prior on f does not allow careful handling of confounding, it also does not allow separate control over the discovery of heterogeneous effects because there is no explicit control over how f varies in Z. Our solution to this problem is a simple re-parametrization that avoids the indirect specification of the prior over the treatment effects:

(11)
$$f(\mathbf{x}_i, z_i) = \mu(\mathbf{x}_i) + \tau(\mathbf{x}_i)z_i.$$

This model can be thought of as a linear regression in z with covariate-dependent functions for both the slope and the intercept. Writing the model this way sacrifices nothing in terms of expressiveness, but permits independent priors to be placed on τ , which is precisely the treatment effect:

(12)
$$E(Y_i \mid x_i, Z_i = 1) - E(Y_i \mid x_i, Z_i = 0) = \{\mu(x_i) + \tau(x_i)\} - \mu(x_i) = \tau(x_i).$$

Under this model $\mu(\mathbf{x}) = E(Y \mid Z = 0, X = x)$ is a prognostic score in the sense of ?, another interpretable quantity, to which we apply an prior distribution independent of τ (as detailed below).

Based on the observations of the previous section, we further propose specifying the model as

(13)
$$f(\mathbf{x}_i, z_i) = \mu(\mathbf{x}_i, \hat{\pi}_i) + \tau(\mathbf{x}_i)z_i,$$

where $\hat{\pi}_i$ is an estimate of the propensity score.

While we will use variants of BART priors for μ and τ (see section 5.2), this parameterization has many advantages in general, regardless of the specific priors. The most obvious advantage is that the treatment effect is an explicit parameter of the model, $\tau(x)$, and as a result we can specify an appropriate prior on it directly. Before turning to the details of our model specification, we first contrast this parameterization with two common alternatives.

5.1. Parameterizing regression models of heterogeneous effects. There are two common modeling strategies for estimating heterogeneous effects. The first we discussed above: treat z as "just another covariate" and specify a prior on $f(\mathbf{x}_i, z_i)$, e.g. as in Hill (2011). The second is to fit entirely separate models to the treatment and control data: $(Y \mid Z = z, \mathbf{x}) \sim N(f_z(\mathbf{x}_i), \sigma_z^2)$ with independent priors over the parameters in the z = 0 and z = 1 models. In this section we argue that neither approach is satisfactory and propose the model in (13) as a reasonable interpolation between the two. (See Künzel et al. (2017) for a related discussion comparing these two approaches in a non-model-based setting.)

It is instructive to consider (11) as a nonlinear regression analogue of the common strategy of parametrizing contrasts (differences) and aggregates (sums) rather than group-specific location parameters. Specifically, consider a two-group difference-in-means problem:

(14)
$$Y_{i1} \stackrel{\text{iid}}{\sim} N(\mu_1, \sigma^2)$$
$$Y_{j2} \stackrel{\text{iid}}{\sim} N(\mu_2, \sigma^2).$$

Although the above parameterization is intuitive, if the estimand of interest is $\mu_1 - \mu_2$, the implied prior over this quantity has variance strictly greater than the variances over μ_1 or μ_2 individually. This is plainly nonsensical if the analyst has no subject matter knowledge regarding the individual levels of the groups, but has strong prior knowledge that $\mu_1 \approx \mu_2$. This is common in a causal inference setting: If the data come from a randomized experiement where Y_1 constitutes a control sample and Y_2 a treated sample, then subject matter considerations will typically limit the plausible range of treatment effects $\mu_1 - \mu_2$.

The appropriate way to incorporate that kind of knowledge is simply to reparametrize as

(15)
$$Y_{i1} \stackrel{\text{iid}}{\sim} N(\mu + \tau, \sigma^2)$$
$$Y_{j2} \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2)$$

whereupon the estimand of interest becomes τ , which can be given an informative prior centered at zero with an appropriate variance. Meanwhile, μ can be given a very vague (perhaps even improper) prior.

While the nonlinear modeling context is more complex, the considerations are the same: our goal is simultaneously to let $\mu(x)$ be flexibly learned (to adequately deconfound and obtain more precise inference), while appropriately regularizing $\tau(x)$, which we expect, a priori, to be relatively small in magnitude and "simple" (minimal heterogeneity). Neither of the two more common parametrizations permit this: Independent estimation of $f_0(x)$ and $f_1(x)$ implies a highly vague prior on $\tau(x) = f_1(x) - f_0(x)$; i.e. a Gaussian process prior on each would imply a twice-as-variable Gaussian process prior on the difference, as in the simple example above. Estimation based on the single response surface f(x, z) often does not allow direct point-wise control of $\tau(x) = f(x, 1) - f(x, 0)$ at all. In particular, with a BART prior on f the induced prior on τ depends on incidental features of the data such as the size and distribution of the covariate vector x.

5.2. Prior specification. With the model parameterized as in (13), we can specify different BART priors on μ and τ . For μ we use the default suggestions in (Chipman, George and McCulloch, 2010) (200 trees, $\beta = 2$, $\eta = 0.95$), except that we place a half-Cauchy prior over the scale of the leaf parameters with prior median equal to twice the marginal standard deviation of Y (Gelman et al., 2006; Polson et al., 2012). We find that inference over τ is typically insensitive to reasonable deviations from these settings, so long as the prior is not so strong that deconfounding does not take place.

For τ , we prefer stronger regularization. First, we use fewer trees (50 versus 200), as we generally believe that patterns of treatment effect heterogeneity are relatively simple. Second, we set the depth penalty $\beta = 3$ and splitting probability $\eta = 0.25$ (instead of $\beta = 2$ and $\eta = 0.95$) to shrink more strongly toward homogeneous effects (the extreme case where none of the trees split at all corresponds to purely homogeneous effects). Finally, we replace the half-Cauchy prior over the scale of τ with a half Normal prior, pegging the prior median to the marginal standard deviation of Y.

6. Empirical evaluations. In this section, we provide a more extensive look at how BCF compares to various alternatives. In Section 6.1 we compare BCF, generalized random forests (Athey, Tibshirani and Wager, 2016), and a linear model with all three-way interactions as plausible methods for estimating heterogeneous treatment effects with measures of uncertainty. We also consider three specifications of BART: the standard response surface BART that considers the treatment variable as "just another covariate", one where separate BART models are fit to the treatment and control arms of the data, and one where an estimate of the propensity score is included as a predictor. In Section 6.2 we report on the results of two separate data analysis challenges, where the entire community was invited to submit methods for evaluation on larger synthetic datasets with heterogeneous treatment effects. In both simulation settings we find that BCF performs well under a wide range of scenarios.

In all cases the estimands of interest are either conditional average treatment effects for individual i accounting for all the variables, estimated by the posterior mean treatment effect $\hat{\tau}(\mathbf{x}_i)$, or sample subgroup average treatment effects estimated by $\sum_{i \in \mathcal{S}} \hat{\tau}(\mathbf{x}_i)$, where \mathcal{S} is the subgroup of interest. Credible intervals are computed from MCMC output.

6.1. Simulation studies. We evaluated three variants of BART, the causal random forest model of Athey, Tibshirani and Wager (2016), and a regularized linear regression with up to three way interactions. We consider eight distinct, but closely related, data generating processes, corresponding to the various combinations of toggling three two-level settings: homogeneous versus heterogeneous treatment effects, a linear versus nonlinear conditional expectation function, and two different sample sizes (n = 250 and n = 500). Five variables comprise x; the first three are continuous, drawn as standard normal random variables, the fourth is a dichotomous variable and the fifth is unordered categorical, taking three levels (denoted 1,2,3). The treatment effect is either

$$\tau(\mathbf{x}) = \begin{cases} 3, & \text{homogeneous} \\ 1 + 2x_2x_5, & \text{heterogeneous}, \end{cases}$$

the prognostic function is either

$$\mu(\mathbf{x}) = \begin{cases} 1 + g(x_4) + x_1 x_3, & \text{linear} \\ -6 + g(x_4) + 6|x_3 - 1|, & \text{nonlinear,} \end{cases}$$

where g(1) = 2, g(2) = -1 and g(3) = -4, and the propensity function is

$$\pi(\mathbf{x}_i) = 0.8\Phi(3\mu(\mathbf{x}_i)/s - 0.5x_1) + 0.05 + u_i/10$$

where s is the standard deviation of μ taken over the observed sample and $u_i \sim \text{Uniform}(0,1)$.

To evaluate each method we consider three criteria, applied to two different estimands. First, we consider how each method does at estimating the (sample) average treatment effect (ATE) according to root mean square error, coverage, and average interval length. Then, we consider the same criteria, except applied to estimates of the conditional average treatment effect (CATE), averaged over the sample. Results are based on 200 independent replications for each DGP. Results are reported in Tables 2 (for the linear DGP) and 3 (for the nonlinear DGP). The important trends are as follows:

- BCF or ps-BART benefit dramatically by explicitly protecting against RIC;
- BART- (f_0, f_1) and causal random forests both exhibit subpar performance in this simulation;
- all methods improve with a larger sample;
- BCF priors are especially helpful at the smaller sample size (when estimation is more difficult);
- the linear model dominates when correct, but fares extremely poorly when wrong;
- BCF's improvements over ps-BART are more pronounced in the nonlinear DGP;
- BCF's average interval length is notably smaller than the ps-BART interval, usually (but not always) with comparable coverage.

Table 2
Simulation study results when the true DGP is a linear model with third order interactions. Root mean square estimation error (rmse), coverage (cover) and average interval length (len) are reported for both the average treatment effect (ATE) estimates and the conditional average treatment effect estimates (CATE).

Homogeneous effect Heterog									erogene	geneous effects			
\overline{n}	Method		ATE			CATE			ATE		CATE		
		rmse	cover	len	rmse	cover	len	rmse	cover	len	rmse	cover	len
	BCF	0.21	0.92	0.91	0.48	0.96	2.0	0.27	0.84	0.99	1.09	0.91	3.3
250	ps-BART	0.22	0.94	0.97	0.44	0.99	2.3	0.31	0.90	1.13	1.30	0.89	3.5
250	BART	0.34	0.73	0.94	0.54	0.95	2.3	0.45	0.65	1.10	1.36	0.87	3.4
	BART (f_0, f_1)	0.56	0.41	0.99	0.92	0.93	3.4	0.61	0.44	1.14	1.47	0.90	4.5
	Causal RF	0.34	0.73	0.98	0.47	0.84	1.3	0.49	0.68	1.25	1.58	0.68	2.4
	LM + HS	0.14	0.96	0.83	0.26	0.99	1.7	0.17	0.94	0.89	0.33	0.99	1.9
	BCF	0.16	0.88	0.60	0.38	0.95	1.4	0.16	0.90	0.64	0.79	0.89	2.4
500	ps-BART	0.18	0.86	0.63	0.35	0.99	1.8	0.16	0.90	0.69	0.86	0.95	2.8
300	BART	0.27	0.61	0.61	0.42	0.95	1.8	0.25	0.76	0.67	0.88	0.94	2.8
	BART (f_0, f_1)	0.47	0.21	0.66	0.80	0.93	3.1	0.42	0.42	0.75	1.16	0.92	3.9
	Causal RF	0.36	0.47	0.69	0.52	0.75	1.2	0.40	0.59	0.88	1.30	0.71	2.1
	LM + HS	0.11	0.96	0.54	0.18	0.99	1.0	0.12	0.93	0.59	0.22	0.98	1.2

Table 3 Simulation study results when the true DGP is nonlinear. Root mean square estimation error (rmse), coverage (cover) and average interval length (len) are reported for both the average treatment effect (ATE) estimates and the conditional average treatment effect estimates (CATE).

Homogeneous effect Heterogeneous effects								ets					
n	Method	ATE			CATE		ATE			CATE			
		rmse	cover	len	rmse	cover	len	rmse	cover	len	rmse	cover	len
	BCF	0.26	0.945	1.3	0.63	0.94	2.5	0.30	0.930	1.4	1.3	0.93	4.5
250	ps-BART	0.54	0.780	1.6	1.00	0.96	4.3	0.56	0.805	1.7	1.7	0.91	5.4
250	BART	0.84	0.425	1.5	1.20	0.90	4.1	0.84	0.430	1.6	1.8	0.87	5.2
	BART (f_0, f_1)	1.48	0.035	1.5	2.42	0.80	6.4	1.44	0.085	1.6	2.6	0.83	7.1
	Causal RF	0.81	0.425	1.5	0.84	0.70	2.0	1.10	0.305	1.8	1.8	0.66	3.4
	LM + HS	1.77	0.015	1.8	2.13	0.54	4.4	1.65	0.085	1.9	2.2	0.62	4.8
	BCF	0.20	0.945	0.97	0.47	0.94	1.9	0.23	0.910	0.97	1.0	0.92	3.4
E00	ps-BART	0.24	0.910	1.07	0.62	0.99	3.3	0.26	0.890	1.06	1.1	0.95	4.1
500	BART	0.31	0.790	1.00	0.63	0.98	3.0	0.33	0.760	1.00	1.1	0.94	3.9
	BART (f_0, f_1)	1.11	0.035	1.18	2.11	0.81	5.8	1.09	0.065	1.17	2.3	0.82	6.2
	Causal RF	0.39	0.650	1.00	0.54	0.87	1.7	0.59	0.515	1.18	1.5	0.73	2.8
	LM + HS	1.76	0.005	1.34	2.19	0.40	3.5	1.71	0.000	1.34	2.2	0.45	3.7

6.2. Atlantic causal inference conference data analysis challenges. The Atlantic Causal Inference Conference (ACIC) has featured a data analysis challenge since 2016. Participants are given a large number of synthetic datasets and invited to submit their estimates of treatment effects along with confidence or credible intervals where available. Specifically, participants were asked to produce estimates and uncertainty intervals for the sample average treatment effect on the treated, as well as conditional average treatment effects for each unit. Methods were evaluated based on a range of criteria including estimation error and coverage of uncertainty intervals. The datasets and ground truths are publicly available, so while BCF was not entered into either the 2016 or 2017 competitions we can benchmark its performance against a suite of methods that we did not choose, design, or implement.

6.2.1. ACIC 2016 competition. The 2016 contest design, submitted methods, and results are summarized in Dorie et al. (2017). Based on an early draft of our manuscript Dorie et al. (2017) also evaluated a version of BART that included an estimate of the propensity score, which was one of the top methods on bias and RMSE for estimating the sample ATT. BART with the propensity score outperformed BART without the propensity score on bias, RMSE, and coverage for the SATT, and was a leading method overall.

Therefore, rather than include results for all 30 methods here we simply include BART and ps-BART as leading contenders for estimating heterogeneous treatment effects in this setting. Using the publicly-available competition datasets (Dorie and Hill, 2017) we implemented two additional methods: BCF and causal random forests as implemented in the R package grf (Athey, Tibshirani and Wager, 2016), using 4,000 trees to obtain confidence intervals for conditional average treatment effects and a doubly robust estimator for the SATT (as suggested in the package documentation).

Table 4 collects the results of our methods (ps-BART and BCF) as well as BART and causal random forests. Causal random forests performed notably worse than BART-based methods on every metric. BCF performed best in terms of estimation error for CATE and SATT, as measured by bias and absolute bias. While the differences in the various metrics are relatively small compared to their standard deviation across the 7,700 simulated datasets, nearly all the pairwise differences between BCF and the other methods are statistically significant as measured by a permutation test (Table 5). The sole exception is the test for a difference in bias between ps-BART and BCF, suggesting the presence of RIC in at least some of the simulated datasets. This is especially notable since the datasets were not intentionally simulated to include targeted selection.

Dorie et al. (2017) note that all submitted methods were "somewhat disappointing" in inference for the SATT (i.e., few methods had coverage near the nominal rate with reasonably sized intervals). However, ps-BART did relatively well, 88% coverage of a 95% credible interval and one of the smallest interval lengths. ps-BART had slightly better coverage than BCF (88% versus 82%), with an average interval length that was 45% larger than BCF. Vanilla BART and BCF had similar coverage rates, but BART's interval length was about 55% larger than BCF. Dorie et al. (2017) found that TMLE-based adjustments could improve the coverage of BART-based estimates of the SATT at significant computational cost; we expect that similar benefits would accrue using BCF with a TMLE adjustment, but obtaining valid confidence intervals for SATT is not our focus so we did not pursue this further.

Table 4
Abbreviated ACIC 2016 contest results. Coverage and average interval length are reported for nominal 95% uncertainty intervals. Bias and |Bias| are average bias and average absolute bias, respectively, over the. PEHE is the average precision in estimating heterogeneous treatment effects (the average root mean squared error of CATE estimates for each unit in a dataset) (Hill, 2011).

	Coverage	Interval Length	Bias	(SD)	Bias	(SD)	PEHE	(SD)
BCF	0.82	0.026	-0.0009	(0.01)	0.008	0.010	0.33	0.18
ps-BART	0.88	0.038	-0.0011	(0.01)	0.010	0.011	0.34	0.16
BART	0.81	0.040	-0.0016	(0.02)	0.012	0.013	0.36	0.19
Causal RF	0.58	0.055	-0.0155	(0.04)	0.029	0.027	0.45	0.21

Table 5

Tests and estimates for differences between BCF and other methods in the ACIC 2016 competition. The p-values are from bootstrapp permutation tests with 100,000 replicates.

	Diff Bias	p	Diff Bias	p	Diff PEHE	p
ps-BART	-0.00020	0.146	0.0011	$< 1e^{-4}$	0.010	$< 1e^{-4}$
BART	-0.00070	$< 1e^{-4}$	0.0031	$< 1e^{-4}$	0.037	$< 1e^{-4}$
Causal RF	-0.01453	$< 1e^{-4}$	0.0204	$< 1e^{-4}$	0.125	$< 1e^{-4}$

6.2.2. ACIC 2017 competition. The ACIC 2017 competition was designed to have average treatment effects that were smaller, with heterogenous treatment effects that were less variable, relative to the 2016 datasets. Arguably, the 2016 competition included many datasets with unrealistically large average treatment effects and similarly unrealistic degrees of heterogeneity¹. Additionally, the 2017 competition explicitly incorporated targeted selection (unlike the 2016 datasets). The ACIC

 $^{^{1}}$ Across the 2016 competition datasets, the interquartile range of the SATT was 0.57 to 0.79 in standard deviations of Y, with a median value of 0.68. The standard deviation of the conditional average treatment effects for the sample units had an interquartile range of 0.24 to 0.93, again in units of standard deviations of Y. A significant fraction of the variability in Y was explained by heterogeneous treatment effects in a large number of the simulated datasets.

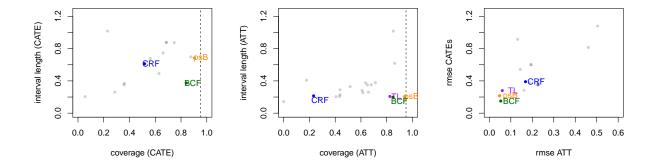


FIG 6. Each data point represents one method. ps-BART (psB, in orange) was submitted by a group independent of the authors based on a draft of this manuscript. TL (purple) is a TMLE-based submission that performed will for estimating SATT, but did not furnish estimates of conditional average treatment effects. BCF (green) and causal random forests (CRF, blue) were not part of the original contest. For descriptions of the other methods refer to Hahn, Dorie and Murray (2018).

2017 competition design and results are summarized completely in Hahn, Dorie and Murray (2018); here we report selected results for the datasets with independent additive errors.

Figure 6.2.2 contains the results of the 2017 competition. The patterns here are largely similar to the 2016 competition, despite some stark differences in the generation of synthetic datasets. ps-BART and BCF have the lowest estimation error for CATE and SATE. The closest competitor on estimation error was a TMLE-based approach. We also see that ps-BART edges BCF slightly in terms of coverage once again, although BCF has much shorter intervals. Causal random forests does not perform well, with coverage for SATT and CATE far below the nominal rate.

7. Application: The effect of smoking on medical expenditures.

7.1. Background and data. As an empirical demonstration of the Bayesian causal forest model, we consider the question of how smoking affects medical expenditures. This question is of interest as it relates to lawsuits against the tobacco industry. The lack of experimental data speaking to this question motivates the reliance on observational data. This question has been studied in several previous papers; see Zeger et al. (2000) and references therein. Here, we follow Imai and Van Dyk (2004) in analyzing data extracted from the 1987 National Medical Expenditure Survey (NMES) by Johnson et al. (2003). The NMES records many subject-level covariates and boasts third-party-verified medical expenses. Specifically, our regression includes the following ten patient attributes:

• age: age in years at the time of the survey

- smoke_age: age in years when the individual started smoking
- gender: male or female
- race: other, black or white
- marriage_status: married, widowed, divorced, separated, never married
- education_level: college graduate, some college, high school graduate, other
- census_region: geographic location, Northeast, Midwest, South, West
- poverty_status: poor, near poor, low income, middle income, high income
- seat_belt: does patient regularly use a seat belt when in a car
- years_quit: how many years since the individual quit smoking.

The response variable is the natural logarithm of annual medical expenditures, which makes the normality of the errors more plausible. Under this transformation, the treatment effect corresponds to a multiplicative effect on medical expenditure. Following Imai and Van Dyk (2004), we restrict our analysis to smokers who had non-zero medical expenditure. Our treatment variable is an indicator of heavy lifetime smoking, which we define to be greater than 17 pack-years, the equivalent of 17 years of pack-a-day smoking. See again Imai and Van Dyk (2004) for more discussion of this variable. We scrutinize the overlap assumption and exclude individuals younger than 28 on the grounds that it is improbable for someone that young to have achieved this level of exposure. After making these restrictions, our sample consists of n = 6,798 individuals.

7.2. Results. Here, we highlight the differences that arise when analyzing this data using standard BART versus using BCF. First, the estimated expected responses from the two models have correlation of 0.98, so that the two models concur on the nonlinear prediction problem. This suggests that, as was intended, BCF will inherit BART's outstanding predictive capabilities. By contrast, the estimated individual treatment effects are only correlated 0.70. The most notable differences between these CATE estimates is that the BCF estimates exhibit a strong trend in the age variable, as shown in Figure 7.2; the BCF estimates suggest that smoking has an pronounced impact on the health expenditures of younger people.

Despite a wider range of values in the CATE estimates (due largely to the inferred trend in the age variable), the ATE estimate of BCF is notably lower than that of BART, the posterior 95% credible intervals being translated by 0.05, (0.00, 0.20) for BCF vs (0.05, 0.25) for BART. The higher estimate of BART is possibly a result of RIC. Figure 7.2 shows a LOWESS trend between

the estimated propensity and prognostic scores (from BCF); the monotone trend is suggestive of targeted selection (high medical expenses are predictive of heavy smoking) and hints at the possibility of RIC-type inflation of the BART ATE estimate (compare to Figures 2 and 4).

Although the vast majority of individual treatment effect estimates are statistically uncertain, as reflected in posterior 95% credible intervals that contain zero (Figure 7.2), the evidence for subgroup heterogeneity is relatively strong, as uncovered by the following posterior exploration strategy. First, we grow a parsimonious regression tree to the point estimates of the individual treatment effects (using the rpart package in R); see the left panel of Figure 7.2. Then, based on the candidate subgroups revealed by the regression summary tree, we plot a posterior histogram of the difference between any two covariate-defined subgroups. The right panel of Figure 7.2 shows the posterior distribution of the difference between men younger than 46 and women over 66; virtually all of the posterior mass is above zero, suggesting that the treatment effect of heavy smoking is discernibly different for these two groups, with young men having a substantially higher estimated subgroup ATE. This approach, although somewhat informal, is a method of exploring the posterior distribution and, as such, any resulting summaries are still valid Bayesian inferences. Moreover, such Bayesian "fit-the-fit" posterior summarization strategies can be formalized from a decision theoretic perspective (Sivaganesan, Müller and Huang, 2017; Hahn and Carvalho, 2015); we do not explore this possibility further here.

From the above we conclude that how a model treats the age variable would seem to have an outsized impact on the way that predictive patterns are decomposed into treatment effect estimates based on this data, as age plausibly has prognostic, propensity and moderating roles simultaneously. Although it is difficult to trace the exact mechanism by which it happens, the BART model clearly de-emphasizes the moderating role, whereas the BCF model is designed specifically to capture such trends. Possible explanations for the age heterogeneity could be a mixed additive-multiplicative effect combined with higher baseline expenditures for older individuals or possibly survivor bias (as also mentioned in Imai and Van Dyk (2004)), but further speculation is beyond the scope of this analysis.

8. Discussion. We have demonstrated the utility of our new model for the estimation of conditional average treatment effects. We conclude by reviewing the contributions made here and positioning them in the existing literature.

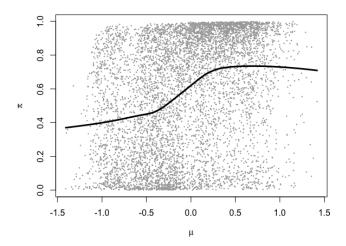


FIG 7. Each gray dot depicts the estimated propensity and prognostic scores for an individual. The solid bold line depicts a LOESS trend fit to these points; the monotonicity is suggestive of targeted selection. Compare to Figures 2 and 4.

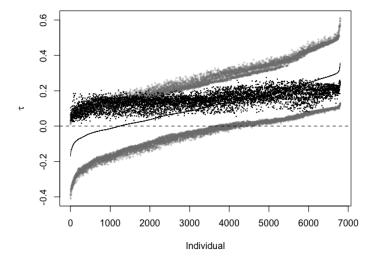


FIG 8. Point estimates of individual treatment effects are shown in black. The smooth line depicts the estimates from BCF, which are ordered from smallest to largest. The unordered points represent the corresponding ITE estimates from BART. Note that the BART estimates seem to be higher, on average, than the BCF estimates. The upper and lower gray dots correspond to the posterior 95% credible interval end points associated with the BCF estimates; most ITE intervals contain zero, especially those with smaller (even negative) point estimates.

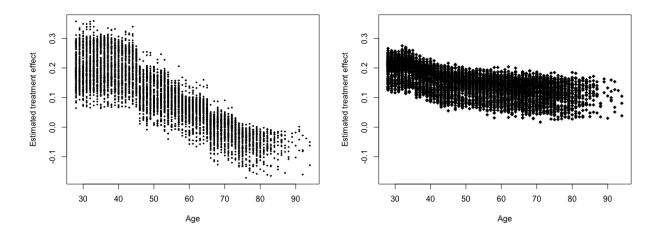


FIG 9. Each point depicts the estimated treatment effect for an individual. The BCF model (left panel) detects pronounced heterogeneity moderated by the age variable, whereas BART (right panel) does not.

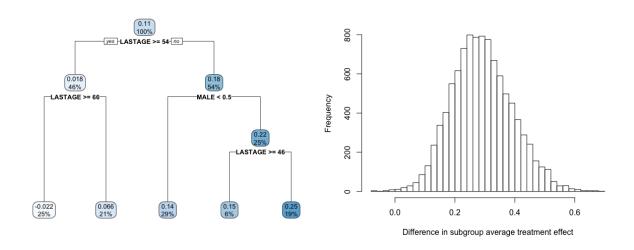


FIG 10. Left panel: a summarizing regression tree fit to posterior point estimates of individual treatment effects. The top number in each box is the average subgroup treatment effect in that partition of the population, the lower number shows the percentage of the total sample constituting that subgroup. Age and gender are flagged as important moderating variables. Right panel: based on the tree in the left panel, we consider the difference in treatment effects between men younger than 46 and women older than 66; a posterior histogram of this difference shows that nearly all of the posterior mass is above zero, indicating that these two subgroups are discernibly different, with young men having substantially higher subgroup average treatment effect.

8.1. Zellner priors for non- and semiparametric Bayesian causal inference. In Section 4 we showed that the current gold standard in nonparametric Bayesian regression models for causal inference (BART) is susceptible to regression induced confounding as described by Hahn et al. (2016). The solution we propose is to include an estimate of the propensity score as a covariate in the outcome model. This induces a prior distribution that treats Z_i and $\hat{\pi}_i$ equitably, discouraging the outcome model from erroneously attributing the effect of confounders to the treatment variable. Here we justify and collect arguments in favor of this approach. We discuss an argument against, namely that it does not incorporate uncertainty in the propensity score, in a later subsection.

Conditioning on an estimate of the propensity score is readily justified: Because our regression model is conditional on Z and X, it is perfectly legitimate to condition our prior on them as well. This approach is widely used in linear regression, the most common example being Zellner's g-prior (Zellner, 1986) which parametrizes the prior covariance of a vector of regression coefficients in terms of a plug-in estimate of the predictor variables' covariance matrix. Nodding to this heritage, we propose to call general predictor-dependent priors "Zellner priors".

In the Bayesian causal forest model, we specify a prior over f by applying an independent BART prior that includes $\hat{\pi}(\mathbf{x}_i)$ as one of its splitting dimensions. That is, because $\hat{\pi}(\mathbf{x}_i)$ is a fixed function of \mathbf{x}_i , f is still a function $f:(\mathcal{X},\mathcal{Z}) \mapsto \mathbb{R}$; the inclusion of $\hat{\pi}(\mathbf{x}_i)$ among the splitting dimensions does not materially change the support of the prior, but it does alter which functions are deemed more likely. Therefore, although writing $f(\mathbf{x}_i, z_i, \hat{\pi}(\mathbf{x}_i))$ is suggestive of how the prior is implemented in practice, we prefer notation such as

(16)
$$Y_i = f(\mathbf{x}_i, z_i) + \epsilon_i, \quad \epsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2),$$
$$f \sim \text{BART}(\mathbf{X}, Z, \hat{\pi}),$$

where $\hat{\pi}$ is itself a function of (\mathbf{X}, Z) . Viewing BART as a prior in this way highlights the fact that various transformations of the data could be computed beforehand, prior to fitting the data with the default BART priors; the choice of transformations will control the nature of the regularization that is imposed. In conventional predictive modeling there is often no particular knowledge of which transformations of the covariates might be helpful. However, in the treatment effect context the propensity score is a natural and, in fact, critical choice.

8.2. Why not use only the propensity score? vs. Why use the propensity score at all?. It has long been recognized that regression on the propensity score is a useful dimension reduction tactic

(Rosenbaum and Rubin, 1983). For the purpose of estimating average treatment effects, a regression model on the one-dimensional propensity score is sufficient for the task, allowing one to side-step estimating high dimensional nuisance parameters. In our notation, if π is assumed known, then one need only infer $f(\pi)$. That said, there are several reasons one should include the control vector \mathbf{x}_i in its entirety (in addition to the propensity score).

The first reason is pragmatic: If one wants to identify heterogeneous effects, one needs to include any potential effect moderating variables anyway, precluding any dimension reduction at the outset.

Second, if we are to take a conditionally-iid Bayesian regression approach to inference and we do not in fact believe the response to depend on \mathbf{X} strictly through the propensity score, we simply must include the covariates themselves and model the conditional distribution $p(Y \mid Z, \mathbf{X})$ (otherwise the error distribution is highly dependent, integrated across \mathbf{X}). The justification for making inference about average treatment effects using regression or stratification on the propensity score alone is entirely frequentist; this approach is not without its merits, and we do not intend to argue frequency calibration is not desirable, but a fully Bayesian approach has its own appeal.

Third, if our propensity score model is inadequate (misspecified or otherwise poorly estimated), including the full predictor vector allows for the possibility that the response surface model remains correctly specified.

The converse question, Why bother with the propensity score if one is doing a high dimensional regression anyway?, has been answered in the main body of this paper. Incorporating the propensity score (or another balancing score) yields a prior that can more readily adapt to complex patterns of confounding. In fact, in the context of response surface modeling for causal effects, failing to include an estimate of the propensity score (or another balancing score) can lead to additional bias in treatment effect estimates, as shown by the simple, low-dimensional example in Section 4.

8.3. Why not joint response-treatment modeling and what about uncertainty in the propensity score?. Using a presumptive model for Z to obtain $\hat{\pi}$ invites the suggestion of fitting a joint model for (Y, Z). Indeed, this is the approach taken in Hahn et al. (2016) as well as earlier papers, including Rosenbaum and Rubin (1983), Robins, Mark and Newey (1992), McCandless, Gustafson and Austin (2009), Wang, Parmigiani and Dominici (2012), and Zigler and Dominici (2014). While this approach is certainly reasonable, the Zellner prior approach would seem to afford all the same benefits while avoiding the distorted inferences that would result from a joint model if the propensity score model is misspecified (Zigler and Dominici, 2014).

One might argue that our Zellner prior approach gives under-dispersed posterior inference in the sense that it fails to account for the fact that $\hat{\pi}$ is simply a point estimate (and perhaps a bad one). However, this objection is somewhat misguided. First, as discussed elsewhere (e.g. Hill (2011)), inference on individual or subgroup treatment effects follows directly from the conditional distribution $(Y \mid Z, \mathbf{X})$. To continue our analogy with the more familiar Zellner g-prior, to model $(Y \mid Z, \mathbf{X})$ we are no more obligated to consider uncertainty in $\hat{\pi}$ than we are to consider uncertainty in $(\mathbf{X}'\mathbf{X})^{-1}$ when using a g-prior for on the coefficients of a linear model. Second, $\hat{\pi}$ appears in the model along with the full predictor vector \mathbf{x} : it is provided as a hint, not as a certainty, and this model is at least as capable of estimating a complex response surface as the corresponding model without $\hat{\pi}$, and the cost incurred by the addition of a one additional "covariate" can be more than offset by the bias reduction in the estimation of treatment effects.

On the other hand, we readily acknowledge that one might be interested in what inferences would obtain if we used different $\hat{\pi}$ estimates. One might consider fitting a series of BCF models with different estimates of $\hat{\pi}$, perhaps from alternative models or other procedures. This is a natural form of sensitivity analysis in light of the fact that the adjustments proposed in this paper only work if $\hat{\pi}$ accurately approximates π . However, it is worth noting that the available (z, x) data speak to this question: a host of empirically proven prediction methods (i.e. neural networks, support vector machines, random forests, boosting, or any ensemble method) can be used to construct candidate $\hat{\pi}$ and cross-validation may be used to gauge their accuracy. Only if a "tie" in generalization error (predicting Z) is encountered must one turn to sensitivity analysis.

8.4. Connections to doubly robust estimation. Our combination of propensity score estimation and outcome modeling is superficially reminiscent of doubly robust estimation (Bang and Robins, 2005), where propensity score and outcome regression models are combined to yield consistent estimates of finite dimensional treatment effects, provided at least one model is correctly specified. We do not claim our approach is doubly robust, however, and in all of our examples above we use the natural Bayesian estimates of (conditional) average treatment effects rather than doubly robust versions.

The motivation behind our approach is in fact quite different than that behind doubly robust estimation: RIC is fundamentally a finite sample phenomenon, and including the estimated propensity score is an effort to improve the finite sample performance of the estimated outcome regression model. The two approaches are complementary – if the outcome regression is more accurate in finite

samples, similar benefits should accrue to a doubly robust estimator computed from the outcome regression estimates.

8.5. Related non-Bayesian work. The Bayesian causal forest model is a flexible semi-parametric prediction model for estimating causal effects. Other recent work also occupies this intersection between "machine learning" and causal inference, each with a somewhat different focus. Targeted maximum likelihood estimation (TMLE) (van der Laan, 2010a,b), double machine learning (Chernozhukov et al., 2016), and generalized boosting (McCaffrey, Ridgeway and Morral, 2004; McCaffrey et al., 2013) all focus on estimation of average treatment effects, whereas our focus is on individual (heterogeneous, subgroup) treatment effects. Like us, Taddy et al. (2016) focuses on estimating heterogeneous effects, but they analyze data from experiments, whereas our data are observational. As we have discussed, this has significant implications for how we should specify prior distributions. Other related contributions include Su et al. (2012) and Lu et al. (2017). Finally, Wager and Athey (2015) prove theoretical results for random forest-based estimation of heterogeneous effects from observational data. However, Wager and Athey (2015) offer little practical guidance on how to deploy their method in practice, in particular on how much regularization to impose, even in the newest iteration of the method as reported in Athey, Tibshirani and Wager (2016). Moreover, as we saw in the simulation studies and the ACIC challenge results, this method offers inferior empirical performance, a finding that is also corroborated in Wendling et al. (2018). Given this landscape, we believe the Bayesian causal forest model presented in this paper represents a beneficial new tool for causal inference from observational data, especially when confounding is suspected to be strong and the general magnitude of treatment effects is thought to be relatively modest.

References.

ATHEY, S., TIBSHIRANI, J. and WAGER, S. (2016). Generalized Random Forests. arXiv preprint arXiv:1610.01271. Bang, H. and Robins, J. M. (2005). Doubly robust estimation in missing data and causal inference models. Biometrics 61 962–973.

Chernozhukov, V., Chetverikov, D., Demirer, M., Duflo, E., Hansen, C. et al. (2016). Double machine learning for treatment and causal parameters. arXiv preprint arXiv:1608.00060.

CHIPMAN, H. A., GEORGE, E. I. and McCulloch, R. E. (1998). Bayesian CART model search. *Journal of the American Statistical Association* **93** 935–948.

Chipman, H. A., George, E. I. and McCulloch, R. E. (2010). BART: Bayesian additive regression trees. *The Annals of Applied Statistics* 266–298.

- DORIE, V. and HILL, J. (2017). aciccomp2016: Atlantic Causal Inference Conference Competition 2016 Simulation R package version 0.1-0.
- DORIE, V., HILL, J., SHALIT, U., SCOTT, M. and CERVONE, D. (2017). Automated versus do-it-yourself methods for causal inference: Lessons learned from a data analysis competition. arXiv preprint arXiv:1707.02641.
- Gelman, A. et al. (2006). Prior distributions for variance parameters in hierarchical models (comment on article by Browne and Draper). *Bayesian Analysis* 1 515–534.
- Giles, D. and Rayner, A. (1979). The mean squared errors of the maximum likelihood and natural-conjugate Bayes regression estimators. *Journal of Econometrics* 11 319–334.
- Gramacy, R. B. and Lee, H. K. (2008). Bayesian treed Gaussian process models with an application to computer modeling. *Journal of the American Statistical Association* **103**.
- Green, D. P. and Kern, H. L. (2012). Modeling heterogeneous treatment effects in survey experiments with Bayesian additive regression trees. *Public opinion quarterly* nfs036.
- Gustafson, P. and Greenland, S. (2006). Curious phenomena in Bayesian adjustment for exposure misclassification. Statistics in Medicine 25 87–103.
- Hahn, P. R. and Carvalho, C. M. (2015). Decoupling shrinkage and selection in Bayesian linear models: a posterior summary perspective. *Journal of the American Statistical Association* **110** 435–448.
- Hahn, P. R., Dorie, V. and Murray, J. S. (2018). Atlantic Causal Inference Conference (ACIC) Data Analysis Challenge 2017.
- Hahn, P. R., Puelz, D., He, J. and Carvalho, C. M. (2016). Regularization and confounding in linear regression for treatment effect estimation. *Bayesian Analysis*.
- HECKMAN, J. J., LOPES, H. F. and PIATEK, R. (2014). Treatment effects: A Bayesian perspective. *Econometric reviews* 33 36–67.
- HILL, J. L. (2011). Bayesian nonparametric modeling for causal inference. Journal of Computational and Graphical Statistics 20.
- HILL, J., Su, Y.-S. et al. (2013). Assessing lack of common support in causal inference using Bayesian nonparametrics: Implications for evaluating the effect of breastfeeding on children's cognitive outcomes. The Annals of Applied Statistics 7 1386–1420.
- IMAI, K. and VAN DYK, D. A. (2004). Causal inference with general treatment regimes: Generalizing the propensity score. *Journal of the American Statistical Association* **99** 854–866.
- Imbens, G. W. and Rubin, D. B. (2015). Causal Inference in Statistics, Social, and Biomedical Sciences. Cambridge University Press.
- Johnson, E., Dominici, F., Griswold, M. and Zeger, S. L. (2003). Disease cases and their medical costs attributable to smoking: an analysis of the national medical expenditure survey. *Journal of Econometrics* **112** 135–151.
- KERN, H. L., STUART, E. A., HILL, J. and GREEN, D. P. (2016). Assessing methods for generalizing experimental impact estimates to target populations. *Journal of Research on Educational Effectiveness* 9 103–127.
- KÜNZEL, S., SEKHON, J., BICKEL, P. and Yu, B. (2017). Meta-learners for Estimating Heterogeneous Treatment Effects using Machine Learning. arXiv preprint arXiv:1706.03461.

- LI, M. and Tobias, J. L. (2014). Bayesian analysis of treatment effect models. In *Bayesian Inference in the Social Sciences* (I. Jeliazkov and X.-S. Yang, eds.) 3, 63–90. Wiley.
- LINERO, A. R. and YANG, Y. (2017). Bayesian Regression Tree Ensembles that Adapt to Smoothness and Sparsity. arXiv preprint arXiv:1707.09461.
- Lu, M., Sadiq, S., Feaster, D. J. and Ishwaran, H. (2017). Estimating Individual Treatment Effect in Observational Data Using Random Forest Methods. arXiv preprint arXiv:1701.05306.
- MCCAFFREY, D. F., RIDGEWAY, G. and MORRAL, A. R. (2004). Propensity score estimation with boosted regression for evaluating causal effects in observational studies. *Psychological Methods* **9** 403.
- McCaffrey, D. F., Griffin, B. A., Almirall, D., Slaughter, M. E., Ramchand, R. and Burgette, L. F. (2013). A tutorial on propensity score estimation for multiple treatments using generalized boosted models. *Statistics in Medicine* **32** 3388–3414.
- McCandless, L. C., Gustafson, P. and Austin, P. C. (2009). Bayesian propensity score analysis for observational data. *Statistics in Medicine* **28** 94–112.
- Murray, J. S. (2017). Log-Linear Bayesian Additive Regression Trees for Categorical and Count Responses. arXiv preprint arXiv:1701.01503.
- Polson, N. G., Scott, J. G. et al. (2012). On the half-Cauchy prior for a global scale parameter. *Bayesian Analysis* 7 887–902.
- ROBINS, J. M., MARK, S. D. and NEWEY, W. K. (1992). Estimating exposure effects by modelling the expectation of exposure conditional on confounders. *Biometrics* 479–495.
- ROCKOVA, V. and VAN DER PAS, S. (2017). Posterior Concentration for Bayesian Regression Trees and their Ensembles. arXiv preprint arXiv:1708.08734.
- ROSENBAUM, P. R. and RUBIN, D. B. (1983). The central role of the propensity score in observational studies for causal effects. *Biometrika* 41–55.
- SIVAGANESAN, S., MÜLLER, P. and HUANG, B. (2017). Subgroup finding via Bayesian additive regression trees. Statistics in Medicine.
- Su, X., Kang, J., Fan, J., Levine, R. A. and Yan, X. (2012). Facilitating score and causal inference trees for large observational studies. *Journal of Machine Learning Research* 13 2955–2994.
- TADDY, M., GARDNER, M., CHEN, L. and DRAPER, D. (2016). A Nonparametric Bayesian Analysis of Heterogenous Treatment Effects in Digital Experimentation. *Journal of Business & Economic Statistics* **34** 661–672.
- VAN DER LAAN, M. J. (2010a). Targeted maximum likelihood based causal inference: Part I. The International Journal of Biostatistics 6.
- VAN DER LAAN, M. J. (2010b). Targeted maximum likelihood based causal inference: Part II. The International Journal of Biostatistics 6.
- Wager, S. and Athey, S. (2015). Estimation and Inference of Heterogeneous Treatment Effects using Random Forests. arXiv preprint arXiv:1510.04342.
- Wang, C., Parmigiani, G. and Dominici, F. (2012). Bayesian effect estimation accounting for adjustment uncertainty. *Biometrics* **68** 661–671.
- WENDLING, T., JUNG, K., CALLAHAN, A., SCHULER, A., SHAH, N. and GALLEGO, B. (2018). Comparing methods

- for estimation of heterogeneous treatment effects using observational data from health care databases. *Statistics in Medicine*.
- Yang, Y., Cheng, G. and Dunson, D. B. (2015). Semiparametric Bernstein-von Mises Theorem: Second Order Studies. arXiv preprint arXiv:1503.04493.
- ZEGER, S. L., WYANT, T., MILLER, L. S. and SAMET, J. (2000). Statistical testimony on damages in Minnesota v. Tobacco Industry. In *Statistical Science in the Courtroom* 303–320. Springer.
- Zellner, A. (1986). On assessing prior distributions and Bayesian regression analysis with g-prior distributions.

 Bayesian inference and decision techniques: Essays in Honor of Bruno De Finetti 6 233–243.
- ZIGLER, C. M. and DOMINICI, F. (2014). Uncertainty in propensity score estimation: Bayesian methods for variable selection and model-averaged causal effects. *Journal of the American Statistical Association* **109** 95–107.