

PAPER TITLE Blah Blah
Estimating Reaction to Welfare
Using Causal Forests

ECMA 31330 Final Project

Spring 2021

Max Bronckers
Veronica Song
Dustin Zhang

1 Introduction

There has been an increasing adoption of Machine Learning (ML) methods in economics for causal inference. While initially ML methods were avoided due to an uncertainty in their consistency, normality, and efficiency, major developments in methodology have allowed a stable large-sample confidence interval to be constructed around treatment effect estimates conditional on multiple covariates—leading to the broader adoption of these methods. [1]

One such method is Bayesian Additive Regression Trees (BART), a non-parametric method that models heterogeneous treatment effects flexibly by building on the concept of ensembles of trees. Using a Markov Chain Monte Carlo (MCMC) algorithm that derives effects from the posterior mean and interval instead of pre-specified tree parameters, BART has a much smoother and adaptive structure than traditional OLS or single tree models popular in economics, and is also resilient to problems with overfitting. [6] However, despite its excellent predictive capacity, BART lacks an asymptotic explanation for its estimates and is not guaranteed to converge in polynomial time [3], weakening its effectiveness as inference tool. Though recent work by Ročková and Saha (2018) suggests modifications to BART that may allow asymptotic concentration of the posterior mean around the true mean, the construction of an asymptotic theory of BART estimates is still an ongoing effort. [9] An alternative method is the Causal Forest (CF), proposed by Wager and Athey (2017). One of the advantages of CF is that their estimates are asymptotically gaussian and unbiased, allowing proper confidence intervals to be constructed around the treatment effect. The construction of adequate confidence intervals is especially relevant in policy applications, as consistent estimates can be produced for the treatment group. [3]

In this paper, we evaluate Causal Forests as a method for heterogeneous treatment effect estimation on both empirical and simulated datasets. Specifically, we compare the CF method to BART in the estimation of heterogeneous treatment effects in a survey experiment on welfare opinions from the General Social Survey as used by Green and Kern. Green and Kern use the GSS survey dataset with individual perception on public spending. To find the impact of the question phrasing on the responses, they use BART to estimate heterogeneous treatment effects by conditioning on a suite of socioeconomic backgrounds of the respondent. Given the theoretical shortcomings of BART and the theoretical benefits of CFs, we apply CFs to same problem and empirical dataset and compare it to the authors’ findings.

In section 2, we fit a CF model on the welfare data and compares our estimates against BART estimates of conditional average treatment effect (CATE) obtained by Green and Kern. Since empirical data offers no ground truth CATEs, we are confined to comparing the methods on interval length and RMSE. In section 3, we remedy this shortcoming with an evaluation of the CF method applied to DGPs that attempt to represent the empirical dataset. Using the ground truth ATE and CATEs of our DGPs, we assess the CF’s performance under a variety of different DGP parameters and assess the conditions under which CFs perform well or poorly.

2 Causal Forest Estimates of CATE in Welfare Dataset

2.1 Discussion of data and motivation for research

In our analysis below, we use a survey experiment from GSS which is also utilized by Green and Kern to investigate interactions between treatment and covariates that may lead to treatment effect heterogeneity. The experiment was conducted in the mid-1980s by GSS to study the negative sentiment Americans carry toward government programs labeled as "welfare". Due to associations with racial connotations and poorly managed welfare programs, respondents were found much more likely to endorse government spending for "the poor" than for public "welfare". [8]

Using BART, Green and Kern illustrates the extent to which such reaction to the question wording as "welfare" varies based on the respondents' background characteristics such as years of education, race, or political alignment. BART has been a popular choice for heterogeneous treatment effect modeling, as its estimates require little parameter tuning, allow accurate detection of interactions between covariates, and are much smoother than those of single tree models. Each individual tree in the forest has only a small effect on the model, by assuming a prior distribution over the tree parameters. [5] Such lack of any individual influential trees allow regular model fit of the BART set up. Then, a MCMC algorithm is used to sample tree parameters iteratively from the posterior distribution as the model is fit. Though this assumption of prior distribution and back-fitting algorithm allows BART to be relatively invariant across, in the presence of confounding variables and treatment effect heterogeneity, such regularization may severely bias the treatment effect estimates [7]. Moreover, BART estimates still lack theoretical explanation on its asymptotic concentration, making the construction of adequate confidence intervals challenging.

2.2 Causal forests and model assumptions

For these reasons, we choose to investigate causal forests to estimate the the treatment effect of question wording on welfare program support. Just like BART, causal forests is fit to model non-linear relationships and interactions between covariates. One advantage that causal forests carry over BART is that under weak assumptions, the estimate are asymptotically standard normal distributed with Gaussian confidence intervals. [3]. Put in context of economic policy, the presence of an asymptotic theory allows hypothesis testing regarding treatment results and thus aids policy decisions to be made.

Causal forests are a specific form of generalized random forests (GRF), which uses adaptive sample splitting criterion taking into account the MSE. To avoid overfitting and reduce bias in the estimates, we also ensure that the tree is 'honest' - the subsample with which we grow a tree is disparate from the subsample with which we drop down and obtain predictions. Additionally, we also note that although causal forests uncover heterogeneous treatment effects with valid confidence intervals for statistical inference, it does not necessarily address the affect of confounding due to the regularization on our trees. The terminal leaves of our tree are not homogenous across covariates for the sake of lowering variance and

thus increasing precision, but with confounding within the leaves, we cannot guarantee that our treatment effect estimates will be unbiased. We thus use the Double/Debiased Machine Learning method (DML) for causal forests proposed by Chernozhukov et al. (2016) that uses orthogonalized treatment on covariates to estimate the treatment effect. [4] We achieve this through the *CausalForestDML* function. In our model, we use 80% of the welfare data to build our tree and estimate treatment effects on the rest with a 5-fold CV for parameters of the model.

2.3 Results of analysis and explanation

We first investigate the average treatment effects (ATE) without conditioning on any covariates for possible effect heterogeneity. We obtain an ATE of 0.336, with a standard deviation of 0.049 and a 95% confidence interval of (0.256, 0.416). Green and Kern (2012) identified seven variables to condition the treatment effects on: *party identification*, *political views*, *age*, *education*, *negative attitude towards Blacks*, and *survey year*. Figure 1 displays the CATE estimates obtained by the causal forest model conditional on each of the seven variables. The blue areas represent the 95% confidence intervals of our estimates.

The top two graphs represents CATE as a function of party identification and self-identified political alignment from liberal to conservative. We see around 5 and 3 percentage points difference in the effect of question wording on support for welfare spending between strong Republicans-Democrats and Conservative-Liberals, respectively, controlling for all other covariates. Conservatives and strong Republicans are more likely to be affected by the framing of the survey question as for welfare. The treatment effect conditional on age, on the other hand, is greatest for those in their 30-40s, and diminishes past that age group. The negative bias toward Blacks has a less pronounced moderation on the treatment effect than the BART estimates. We note that the trends in CATE estimates obtained via causal forest are generally similar to those obtained using BART for most of the covariate groups except education. There seems to be an increasing effect of question wording as the respondent receives greater years of education, with the effect peaking at around 11-13 years (with college education). Whereas Green and Kern found no distinct moderation of treatment effects based on education, we find that there is around 5 percentage point difference in treatment effect estimates between those who received no education and those who received post-Graduate education. This indicates that more educated individuals respond more favorably to the question worded as 'assistance to the poor'. Lastly, we examine that change in CATE through time (survey year) and find that the treatment effect was strongest during years 1993-1996. Our results once again largely agree with the results obtained through BART, and illustrate that response to "welfare" is highly associated with the respondent's personal background characteristics.

We note that our CATE estimates from causal forest are on average lower than those of BART. The higher estimates of the BART estimates can be attributed to the regularization-induced confounding (RIC) identified by Hahn, Murray, and Carvalho (2019). [7] The treatment effect in BART is obtained by taking the conditional expectation of the outcome

conditional on treatment and certain covariates: $E[Y|x, Z = 1] - E[Y|x, Z = 0]$, which in the presence of confounding in finite sample size, outcome may be dependent mostly on x rather than the treatment Z . RIC thus states that the obtained CATEs are heavily dependent on the regularization by the prior distribution in samples of defined sizes.

Figure 2 shows the overall presence of treatment effect heterogeneity in the sample. The histogram shows that treatment effect ranges from 10 percentage points to 52 percentage points, with the median estimated CATE of 34 percentage points. We note that compared to the BART median of 37 percentage points, we again obtain lower estimates on average. All estimates of CATE were positive, indicating that although the degrees of the response varies based on personal background covariates, all respondents react more favorably to a question framed for 'the poor' than for 'welfare'.

Add Interpretation of results, Comparison to BART (Bias, RMSE, Coverage, Int length)?? TBD

2.4 Suggestions for further analysis

The causal forests approach also allows us to gauge covariate importance by calculating SHAP values, as indicated in figure 3. We note that aside from the covariates estimated by Green and Kern, we find that variables as work status and racial backgrounds are also significant in the model. Although in this paper we limit the scope of our research to comparisons with the authors' estimates through BART, it may be meaningful in future research to estimate CATE conditional on those variables.

3 Testing Parameter Robustness of Causal Forest Estimates on Synthetic DGP

The analysis using the welfare data in the section above gives us an interesting real-life application of causal forest method; however, given the lack of ground truth in the empirical situation, we are unable to test the performance of our model in terms of metrics as bias and MSE. In this section, we explore multiple synthetic DGPs to estimate how causal forest performs under different data parameters.

3.1 Baseline DGP specification

We specify a DGP that resembles our empirical data using a simplified model of N individuals and 120 covariates. Each variable is drawn from a gaussian normal with mean and covariance obtained from the empirical data. We then take each randomly generated covariate and convert to its corresponding form. We assume that the outcome data is generated as:

$$Y = B\Gamma + T \cdot \Theta(X) \tag{1}$$

where $\Gamma = X + \sum_{i=0}^I X_i^\omega$ is the full list of covariates including their higher-order interactions with the i most important features in the model. This parameter models the degree of linearity in the response surface: i features are selected based on SHAP values obtained in the empirical analysis from section 1, and raised to the order of ω to model non-linear interactions and increased importance of these variables. In our baseline model, we sample $N=5000$ individuals with parameters $\omega = 3, i = 4$ for ‘med’-degree linearity. T , our treatment vector, models both binary and continuous treatment assignment; in the binary case, $T \sim \text{Binomial}(n, p)$ where p is the probability of receiving treatment (propensity score). We assume there is a 50% probability of treatment in the baseline model. In the continuous case, we sample treatment from $T \sim \text{Unif}(0, 1)$. We model heterogeneous or homogeneous treatment effects through $\Theta(X)$ as defined in equation (2).

$$\Theta(X) = \begin{cases} \delta_H & \{H = 0\} \\ \delta_H \cdot (\sum_{j=0}^J X_j + (\sum_{j=0}^J X_j)^\omega) & \{H = 1\} \end{cases} \quad (2)$$

Variable H indicates whether heterogenous treatment effects are present in the model ($H = 1$), and δ_H is the size of the effect. We take $\delta_0 = 10, \delta_1 = 2$ in the base model for each homogeneous and heterogeneous cases. X_j parameters are again the j most influential features in the empirical data as obtained by SHAP values.

Also add about ATE stuff etc in baseline DGP In the following analysis on parameter robustness of the model, we omit CATE estimation due to computational constraints.

3.2 DGP parameter modification and expectations of CF performance

Based on the synthetic DGP above, we test the causal forest model for robustness by modifying the parameters. We then estimate the model $K = 50$ times for each parameter value in a range of different parameters. We consider the following dimensions for parameter tweaking: 1) *sample size*, 2) *non-linearity in covariates*, 3) *propensity score*, 4) *overlap*, 5) *degree of treatment effect heterogeneity*, and 6) *number of estimators*.

3.2.1 Sample size

In order to test the asymptotic theory of the causal forest estimates, we test the model performance on multiple sample sizes, $N = \{1000, 5000, 10000\}$. We expect from standard statistics for both our bias and variance to decrease with increasing sample size, given our estimate converges to the true CATE value. Accordingly, interval length would also decrease.

3.2.2 Linearity in response surface

We tweak the degree of linearity in our covariates to check if the causal forests model is able to handle higher degree relationships and complex interactions between variables. 4 degrees of linearity are explored, specified by i in our definition of Γ in equation (1) - which

corresponds to *full* ($i = 0$), *high* ($i=2$), *med* ($i=4$), *low* ($i=8$). As causal forests are expected to be apt at detecting non-linear relationships than traditional treatment effect estimation methods as OLS, we would expect the model to be robust to the introduction of a non-linear response surface.

3.2.3 Propensity score

In this scenario, the data has an imbalanced treatment and control group sizes. Since treatment assignment is random with the propensity score $p = \pi(X)$ being constant across all individuals we do not expect this change to introduce any selection bias. However, as the ATE estimates are obtained by taking differences in means at the terminal node, the deviation in group sizes from a 50/50 split will likely harm the power of our model and increase the variance in our estimates. We estimate treatment probabilities of $p = \{0.1, 0.5, 0.9\}$ for extreme cases of imbalance.

3.2.4 Overlap

We explore situations in which the propensity score p satisfies or fails the overlap condition: $0 < p = \pi(X) < 1$. Given that in traditional statistical settings, overlap ensures there are observations on which we can estimate credible counterfactuals, we expect our causal forest estimate variances to increase without the condition. When overlap does not hold, we randomly sample half of the observations and force $T = 0$ to ensure that half the sample always has a 0% probability of receiving treatment.

3.2.5 Degree of treatment effect heterogeneity

Causal forest optimizes sample-split by preferring leaves with heterogeneity in a key parameter and penalizing those with greater variance [2]. The model is expected to show stable performance across the presence of complex heterogeneous effects. For this reason, we decide to alter the degree of treatment effect heterogeneity in the model, $\Theta(X)$, to check if the causal forest estimates are robust under such changes. We specify 4 different parameter values of $j = \{0, 2, 4, 8\}$ for j in equation (2).

3.2.6 Number of estimators

We vary the number of trees we fit in our causal forest estimator to see if our estimates are sensitive to tuning parameters. The baseline model had 1000 trees in each forest, and we expect with increasing number of trees we will be able to reduce overfitting and the importance of each tree in the estimate. Naturally, we believe the increase in this parameter will also show an increase in bias for CATE as our estimates are averaged out over multiple trees. However, this may not be as pronounced in the ATE estimate, which takes the average across the entire sample. We take the *number of trees* = $\{100, 500, 1000, 5000\}$ in each forest.

3.3 Results and discussion

Figure (4) and Table (1) illustrates the parameter resilience of causal forest ATE estimates for our specified DGP.

4 Conclusion

WRITESOME SHIT

Table 1: Model Performance with Parameter Variation

(a) Sample Size

N	True ATE	Estimated ATE	Absolute Bias	Relative Bias	RMSE	Relative RMSE	Relative Int. Length	Int. Length	Coverage
1000	99999	99999	243.20	99999	287.11	99999	1528.25	99999	1.0
5000	99999	99999	101.21	99999	126.20	99999	881.80	99999	1.0
10000	99999	99999	73.12	99999	95.11	99999	698.63	99999	1.0

Relative values calculated as percentage of estimated value compared to true population ATE

(b) Linearity

i	True ATE	Estimated ATE	Absolute Bias	Relative Bias	RMSE	Relative RMSE	Relative Int. Length	Int. Length	Coverage
full	99999	99999	243.20	99999	287.11	99999	1528.25	99999	1.0
high	99999	99999	101.21	99999	126.20	99999	881.80	99999	1.0
med	99999	99999	73.12	99999	95.11	99999	698.63	99999	1.0
low	99999	99999	73.12	99999	95.11	99999	698.63	99999	1.0

Relative values calculated as percentage of estimated value compared to true population ATE

(c) Propensity Score

p	True ATE	Estimated ATE	Absolute Bias	Relative Bias	RMSE	Relative RMSE	Relative Int. Length	Int. Length	Coverage
0.1	99999	99999	243.20	99999	287.11	99999	1528.25	99999	1.0
0.5	99999	99999	101.21	99999	126.20	99999	881.80	99999	1.0
0.9	99999	99999	73.12	99999	95.11	99999	698.63	99999	1.0

Relative values calculated as percentage of estimated value compared to true population ATE

(d) Overlap

overlap	True ATE	Estimated ATE	Absolute Bias	Relative Bias	RMSE	Relative RMSE	Relative Int. Length	Int. Length	Coverage
True	99999	99999	243.20	99999	287.11	99999	1528.25	99999	1.0
False	99999	99999	101.21	99999	126.20	99999	881.80	99999	1.0

Relative values calculated as percentage of estimated value compared to true population ATE

(e) Degree of Heterogeneity

j	True ATE	Estimated ATE	Absolute Bias	Relative Bias	RMSE	Relative RMSE	Relative Int. Length	Int. Length	Coverage
0	99999	99999	243.20	99999	287.11	99999	1528.25	99999	1.0
2	99999	99999	101.21	99999	126.20	99999	881.80	99999	1.0
4	99999	99999	73.12	99999	95.11	99999	698.63	99999	1.0
8	99999	99999	73.12	99999	95.11	99999	698.63	99999	1.0

Relative values calculated as percentage of estimated value compared to true population ATE

(f) Number of Estimators

Num_Estimator	True ATE	Estimated ATE	Absolute Bias	Relative Bias	RMSE	Relative RMSE	Relative Int. Length	Int. Length	Coverage
100	99999	99999	243.20	99999	287.11	99999	1528.25	99999	1.0
500	99999	99999	101.21	99999	126.20	99999	881.80	99999	1.0
1000	99999	99999	73.12	99999	95.11	99999	698.63	99999	1.0
5000	99999	99999	73.12	99999	95.11	99999	698.63	99999	1.0

Relative values calculated as percentage of estimated value compared to true population ATE

References

- [1] S. Athey and G. Imbens. Machine learning methods economists should know about, 2019.
- [2] S. Athey, J. Tibshirani, and S. Wager. Generalized random forests, 2018.
- [3] S. Athey and S. Wager. Estimating treatment effects with causal forests: An application, 2019.
- [4] V. Chernozhukov, D. Chetverikov, M. Demirer, E. Duflo, C. Hansen, W. Newey, and J. Robins. Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal*, 21(1):C1–C68, 2014.
- [5] H. A. Chipman, E. I. George, and R. E. McCulloch. Bart: Bayesian additive regression trees. 4(1), Mar 2010.
- [6] D. P. Green and H. L. Kern. Modeling heterogeneous treatment effects in survey experiments with bayesian additive regression trees. *The Public Opinion Quarterly*, 76(3):491–511, 2012.
- [7] P. R. Hahn, J. S. Murray, and C. M. Carvalho. Bayesian Regression Tree Models for Causal Inference: Regularization, Confounding, and Heterogeneous Effects (with Discussion). *Bayesian Analysis*, 15(3):965 – 1056, 2020.
- [8] K. A. Rasinski. The effect of question wording on public support for government spending. 53(3):388–394. eprint: <https://academic.oup.com/poq/article-pdf/53/3/388/5301247/53-3-388.pdf>.
- [9] V. Rockova and E. Saha. On theory for bart, 2018.

Figure 1: CATE Estimates by Covariate

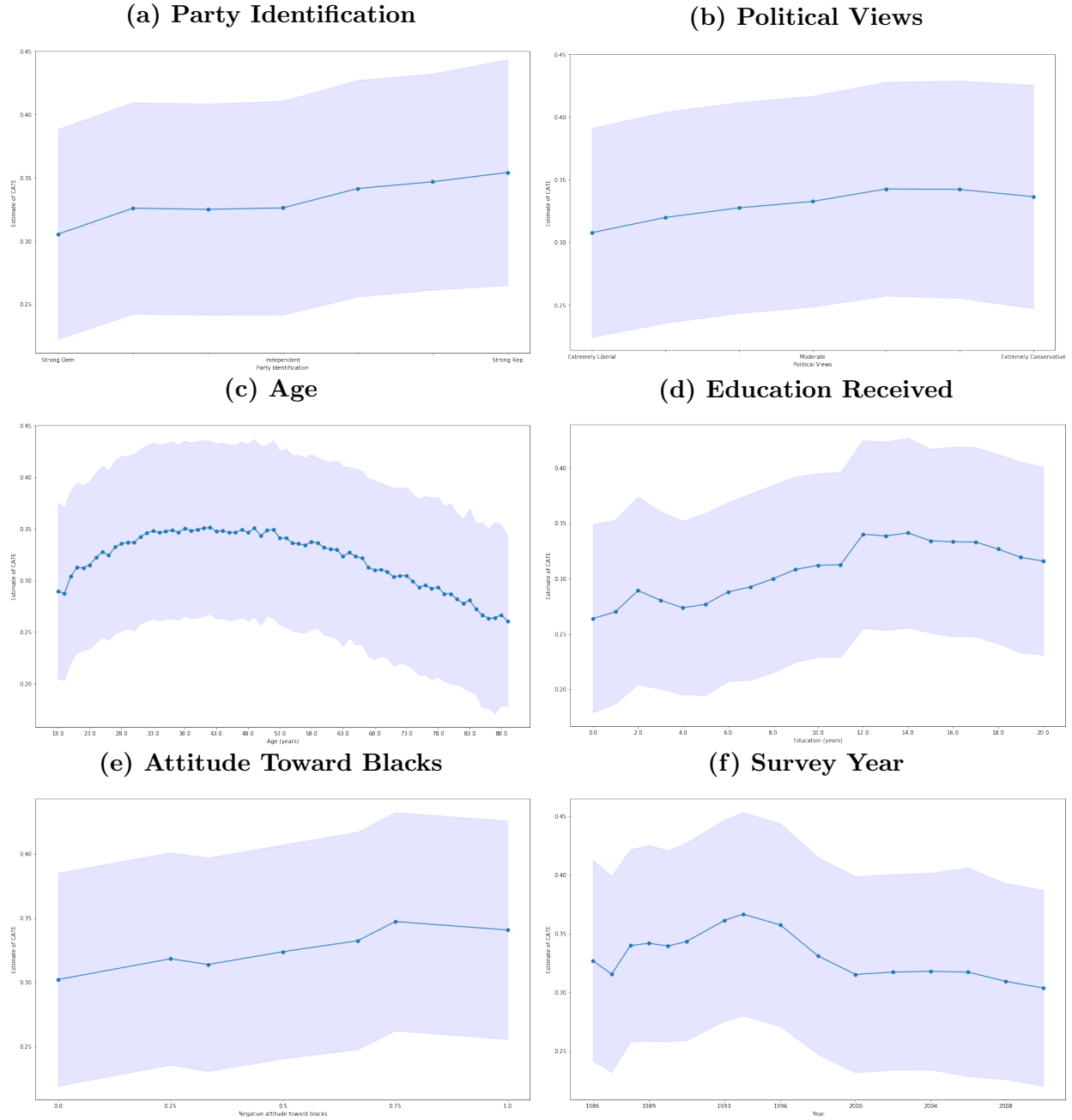


Figure 2: Histogram of CATE Estimates

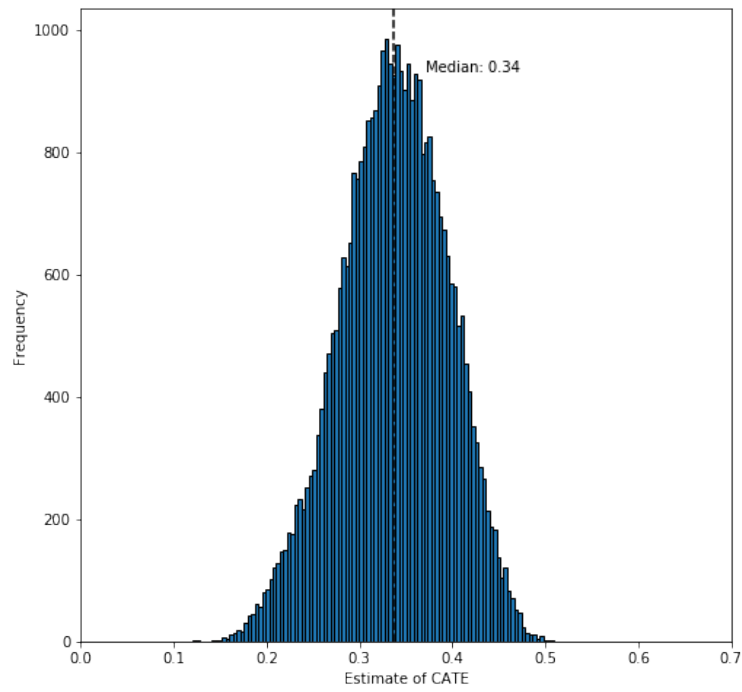


Figure 3: Summary of SHAP Values of Covariate Importance

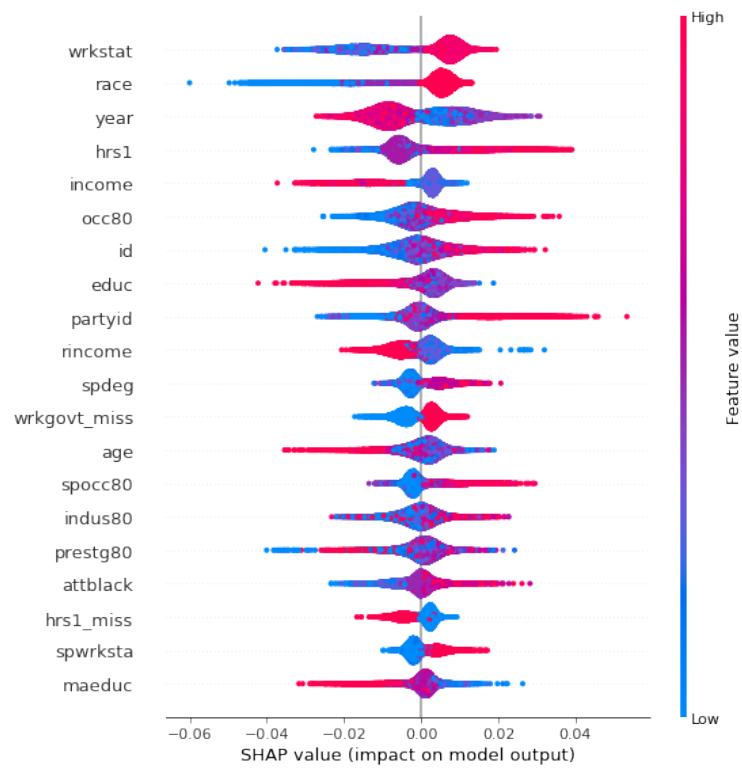
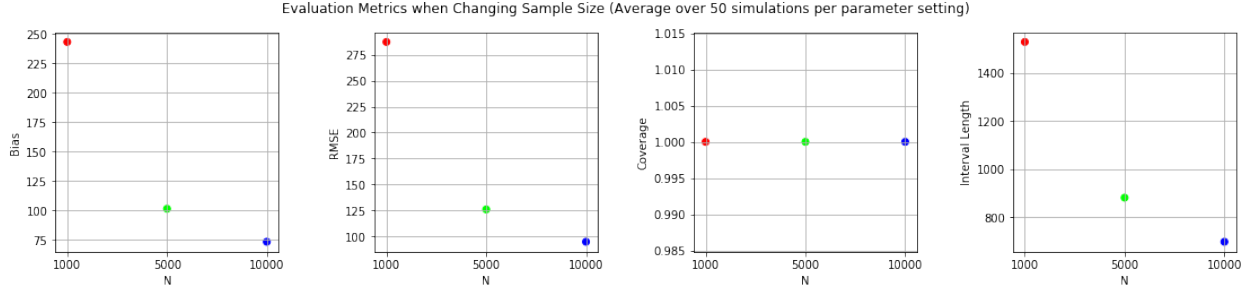
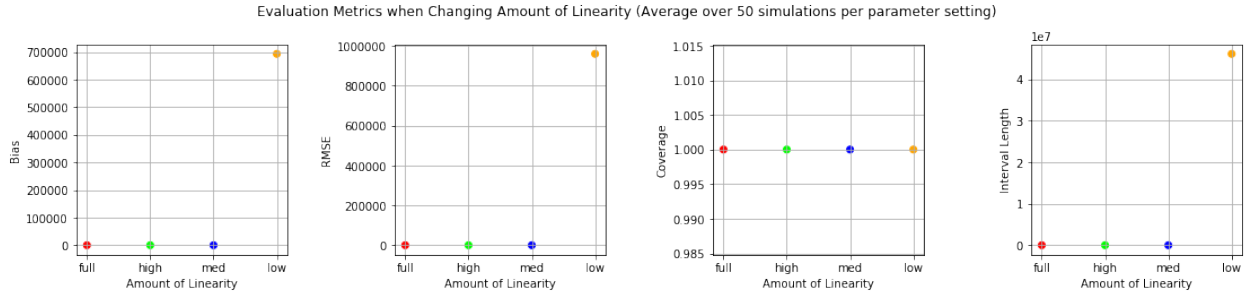


Figure 4: Model Performance by Parameter Variation

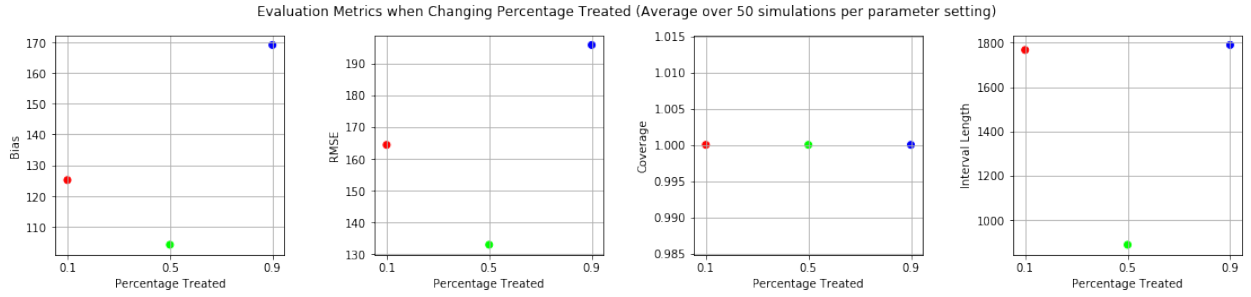
(a) Sample Size



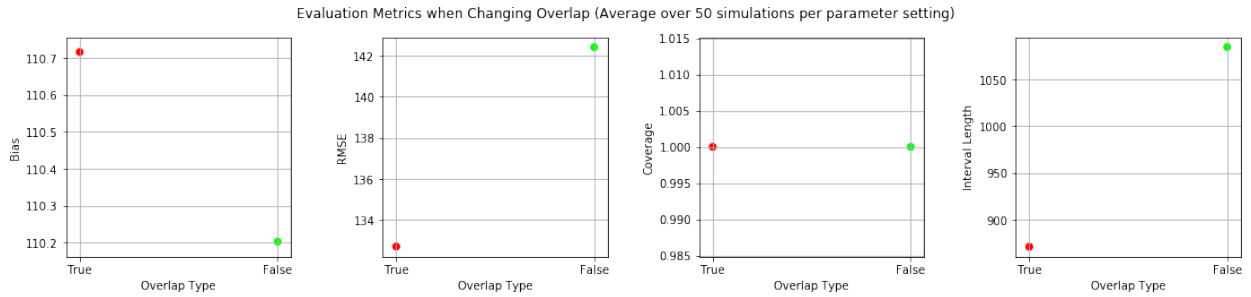
(b) Linearity



(c) Propensity Score

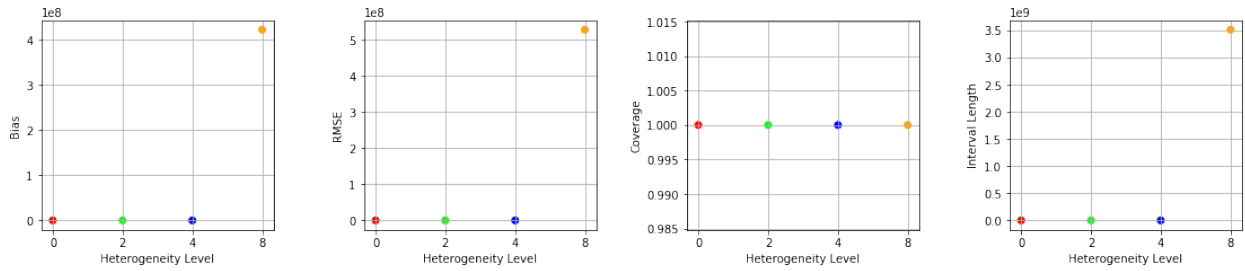


(d) Overlap



(e) Degree of Heterogeneity

Evaluation Metrics when Changing Heterogeneity Level (Average over 50 simulations per parameter setting)



(f) Number of Estimators

Evaluation Metrics when Changing Heterogeneity Level (Average over 50 simulations per parameter setting)

