

# Estimating population average treatment effects from experiments with noncompliance\*

Kellie Ottoboni<sup>†</sup>      Jason Poulos<sup>‡</sup>

August 8, 2019

## Abstract

This paper extends a method of estimating population average treatment effects to settings with noncompliance. Simulations show the proposed compliance-adjusted estimator performs better than its unadjusted counterpart when compliance is relatively low and can be predicted by observed covariates. We apply the proposed estimator to measure the effect of Medicaid coverage on health care use for a target population of adults who may benefit from expansions to the Medicaid program. We draw randomized control trial data from a large-scale health insurance experiment in which a small subset of those randomly selected to receive Medicaid benefits actually enrolled.

---

\*The authors thank Jon McAuliffe and Jas Sekhon and his research group at UC Berkeley for valuable comments. Poulos acknowledges support from the National Science Foundation Graduate Research Fellowship [grant number DGE 1106400]. This work used the Extreme Science and Engineering Discovery Environment (XSEDE) supercomputer Stampede2 at the Texas Advanced Computing Center (TACC) through allocation SES180010.

<sup>†</sup>Department of Statistics, University of California, Berkeley. email: [kellieotto@berkeley.edu](mailto:kellieotto@berkeley.edu).

<sup>‡</sup>*Corresponding author:* [Department of Statistical Science, Duke University, and Statistical and Applied Mathematical Sciences Institute](#). email: [jpoulos@samsi.info](mailto:jpoulos@samsi.info).

# 1 Introduction

Randomized control trials (RCTs) are the gold standard for estimating the causal effect of a treatment. An RCT may give an unbiased **estimates of sample average treatment effects**, but external validity is an issue when the individuals in the RCT are unrepresentative of the actual population of interest. For example, the participants in an RCT in which individuals volunteer to sign up for health insurance may be in poorer health at baseline than the overall population. External validity is particularly relevant to policymakers who want to know how the treatment effect would generalize to the broader population.

**A new research frontier for causal inference focuses on developing methods for extrapolating RCT results to a population (e.g., Imai et al., 2008). Existing approaches to this problem are based in settings where there is full compliance with treatment.** However, noncompliance is a prevalent issue in RCTs and occurs when individuals who are assigned to the treatment group do not comply with the treatment. For individuals assigned to control, we are unable to observe who would have complied had they been assigned treatment. Noncompliance biases the intention-to-treat (ITT) estimate of the effect of treatment assignment toward zero.

We propose a method to estimate the complier-average causal effects for the target population from RCT data with noncompliance, and refer to this **estimator as the Population Average Treatment Effect on Treated Compliers (PATT-C)**. PATT-C involves the expectation of the response of RCT compliers, conditional on their covariates, where the expectation is taken over the distribution of covariates for population members receiving treatment. Our estimation strategy differs from reweighting methods that use propensity scores to adjust the RCT data. **In Stuart et al. (2011), for example,** a propensity score model is used to predict participation in the RCT, given pretreatment covariates common to both the RCT and population data. Individuals in the RCT and population are then weighted according to the inverse of the estimated propensity score. **Similarly, Hartman et al. (2015) propose a method of reweighting the responses of individuals in an RCT according to the covariate**

distribution of the population.

Our approach for estimating PATT-C differs from previous approaches because we only need the potential outcomes for RCT compliers and we cannot observe who in the control group would have complied had they been assigned treatment. We propose an alternative approach of modeling compliance in the RCT and using the compliance model to predict the likely compliers in the RCT control group.<sup>1</sup> Assuming that the response surface is the same for compliers in the RCT and population members who received treatment, we then predict the response surface for all RCT compliers and use the predicted values from the response surface model to estimate the potential outcomes of population members who received treatment, given their covariates.

When estimating the average causal effect **for compliers** from an RCT, researchers typically scale the estimated ITT effect by the compliance rate, **assuming that there is only single crossover from treatment to control**.<sup>2</sup> When extrapolating RCT results to a population, one might simply **reweight the ITT effect according to the covariate distribution of the population and then divide by the proportion of treated compliers in the population** in order to yield a population average effect of treatment on treated compliers. However, **we do not observe the population compliance rate. Moreover, the population compliance rate is likely to differ across subgroups based on pretreatment covariates.** By explicitly modeling compliance, this approach allows researchers to decompose population estimates by covariate group, which is useful for policymakers in evaluating the efficacy of policy interventions for subgroups of interest in a population.

We apply the proposed estimator to measure the effect of Medicaid coverage on health care use for a target population of adults who may benefit from government-backed expan-

---

<sup>1</sup>Reweighting methods typically leverage exchangeability of potential outcomes between the covariate-adjusted treated and controls in the RCT. In our approach, the potential outcomes between the complier treated and complier controls are not exchangeable by design, since we need to assume we know the compliance model.

<sup>2</sup>Alternative approaches include estimating sharp bounds to the ITT effect in the presence of noncompliance (Balke and Pearl, 1997; Imai et al., 2013), adjustment for treatment noncompliance using principal stratification (Frangakis and Rubin, 2002; Frumento et al., 2012), and maximum-likelihood and Bayesian inferential methods (Yau and Little, 2001).

sions to the Medicaid program. We are particularly interested in measuring the effect of Medicaid on emergency room (ER) use because it is the main delivery system through which the uninsured receive health care. The uninsured could potentially receive higher quality health care through primary care visits. An important policy question is whether Medicaid expansions will decrease ER utilization and increase primary care visits by the previously uninsured. We draw RCT data from a large-scale health insurance experiment, in which only 30% of those randomly selected to receive Medicaid benefits actually enrolled. We find substantial differences between sample and population estimates in terms of race, education, and health status subgroups.

The paper proceeds as follows: Section 2 presents the proposed estimator and the necessary assumptions for its identifiability; Section 3 describes the estimation procedure; Section 4 reports the estimator’s performance in simulations; Section 5 uses the estimator to identify the effect of extending Medicaid coverage to the low-income adult population in the U.S; Section 6 discusses the results and offers direction for future research.

## 2 Estimator

We are interested in using the outcomes from an RCT to estimate complier-average causal effects for a target population. Compliance with treatment in the population is not assigned at random, but rather may depend on unobserved variables, confounding the effect of treatment received on the outcome of interest. RCTs are needed to isolate the effect of treatment received.

Ideally, we would take the results of an RCT and reweight the sample such that the reweighted covariates match the those in the population. In practice, one rarely knows the true covariate distribution in the target population. Instead, we consider data from a nonrandomized, observational study in which participants are representative of the target population. The proposed estimator combines RCT and observational data to overcome

these issues.

## 2.1 Assumptions

Let  $Y_{isd}$  be the potential outcome for individual  $i$  in group  $s$  and treatment **received**  $d$ . Let  $S_i$  denote the sample assignment, where  $s = 0$  is the population and  $s = 1$  is the RCT.  $T_i$  indicates treatment assignment and  $D_i$  indicates whether treatment was actually received. Treatment is assigned at random in the RCT, so we observe both  $D_i$  and  $T_i$  when  $S_i = 1$ . For compliers in the RCT,  $D_i = T_i$ .

Let  $W_i$  be individual  $i$ 's observable pretreatment covariates that are related to the sample selection mechanism for membership in the RCT, treatment assignment in the population, and complier status. Let  $C_i$  be an indicator for individual  $i$ 's compliance **with** treatment, which is only observable for individuals in the RCT treatment group.

In the population, we suppose that treatment is made available to individuals based on their covariates  $W_i$ . Individuals with  $T_i = 0$  do not receive treatment, while those with  $T_i = 1$  may decide whether or not to accept treatment. For individuals in the population, we only observe  $D_i$  — not  $T_i$ .<sup>3</sup>

**Assumption 1.** *Consistency under parallel studies:*

$$Y_{i0d} = Y_{i1d} \quad \forall i, d = \{0, 1\}.$$

Assumption (1) requires that each individual  $i$  has the same response to treatment, whether  $i$  is in the RCT or not. Compliance status  $C_i$  is not a factor in this assumption because we assume that compliance is conditionally independent of sample and treatment assignment for all individuals with covariates  $W_i$ .

---

<sup>3</sup>We frame Assumptions (3) and (4) in terms of  $C_i$  and  $T_i$  in order to distinguish among the population controls who should have received treatment (i.e., individuals with  $T_i = 1$  and  $D_i = 0$ ) from noncompliers assigned to control (i.e., individuals with  $T_i = 0$  and  $D_i = 0$ ).

**Assumption 2.** *Conditional independence of compliance and assignment:*

$$C_i \perp\!\!\!\perp S_i, T_i \mid W_i, \quad 0 < \mathbb{P}(C_i = 1 \mid W_i) < 1.$$

Assumption (2) implies that  $P(C_i = 1 \mid S_i = 1, T_i = 1, W_i) = P(C_i = 1 \mid S_i = 1, T_i = 0, W_i)$ , which is useful when predicting the probability of compliance as a function of covariates  $W_i$  in the first step of the estimation procedure. Together, Assumptions (1) and (2) ensure that potential outcomes do not differ based on sample assignment or receipt of treatment.

**Assumption 3.** *Strong ignorability of sample assignment for treated:*

$$(Y_{i01}, Y_{i11}) \perp\!\!\!\perp S_i \mid (W_i, T_i = 1, C_i = 1), \quad 0 < \mathbb{P}(S_i = 1 \mid W_i, T_i = 1, C_i = 1) < 1.$$

Assumption (3) ensures the potential outcomes for treatment are independent of sample assignment for individuals with the same covariates  $W_i$  and assignment to treatment.<sup>4</sup> We make a similar assumption for the potential outcomes under control:

**Assumption 4.** *Strong ignorability of sample assignment for controls:*

$$(Y_{i00}, Y_{i10}) \perp\!\!\!\perp S_i \mid (W_i, T_i = 1, C_i = 1), \quad 0 < \mathbb{P}(S_i = 1 \mid W_i, T_i = 1, C_i = 1) < 1.$$

Restrictive exclusion criteria in RCTs can result in a sample covariate distribution that differs substantially from the population covariate distribution, thereby reducing the external validity of RCTs (Rothwell, 2005). High rates of exclusion also poses a threat to strong ignorability assumptions if exclusion increases the likelihood that there are unobserved differences between the RCT and target population that are correlated with potential outcomes. For example, the RCT described in Section 5 required enrolled participants to recertify their eligibility status every six months during the study period. The exclusion of participants who

---

<sup>4</sup>Throughout, we assume individuals are sampled randomly from an infinite population.

failed to recertify because their household income exceeded a given cutoff threatens strong ignorability if the factors that contributed to the failure to recertify are correlated with unobservables that are also correlated with potential outcomes. The placebo tests described in Section 5.3.1 are designed to detect bias arising from violations of the strong ignorability assumptions.<sup>5</sup>

Figure 1 shows Assumptions (2), (3), and (4) in a directed acyclic graph. Treatment assignment  $T_i$  may only depend on  $C_i$  through  $W_i$ , and the potential outcomes  $(Y_{is0}, Y_{is1})$  may only depend on  $S_i$  through  $W_i$ . From the internal validity standpoint, the role of  $W_i$  is critical: if any relevant observed covariates are not controlled, then there is a backdoor pathway from  $T_i$  back to  $W_i$  and into  $Y_{isd}$ .<sup>6</sup>

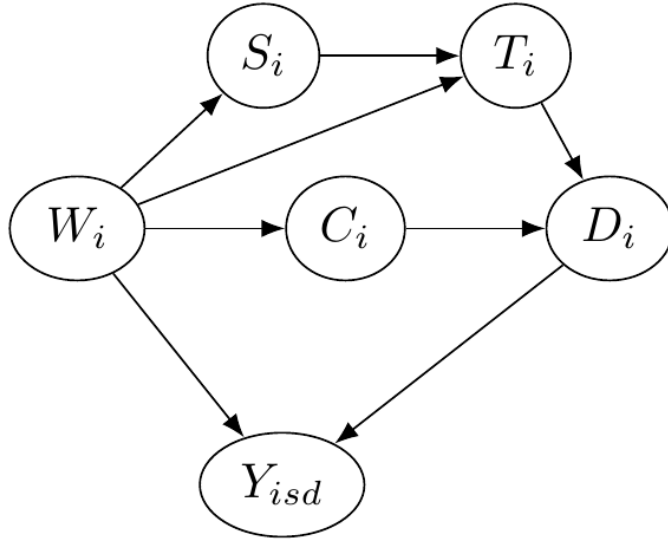


Figure 1: Causal diagram indicating the conditional independence assumptions needed to estimate the PATT-C.

Interference undermines the framework because it creates more than two potential out-

<sup>5</sup>Note that **Assumptions (3) and (4)** also imply strong ignorability of sample assignment for treated and control noncompliers since we assume in that compliance is also independent of sample and treatment assignment, conditional on  $W_i$  (Assumption (2)). However, we are interested only on modeling the response surfaces for compliers.

<sup>6</sup>We use the same  $W_i$  across all identifying assumptions, which implicitly assumes that the observable covariates that determine sample selection also determine population treatment assignment and complier status. This choice reflects a modeling assumption of the estimation procedure described in Section 3.

comes per participant, depending on the **received treatment** of other participants (Rubin, 1990). We therefore assume no interference between units:

**Assumption 5.** *The potential outcomes  $Y_{isd}$  do not depend on  $D_j$ ,  $\forall j \neq i$ .*

We also include the following assumptions made by Angrist et al. (1996) to identify the **average causal effect for compliers**. Assumption (6) ensures that crossover is only possible from treatment to control:

**Assumption 6.** *No defiers:*

$$T_i \geq D_i, \quad \forall i, d, t = \{0, 1\}.$$

Assumption (7) ensures treatment assignment affects the response only through the treatment received. In particular, the treatment effect may only be nonzero for compliers.

**Assumption 7.** *Exclusion restriction: For noncompliers,*

$$Y_{i11} = Y_{i10}, \quad \forall i.$$

## 2.2 PATT-C

PATT-C is interpreted as the complier-average causal effect estimated on the RCT sample extrapolated to what we would have observed in the population if treatment received  $D_i$  is the same. **It is written as follows:**

$$\tau_{\text{PATT-C}} = \mathbb{E}(Y_{i01} - Y_{i00} \mid S_i = 0, D_i = 1). \quad (1)$$

The following theorem relates the treatment effect in the RCT to the treatment effect in the population.



**Theorem 1.** *Under Assumptions (1) – (7),*

$$\tau_{PATT-C} = \mathbb{E}_{01} [\mathbb{E}(Y_{i11} \mid S_i = 1, D_i = 1, W_i)] - \mathbb{E}_{01} [\mathbb{E}(Y_{i10} \mid S_i = 1, T_i = 0, C_i = 1, W_i)] \quad (2)$$

where  $\mathbb{E}_{01} [\mathbb{E}(\cdot \mid \dots, W_i)]$  denotes the expectation with respect to the distribution of  $W_i$  for population members who received treatment.

*Proof.* We separate the expectation linearly into two terms and consider each individually.

$$\begin{aligned} \mathbb{E}(Y_{i01} \mid S_i = 0, D_i = 1) &= \mathbb{E}(Y_{i11} \mid S_i = 0, D_i = 1) && \text{by Assumption (1)} \\ &= \mathbb{E}(Y_{i11} \mid S_i = 0, T_i = 1, C_i = 1) && \text{by Assumption (6)} \\ &= \mathbb{E}_{01} [\mathbb{E}(Y_{i11} \mid S_i = 0, T_i = 1, C_i = 1, W_i)] \\ &= \mathbb{E}_{01} [\mathbb{E}(Y_{i11} \mid S_i = 1, T_i = 1, C_i = 1, W_i)] && \text{by Assumption (3)} \\ &= \mathbb{E}_{01} [\mathbb{E}(Y_{i11} \mid S_i = 1, D_i = 1, W_i)] \end{aligned}$$

Intuitively, conditioning on  $W_i$  makes sample selection ignorable under Assumption (3). This is the critical connector between the third and fourth lines of the first expectation derivation.

$$\begin{aligned} \mathbb{E}(Y_{i00} \mid S_i = 0, D_i = 1) &= \mathbb{E}(Y_{i10} \mid S_i = 0, D_i = 1) && \text{by Assumption (1)} \\ &= \mathbb{E}(Y_{i10} \mid S_i = 0, T_i = 1, C_i = 1) && \text{by Assumption (6)} \\ &= \mathbb{E}_{01} [\mathbb{E}(Y_{i10} \mid S_i = 1, T_i = 1, C_i = 1, W_i)] && \text{by Assumption (4)} \\ &= \mathbb{E}_{01} [\mathbb{E}(Y_{i10} \mid S_i = 1, T_i = 0, C_i = 1, W_i)] && \text{by Assumption (2)} \end{aligned}$$

The last line follows because Assumption (2) allows us to use RCT controls who would have complied had they been assigned to treatment. Finally, the result follows by plugging

these two expressions into Eq. (1). □

### 3 Estimation procedure

There are two challenges in turning Theorem (1) into an estimator of  $\tau_{\text{PATT-C}}$  in practice. First, we must estimate the inner expectation over potential outcomes of compliers in the RCT. In the empirical example, we use an ensemble of algorithms (van der Laan et al., 2007) to estimate the response surface for compliers in the RCT, given their covariates. Thus, the first term in the expression for  $\tau_{\text{PATT-C}}$  is estimated by the weighted average of points on the response surface, evaluated for each treated population member’s potential outcome under treatment. The second term is estimated by the weighted average of points on the response surface, evaluated for each treated population member’s potential outcome under control.

The second challenge is that we cannot observe which individuals are included in the estimation of the second term. In the RCT control group,  $C_i$  is unobservable, as they always receive no treatment ( $D_i = 0$ ). We must estimate the second term of Eq. (2) by predicting who in the control group would be a complier had they been assigned to treatment. **Explicitly modeling compliance allows us to decompose PATT-C estimates by subgroup according to covariates common to both RCT and observational datasets. This approach also accounts for settings where the compliance rate differs between the sample and population, as well as across subgroups.**

The procedure for estimating  $\tau_{\text{PATT-C}}$  using Theorem (1) is as follows:

- S.1** Using the group assigned to treatment in the RCT ( $S_i = 1, T_i = 1$ ), train a model (or an ensemble of models) to predict the probability of compliance as a function of covariates  $W_i$ .
- S.2** Using the model from S.1, predict who in the RCT assigned to control *would have* complied to treatment had they been assigned to the treatment group.<sup>7</sup>

---

<sup>7</sup>We use a standard prediction threshold of 50% in order classify compliers,  $C_i = 1$ . Adjusting the

**S.3** For the observed compliers assigned to treatment and predicted compliers assigned to control, train a model to predict the response using  $W_i$  and  $D_i$ , which gives  $\mathbb{E}(Y_{i1d} \mid S_i = 1, D_i = d, W_i)$  for  $d \in \{0, 1\}$ .

**S.4** For all individuals who received treatment in the population ( $S_i = 0, D_i = 1$ ), estimate their potential outcomes using the model from S.3, which gives  $Y_{i1d}$  for  $d \in \{0, 1\}$ . The mean counterfactual  $Y_{i11}$  minus the mean counterfactual  $Y_{i10}$  is the estimate of  $\tau_{\text{PATT-C}}$ .

Assumptions (3) and (4) are particularly important for estimating  $\tau_{\text{PATT-C}}$ : the success of the proposed estimator hinges on the assumption that the response surface is the same for compliers in the RCT and target population. If this does not hold, then the potential outcomes  $Y_{i10}$  and  $Y_{i11}$  for target population individuals cannot be estimated using the model from S.3. Section 5.3 discusses whether the strong ignorability assumptions are plausible in the empirical application.

### 3.1 Modeling assumptions

In addition to the identification assumptions, we require additional modeling assumptions for the estimation procedure. As pointed out in Section 2.1, we require that  $W_i$  is complete because if any relevant elements of  $W_i$  are not controlled, then there is a backdoor pathway from  $T_i$  back to  $W_i$  and into  $Y_{isd}$ . **Additionally**, we assume that the compliance model is accurate in predicting compliance in the training sample of RCT participants assigned to treatment and also generalizable to RCT participants assigned to control (S.1 and S.2). Section 3.2 below describes the method of evaluating the generalizability of the compliance model.

---

prediction threshold upward would result in more accurate classifications, although we do not explore this approach.

## 3.2 Ensemble method

In the empirical application, we use the weighted ensemble method described in van der Laan et al. (2007) for S.1 and S.3 of the estimation procedure. This ensemble method combines algorithms with a convex combination of weights based on minimizing cross-validated error. It is shown to control for overfitting and outperforms single algorithms selected by cross-validation (Polley and Van Der Laan, 2010).

We choose a variety of candidate algorithms to construct the ensemble, with a preference towards algorithms that tend to outperform in supervised classification tasks. We also have a preference for algorithms that have a built-in variable selection property. The idea is that we input the same  $W_i$  and each candidate algorithm selects the most important covariates for predicting compliance status or potential outcomes.<sup>8</sup> We select three types of candidate algorithms: nonparametric additive regression models (Buja et al., 1989); L1 or L2-regularized linear models (i.e., Lasso or ridge regression, respectively) (Tibshirani et al., 2012); and ensembles of decision trees (i.e., random forests) (Breiman, 2001). L1-regularized linear models are important for the application due to their variable selection properties: Lasso is particularly attractive because it tends to shrink all but one of the coefficients of correlated covariates to zero.

## 4 Simulations

We conduct a simulation study comparing the performance of the **PATT-C estimator against its unadjusted analogue, which we refer to as the Population Average Treatment Effect on the Treated (PATT)**:

$$\tau_{\text{PATT}} = \mathbb{E}(Y_{i01} - Y_{i00} \mid S_i = 0, D_i = 1). \quad (3)$$

---

<sup>8</sup>A potential concern when predicting potential outcomes is that the algorithm might shrink the treatment received predictor to zero, which would result in no difference between counterfactual potential outcomes.

Eq. (3) identifies the population-average causal effect of taking up treatment, adjusted according to the covariate distribution of population members who received treatment. We estimate the response curve for all RCT participants, conditional on their covariates and actual treatment received. Identical to S.4 in the estimation procedure for PATT-C, we use the response model to estimate the outcomes of population members who received treatment given their covariates, which are then used to estimate Eq. (3). Like the PATT-C estimator, the PATT estimator crucially relies on the assumption that the response surface is the same for RCT participants and population members who received treatment.

We compare the population estimators against the sample Complier Average Causal Effect (CACE) (Imbens and Rubin, 1997), which is commonly referred to the Local Average Treatment Effect in the econometrics literature (Imbens and Angrist, 1994; Angrist et al., 1996). In the context of program evaluation, it is a more relevant treatment effect of interest than the ITT effect because only RCT participants who received treatment would have their outcomes affected by treatment in the presence of a nonnegative treatment effect.

CACE is defined as the average causal effect of treatment received restricted to sample compliers:

$$\tau_{\text{CACE}} = \mathbb{E}(Y_{i11} - Y_{i10} \mid S_i = 1, C_i = 1). \quad (4)$$

In other words, CACE is the treatment effect for RCT participants who would comply regardless of treatment assignment. However, we are unable to observe the compliance status of RCT participants assigned to control because we do not know if they would have complied if they had been assigned to treatment. A generalization of the instrumental variables estimator of the CACE in the presence of noncompliers is given by:

$$\hat{\tau}_{\text{CACE}} = \frac{\mathbb{E}(Y_{i11} - Y_{i10} \mid S_i = 1)}{\mathbb{P}(T_i = D_i = 1 \mid S_i = 1)}, \quad (5)$$

which is equivalent to scaling the ITT effect by the sample proportion of treated compliers (e.g., Freedman, 2006). Eq. (5) is identified under Assumptions (5), (6), and (7).

## 4.1 Simulation design

The simulation is designed so that the effect of treatment is heterogeneous and depends on covariates which are different in the RCT and target population. The design satisfies the conditional independence assumptions in Figure 1.

In the simulation, RCT eligibility, complier status, and treatment assignment in the population depend on multivariate normal covariates  $(W_i^1, W_i^2, W_i^3, W_i^4)$  with means  $(0.5, 1, -1, -1)$  and covariances  $\text{Cov}(W_i^1, W_i^2) = \text{Cov}(W_i^1, W_i^4) = \text{Cov}(W_i^2, W_i^4) = \text{Cov}(W_i^3, W_i^4) = 1$  and  $\text{Cov}(W_i^1, W_i^3) = \text{Cov}(W_i^2, W_i^3) = 0.5$ . The first three covariates are observed by the researcher and  $W_i^4$  is unobserved.  $U_i, V_i, R_i$ , and  $Q_i$  are standard normal error terms.  $U_i, V_i, R_i, Q_i$ , and  $(W_i^1, W_i^2, W_i^3, W_i^4)$  are mutually independent.

The equation for selection into the RCT is

$$S_i = \mathbb{I}(e_2 + g_1 W_i^1 + g_2 W_i^2 + g_3 W_i^3 + e_4 W_i^4 + R_i > 0).$$

The parameter  $e_2$  varies the fraction of the population eligible for the RCT and  $e_4$  varies the degree of confounding with sample selection. We set the constants  $g_1, g_2$ , and  $g_3$  to be 0.5, 0.25, and 0.75, respectively.

Complier status is determined by

$$C_i = \mathbb{I}(e_3 + h_2 W_i^2 + h_3 W_i^3 + e_5 W_i^4 + Q_i > 0),$$

where  $e_3$  varies the fraction of compliers in the population, and  $e_5$  varies the degree of confounding with treatment assignment. We set the constants  $h_2$  and  $h_3$  to 0.5.

For individuals in the population ( $S_i = 0$ ), treatment is assigned by

$$T_i = \mathbb{I}(e_1 + f_1 W_i^1 + f_2 W_i^2 + e_6 W_i^4 + V_i > 0),$$

where  $e_1$  varies the fraction eligible for treatment in the population and  $e_6$  varies the degree

of confounding with sample selection. We set the constants  $f_1$  and  $f_2$  to 0.25 and 0.75, respectively. For individuals in the RCT ( $S_1 = 1$ ), treatment assignment  $T_i$  is a sample from a Bernoulli distribution with probability  $p = 0.5$ .

Finally, the response is determined by

$$Y_{isd} = a + bD_i + c_1W_i^1 + c_2W_i^2 + c_3W^4 + dU_i,$$

where we set  $a, c_1, c_3$ , and  $d$  to 1 and  $c_2$  to 2. The treatment effect  $b$  is heterogeneous:

$$b = \begin{cases} 1, & \text{if } W_i^1 > 0.75 \\ -1, & \text{if } W_i^1 \leq 0.75 \end{cases}$$

We generate a population of 30,000 individuals and randomly sample 5,000. Those among the 5,000 who are eligible for the RCT ( $S_i = 1$ ) are selected. Similarly, we sample 5,000 individuals from the population and select those who are not eligible for the RCT ( $S_i = 0$ ) to be our observational study participants.<sup>9</sup> We set each individual's treatment received  $D_i$  according to their treatment assignment and complier status and observe their responses  $Y_{isd}$ . In this design, the manner in which  $S_i$ ,  $T_i$ ,  $D_i$ ,  $C_i$ , and  $Y_{isd}$  are simulated ensures that Assumptions (1) – (7) hold.

In the assigned-treatment RCT group ( $S_i = 1, T_i = 1$ ), we train a gradient boosting algorithm (Friedman, 2001) on the covariates to predict who in the control group ( $S_i = 1, T_i = 0$ ) would comply with treatment ( $C_i = 1$ ), which is unobservable. These individuals *would have* complied had they been assigned to the treatment group. For this group of observed compliers to treatment and predicted compliers from the control group of the RCT, we estimate the response surface using gradient boosting with features  $(W_i^1, W_i^2, W_i^3)$  and  $D_i$ . The PATT-C is estimated according to the estimation procedure outlined above.

---

<sup>9</sup>This set-up mimics the reality that a population census is usually impossible.

## 4.2 Simulation results

We vary each of the parameters  $e_1$ ,  $e_2$ ,  $e_3$ ,  $e_4$ ,  $e_5$ , and  $e_6$  along a grid of five random standard normal values in order to generate different combinations of rates of compliance, treatment eligibility, RCT eligibility in the population, and confounding. For each possible combination of the six parameters, and holding all other parameters constant, we compute over 10 simulation runs the average root mean squared error (RMSE) between the true population average treatment effect and the PATT-C, PATT, or CACE estimates. Averaging across combinations, the unadjusted population estimator yields the highest average RMSE (1.06), followed by the CACE (0.89), and the PATT-C (0.76).

Figures 2 and A1 show the average RMSE of the estimators as a function of the population compliance rate and the share of population members eligible to participate in the RCT or the population treatment rate, respectively. The PATT estimator does not correct for bias resulting from noncompliance in the population and consequently performs poorly when the population compliance rate is relatively low (i.e.,  $\leq 60\%$ ). The PATT-C estimator corrects for noncompliance in the population and thus outperforms the PATT in low-compliance settings. The CACE corrects for noncompliance in the sample, and underperforms compared to the PATT-C due to differences between the sample and population.

Figure 3 compares the average RMSE of the estimators at varying levels of compliance in the population. The error for each of the estimators predictably decreases as a greater share of population members comply with treatment. PATT-C outperforms both PATT and CACE in terms of minimizing RMSE when the population compliance rate is below 90%. The PATT outperforms the PATT-C only at nearly perfect population compliance (i.e., 90% compliance rate), and the CACE outperforms the PATT when the population compliance rate is at 60% or below.

Figures A2, A3, and A4 plot the relationships between estimation error and the degrees of confounding in the mechanisms that determine compliance, treatment assignment, and sample selection, respectively. The estimation error of PATT-C is comparatively less in-



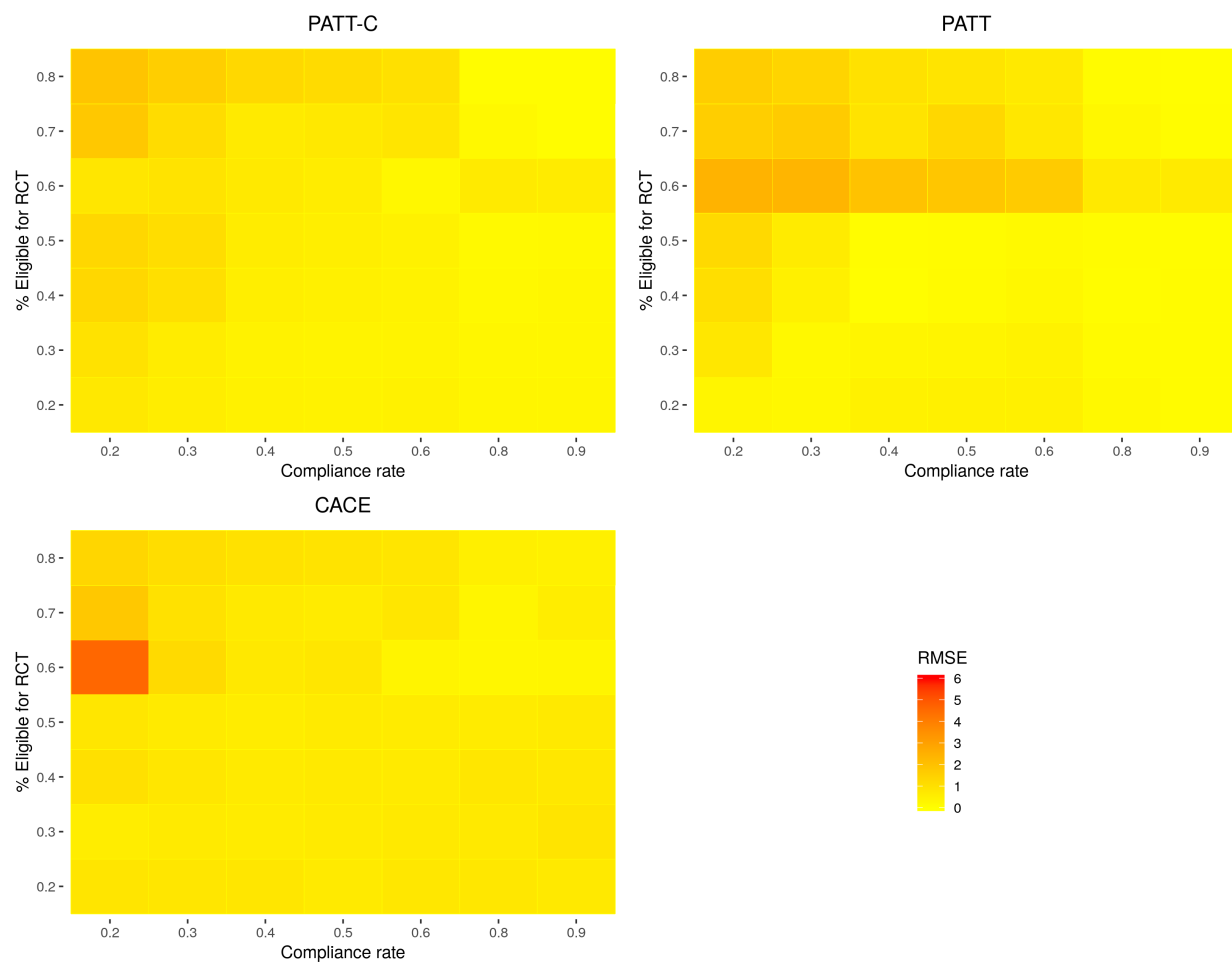


Figure 2: Average RMSE, binned by compliance rate and percent eligible for the RCT. Darker tiles correspond to higher errors and white tiles correspond to missing simulated data.

variant to increases in the degree of confounding in the three mechanisms compared to its unadjusted counterpart. The estimation error of CACE is generally more variable than that of the population estimators due to CACE's inability to account for differences between the sample and target population.

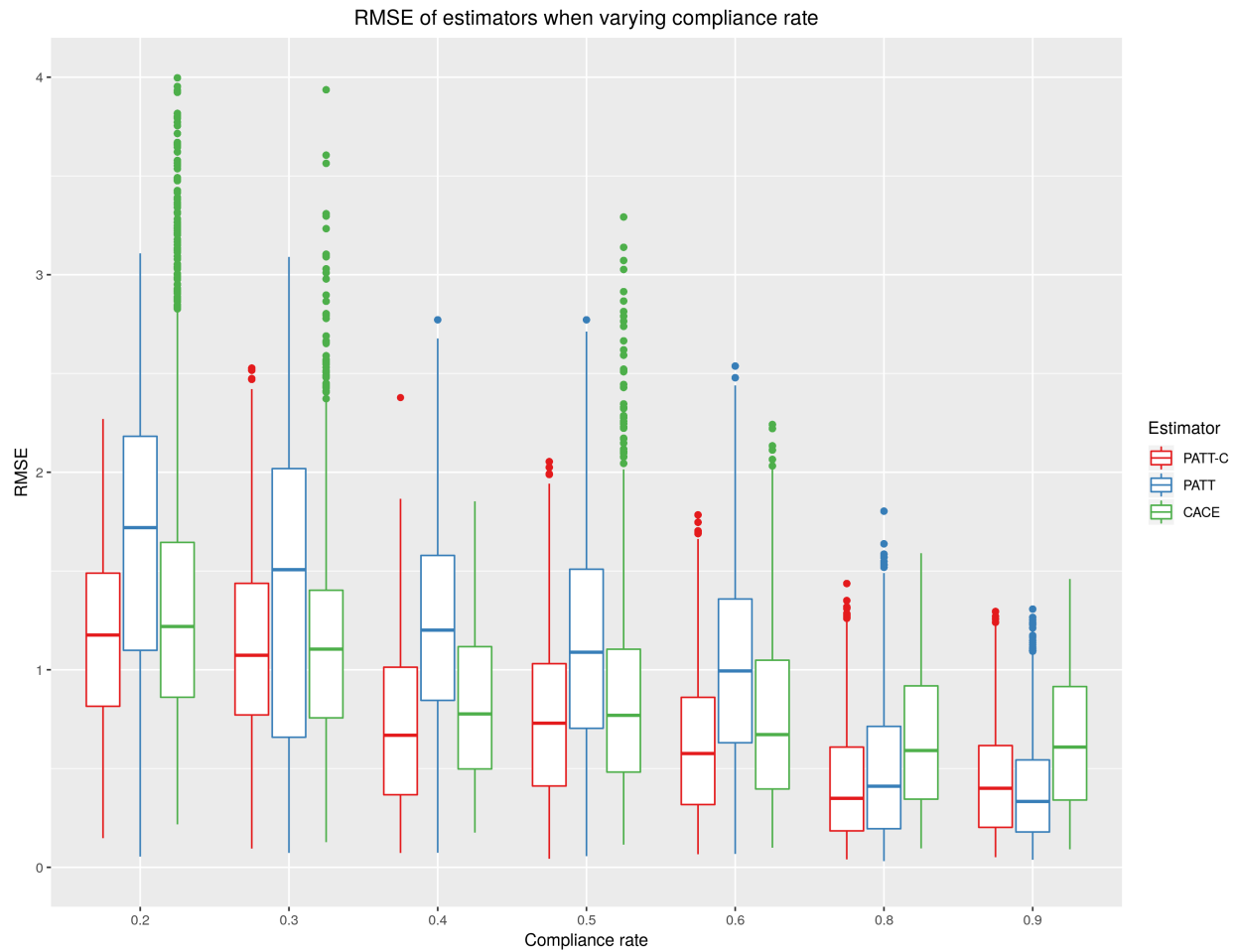


Figure 3: Average RMSE according to compliance rates in the population.

## 5 Application: Medicaid and health care use

We apply the proposed estimator to measure the effect of Medicaid coverage on health care use for a target population of adults who may benefit from expansions to the Medicaid program. In particular, we examine the population of nonelderly adults in the U.S. with household incomes at or below 138% of the Federal Poverty Level (FPL) — which amounts to \$32,913 for a four-person household in 2014 — who may be eligible for Medicaid following the Affordable Care Act (ACA) expansion.

We draw RCT data from the Oregon Health Insurance Experiment (OHIE) (Finkelstein et al., 2012; Taubman et al., 2014), **which randomly assigned Medicaid coverage to the uninsured and examined their subsequent health care use.** Subsequent research calls in to question the external validity of the OHIE, which resulted in the counterintuitive finding that Medicaid increased ER use among RCT participants. For example, quasi-experimental studies on the impact of the 2006 Massachusetts health reform — which served as a model for the ACA — show that ER use decreased or remained constant following the reform (Miller, 2012; Kolstad and Kowalski, 2012). A challenge to the external validity of the OHIE is that its exclusion criteria was likely more restrictive than government health insurance expansions.

### 5.1 RCT sample

In 2008, **a group of uninsured low-income adults participated in the OHIE for the chance to apply to receive health insurance through a state Medicaid program.** In line with program eligibility requirements, participants were restricted to Oregon residents aged 19 to 64 who were not otherwise eligible for public insurance, who had been without insurance for six months, had income below the FPL, and held assets below \$2,000. Treatment assignment occurred at the household level: participants selected by the lottery won the opportunity for themselves and any household member to apply for Medicaid. Within a sample of 74,922

individuals representing 66,385 households, 29,834 participants were selected by the lottery; the remaining 45,008 participants served as controls in the experiment.

Participants in selected households **were enrolled in Medicaid** if they returned an enrollment application within 45 days of receipt. **Only 30% of participants in selected households successfully enrolled.** The low compliance rate is primarily due to failure to return an application or demonstrate income below the FPL. Compliance is measured using a binary variable indicating whether the participant was enrolled in any Medicaid program during the study period.

We include as covariates in our response and complier models (S.1 and S.3, respectively) pretreatment information on participant age, race, gender, education, marital status, number of children in the household, employment status, health status, and household income. We also include indicator variables on household size because lottery selection was random conditional on the number of household members. Because treatment occurs at the household level, all analyses cluster standard errors at the household level.

The response data originate from a mail survey **containing questions about health insurance and health care use that elicited responses from 23,741 OHIE participants.**<sup>10</sup> The response variables measure health care use in terms of the number of ER and primary care (i.e., outpatient) visits in the past six months.

## 5.2 Observational data

We acquire data on the target population from the National Health Interview Study (NHIS) for years 2008 to 2017.<sup>11</sup> We restrict the sample to respondents with income below 138% of the FPL and who are uninsured or on Medicaid and select covariates on respondent characteristics that match the OHIE pretreatment covariates. We use a recoded variable

---

<sup>10</sup>Following Finkelstein et al. (2012), indicator variables for survey wave and interactions with household size indicators are also included as covariates in the response and complier models because the proportion of treated participants varies across the survey waves.

<sup>11</sup>A possible limitation of this application is that it ignores the complex sampling techniques of the NHIS sample design such as differential sampling, which is discussed in detail in Parsons et al. (2014).

that indicates whether respondents are on Medicaid as an analogue to the OHIE compliance measure. The outcomes of interest from the NHIS are based on questions that are virtually identical to the OHIE mail survey questions, except that the utilization questions in the NHIS are asked with a 12 month rather than a 6 month look-back period.<sup>12</sup>

### 5.3 Verifying assumptions

In order for  $\tau_{\text{PATT-C}}$  to be identified, Assumptions (1) – (7) must be met. Assumption (1) ensures that potential outcomes for participants in the target population would be identical to their outcomes in the RCT if they had been randomly assigned their observed treatment. In the empirical application, Medicaid coverage for uninsured individuals was applied in the same manner in the RCT as it is in the population. Differences in potential outcomes due to sample selection might arise, however, if there are differences in the mail surveys used to elicit health care use responses between the RCT and the nonrandomized study.

We cannot directly test Assumptions (3) and (4), which state that potential outcomes for treatment and control are independent of sample assignment for individuals with the same covariates and assignment to treatment. The assumptions are only met if every possible confounder associated with the response and the sample assignment is accounted for. In estimating the response surface, we use all demographic, socioeconomic, and pre-existing health condition data that were common in the OHIE and NHIS data. Potentially important unobserved confounders include the number of hospital and outpatient visits in the previous year, proximity to health services, and enrollment in other federal programs.

The final two columns of Table A1 compares RCT participants selected for Medicaid with population members on Medicaid. Compared to the RCT compliers, the population members who received treatment are predominantly female, younger, more racially and ethnically diverse, less educated, and live in higher income households. Diagnoses of diabetes, asthma, high blood pressure, and heart disease are more common among the population on Medicaid

---

<sup>12</sup>Following Finkelstein et al. (2012), we resolve this discrepancy by halving the NHIS responses in order to make them comparable to the OHIE outcomes.

then the RCT treated. These summary statistics and the analyses that follow use survey weights to account for the probability of being sampled and non-response.

Strong ignorability assumptions may also be violated due to the fact that the OHIE applied a more stringent exclusion criteria compared to the NHIS sample. While the RCT and population sample both screened for individuals below the FPL, only the RCT required those enrolled to recertify their household income eligibility during the study period. Strong ignorability would not hold if the failure to recertify is correlated with unobserved variables.

Following Assumption 5, we assume no interference between households in the OHIE because treatment assignment occurred at the household level. Within-household interference is not possible in this RCT because household members share the same treatment status. Interference between households would threaten the no-interference assumption in the unlikely case that the Medicaid coverage of individuals in treated households affects the health care use of individuals in households assigned to control.

Assumption (2) is violated if assignment to treatment influences the compliance status of individuals with the same covariates. The compliance ensemble can accurately classify compliance status for 77% of treated RCT participants with only the covariates — and not treatment assignment — as model inputs.<sup>13</sup> This gives evidence in favor of the conditional independence assumption.

The exclusion restriction assumption (7) ensures treatment assignment affects the response only through enrollment in Medicaid. It is reasonable that a person’s enrollment in Medicaid, not just their eligibility to enroll, would affect their hospital use. For private health insurance one might argue that eligibility may be negatively correlated with hospital use, as people with pre-existing conditions are less often eligible yet go to the hospital more frequently. This should not be the case with a federally funded program such as Medicaid.

---

<sup>13</sup>The compliance ensemble is evaluated in terms of 10-fold cross-validated MSE. The distribution of MSE for the ensemble and its candidate algorithms are provided in Table A5.

### 5.3.1 Placebo tests

Similar to the procedure proposed by (Hartman et al., 2015), we conduct placebo tests to check whether the average outcomes differ between the RCT compliers on Medicaid and the adjusted **population members who received Medicaid**.<sup>14</sup> If the placebo tests detect a significant difference between the mean outcomes of these groups, it would indicate that either Assumption (1) (for  $d = 1$ ), or Assumptions (3) and (4) are violated.

Table A3 reports the results of placebo tests, comparing the mean outcomes of RCT compliers against the mean outcomes of adjusted **population members who received treatment**. The former quantity is calculated from the observed RCT sample and the latter quantity is the mean counterfactual  $Y_{i11}$  estimated from S.4 of the estimation procedure. **In addition to per-comparison  $p$ -values, we report family-wise  $p$ -values to adjust for multiple comparisons.**<sup>15</sup>

Tests of equivalence between the two groups indicate that the differences across each outcome are not statistically significant. These results imply that the PATT-C estimator is not biased by differences in how Medicaid is delivered or health outcomes are measured between the RCT and population, or by differences in sample or population members' unobserved characteristics.

### 5.3.2 Sensitivity to no defiers assumption

Angrist et al. (1996) show that the bias due to violations of Assumption (6) is equivalent to the difference of average causal effects of treatment received for compliers and defiers, multiplied by the relative proportion of defiers,  $\mathbb{P}(i \text{ is a defier})/(\mathbb{P}(i \text{ is a complier})) - \mathbb{P}(i \text{ is a defier})$ .

Table A2 reports the distribution of participants in the OHIE by status of treatment

---

<sup>14</sup>Note that a placebo test for Assumption (2) is not possible because we never observe whether RCT controls would actually take-up treatment if assigned.

<sup>15</sup>In the context of the placebo tests, the family-wise  $p$ -value is the probability of rejecting the null hypothesis of no difference between mean outcomes on a given outcome under the family of null hypotheses on any outcome in the domain of health care utilization. We calculate family-wise  $p$ -values based on 10,000 bootstrap iterations of Westfall et al.'s [1993] step-down procedure.

assignment and treatment received. Assumption (6) does not hold due to the presence of defiers; i.e., participants who were assigned to control and enrolled in Medicaid during the study period. About 6.7% of the RCT sample were assigned to control but were enrolled in Medicaid ( $T_i < D_i$ ) and 65.5% of the sample complied with treatment assignment ( $D_i = T_i$ ), which results in a bias multiplier of 0.11. Suppose that the difference of average causal effects of Medicaid received on ER use for compliers and defiers is 1.2%. The resulting bias is only 0.1%, which would not meaningfully alter the interpretation of the CACE or PATT-C estimates reported below.

## 5.4 Empirical results

We compare PATT-C and PATT estimates for ER and outpatient use. We obtain estimates for the overall group of participants and subgroups according to sex, age, race, health status, education, and household income. Subgroup treatment effects are estimated by taking differences across response surfaces for a given covariate subgroup, and response surfaces are estimated with the ensemble mean predictions. We use treatment received, number of household members, and the subgroup covariates as features in the response models. We generate 95% confidence intervals for these estimates using 1,000 bootstrap samples.

Table A4 presents the PATT-C estimates, which indicate that Medicaid coverage has a positive, but considerably smaller effect on the number of ER and outpatient visits. For comparison, Finkelstein et al. (2012) reports population estimates of the effect of Medicaid coverage on the number of ER and out-patient visits using 2004–2009 NHIS data on adults aged 19–64 below 100 percent of the federal poverty line ( $n = 15,528$ ). Finkelstein et al. (2012) estimates Medicaid coverage significantly increases the number of ER visits by 0.08 [0.05, 0.12] and increases the number of outpatient visits by 1.45 [1.33, 1.57].

Figures A5, A6, and A7 examine heterogeneous treatment effect estimates on ER and outpatient use in the population. While this study is the first to our knowledge to estimate heterogeneity in treatment effects for the target population, Taubman et al. (2014) and



Kowalski (2016) perform subgroup analyses on the RCT sample. Similar to the PATT-C estimates, Taubman et al.’s [2014] subgroup analyses indicate that increases in ER use due to Medicaid are significantly larger for younger individuals and those with high school-level education.<sup>16</sup>

## 6 Discussion

The simulation results presented in Section 4 show that the PATT-C estimator outperforms its unadjusted counterpart when the compliance rate is low. Of course, the simulation results depend on the particular way we parameterized the compliance, selection, treatment assignment, and response schedules.

In particular, the strength of correlation between the covariates and compliance governs how well the estimator will perform, since S.1 of the estimation procedure is to predict who *would be* a complier in the RCT control group, had they been assigned to treatment. If it is difficult to predict compliance using the observed covariates, then the estimator will perform badly because of noise introduced by incorrectly treating noncompliers as compliers. Further research should be done into ways to test how well the model of compliance works in the population or explore models to more accurately predict compliance in RCTs. Accurately predicting compliance is not only essential for yielding unbiased estimates of the average causal effects for target populations, it is also useful for researchers and policymakers to know which groups of individuals are unlikely to comply with treatment.

In the OHIE trial, only about 30% of those selected to receive Medicaid benefits actually enrolled. The compliance ensemble accurately classified compliance status for 77% of treated RCT participants using only the pretreatment covariates as features. While we don’t know how well the compliance ensemble predicts for the control group, the control group should be similar to the treatment group on pretreatment covariates because of the RCT

---

<sup>16</sup>Kowalski (2016) perform subgroup analyses on OHIE sample data and find larger increases in ER use as a result of Medicaid for men, English speakers, and individuals enrolled in a food stamp program prior to the lottery.

randomization. The model’s performance on the training set suggests that compliance is not purely random and depends on observed covariates. This gives evidence in favor of using the proposed estimator.

In the empirical application, the sample population differs in several dimensions from the target population of individuals who will be covered by other Medicaid expansions, such as the ACA expansion to cover all adults up to 138% of the FPL. For instance, the RCT participants are disproportionately white urban-dwellers (Taubman et al., 2014). The RCT participants volunteered for the study and therefore may be in poorer health compared to the target population. These differences in baseline covariates make reweighting or response surface methods necessary to extend the RCT results to the population.

Explicitly modeling compliance allows us to decompose population estimates by subgroup according to pretreatment covariates common to both RCT and observational datasets; e.g, demographic variables, pre-existing conditions, and insurance coverage. We find substantial differences between sample and population estimates in terms of race, education, and health status subgroups. This pattern is expected because RCT participants volunteered for the study and are predominately white and educated.

## References

- Angrist, J. D., G. W. Imbens, and D. B. Rubin (1996, June). Identification of causal effects using instrumental variables. *Journal of the American Statistical Association* 91(434), 444–455.
- Balke, A. and J. Pearl (1997). Bounds on treatment effects from studies with imperfect compliance. *Journal of the American Statistical Association* 92(439), 1171–1176.
- Breiman, L. (2001). Random forests. *Machine Learning* 45(1), 5–32.
- Buja, A., T. Hastie, and R. Tibshirani (1989). Linear smoothers and additive models. *The Annals of Statistics*, 453–510.
- Finkelstein, A., S. Taubman, B. Wright, M. Bernstein, J. Gruber, J. P. Newhouse, and H. Allen (2012). The oregon health insurance experiment: Evidence from the first year. *The Quarterly Journal of Economics* 127(3), 1057.
- Frangakis, C. E. and D. B. Rubin (2002). Principal stratification in causal inference. *Biometrics* 58(1), 21–29.
- Freedman, D. A. (2006). Statistical models for causation what inferential leverage do they provide? *Evaluation Review* 30(6), 691–713.
- Friedman, J. H. (2001). Greedy function approximation: A gradient boosting machine. *Annals of Statistics*, 1189–1232.
- Frumento, P., F. Mealli, B. Pacini, and D. B. Rubin (2012). Evaluating the effect of training on wages in the presence of noncompliance, nonemployment, and missing outcome data. *Journal of the American Statistical Association* 107(498), 450–466.
- Hartman, E., R. Grieve, R. Ramsahai, and J. S. Sekhon (2015). From sate to patt: Combining experimental with observational studies to estimate population treatment effects. *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 10, 1111.

- Imai, K., G. King, and E. Stuart (2008). Misunderstandings among experimentalists and observationalists about causal inference. *Journal of the Royal Statistical Society, Series A 171, part 2*, 481–502.
- Imai, K., D. Tingley, and T. Yamamoto (2013). Experimental designs for identifying causal mechanisms. *Journal of the Royal Statistical Society: Series A (Statistics in Society) 176*(1), 5–51.
- Imbens, G. W. and J. D. Angrist (1994). Identification and estimation of local average treatment effects. *Econometrica* 62(2), 467–475.
- Imbens, G. W. and D. B. Rubin (1997). Bayesian inference for causal effects in randomized experiments with noncompliance. *The Annals of Statistics*, 305–327.
- Kolstad, J. T. and A. E. Kowalski (2012). The impact of health care reform on hospital and preventive care: evidence from Massachusetts. *Journal of Public Economics* 96(11-12), 909–929.
- Kowalski, A. E. (2016, June). Doing more when you’re running LATE: Applying marginal treatment effect methods to examine treatment effect heterogeneity in experiments. Working Paper 22363, National Bureau of Economic Research.
- Miller, S. (2012). The effect of insurance on emergency room visits: an analysis of the 2006 Massachusetts health reform. *Journal of Public Economics* 96(11-12), 893–908.
- Parsons, V. L., C. L. Moriarity, K. Jonas, T. F. Moore, K. E. Davis, and L. Tompkins (2014). Design and estimation for the National Health Interview Survey, 2006–2015. *National Center for Health Statistics. Vital and Health Statistics* 2(165), 1–53.
- Polley, E. C. and M. J. Van Der Laan (2010). Super learner in prediction. Working Paper 266, Division of Biostatistics, University of California, Berkeley.

- Rothwell, P. M. (2005). External validity of randomised controlled trials: “to whom do the results of this trial apply?”. *The Lancet* 365(9453), 82–93.
- Rubin, D. B. (1990). Comment: Neyman (1923) and causal inference in experiments and observational studies. *Statistical Science* 5(4), 472–480.
- Stuart, E. A., S. R. Cole, C. P. Bradshaw, and P. J. Leaf (2011). The use of propensity scores to assess the generalizability of results from randomized trials. *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 174(2), 369–386.
- Taubman, S., H. Allen, B. Wright, K. Baicker, and A. Finkelstein (2014, January). Medicaid increases emergency department use: Evidence from oregon’s health insurance experiment. *Science* 343(6168), 263–268.
- Tibshirani, R., J. Bien, J. Friedman, T. Hastie, N. Simon, J. Taylor, and R. J. Tibshirani (2012). Strong rules for discarding predictors in lasso-type problems. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 74(2), 245–266.
- van der Laan, M. J., E. C. Polley, and A. E. Hubbard (2007). Super learner. *Statistical Applications in Genetics and Molecular Biology* 6(1).
- Westfall, P. H., S. S. Young, et al. (1993). *Resampling-based multiple testing: Examples and methods for p-value adjustment*, Volume 279. John Wiley & Sons.
- Yau, L. H. and R. J. Little (2001). Inference for the complier-average causal effect from longitudinal data subject to noncompliance and missing data, with application to a job training assessment for the unemployed. *Journal of the American Statistical Association* 96(456), 1232–1244.

# Appendix

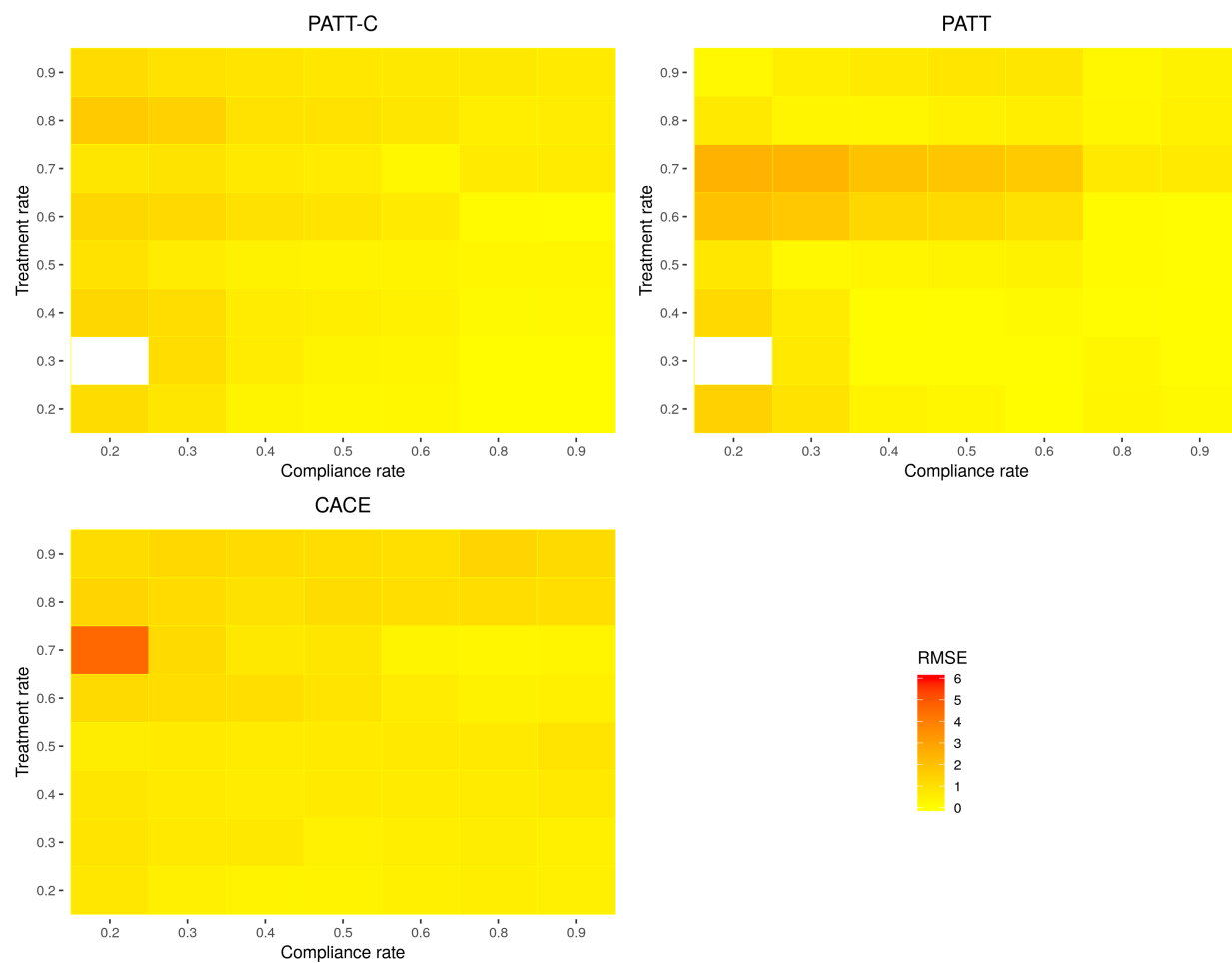


Figure A1: Average RMSE binned by compliance rate and treatment rate.

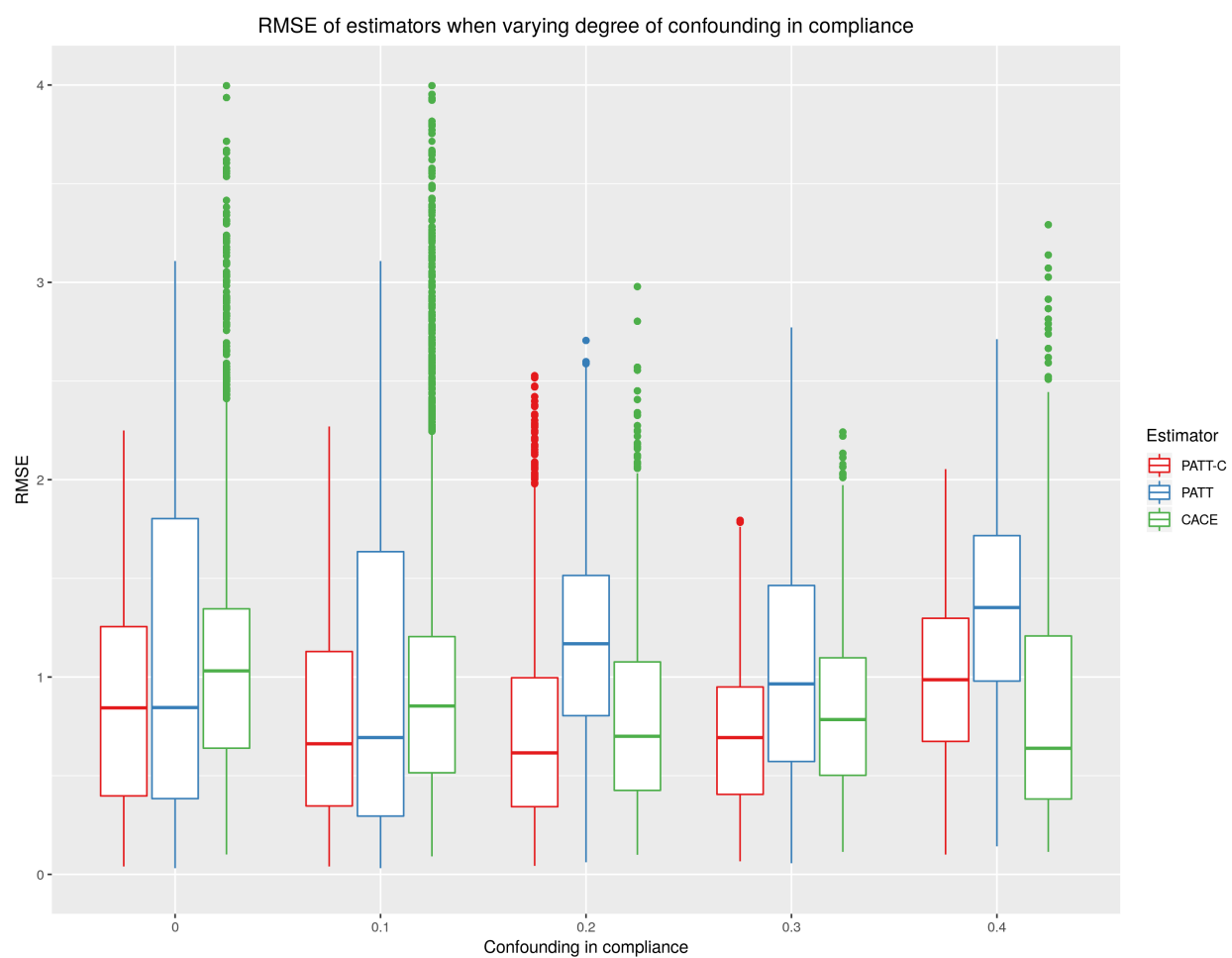


Figure A2: Average RMSE according to degree of confounding in compliance.

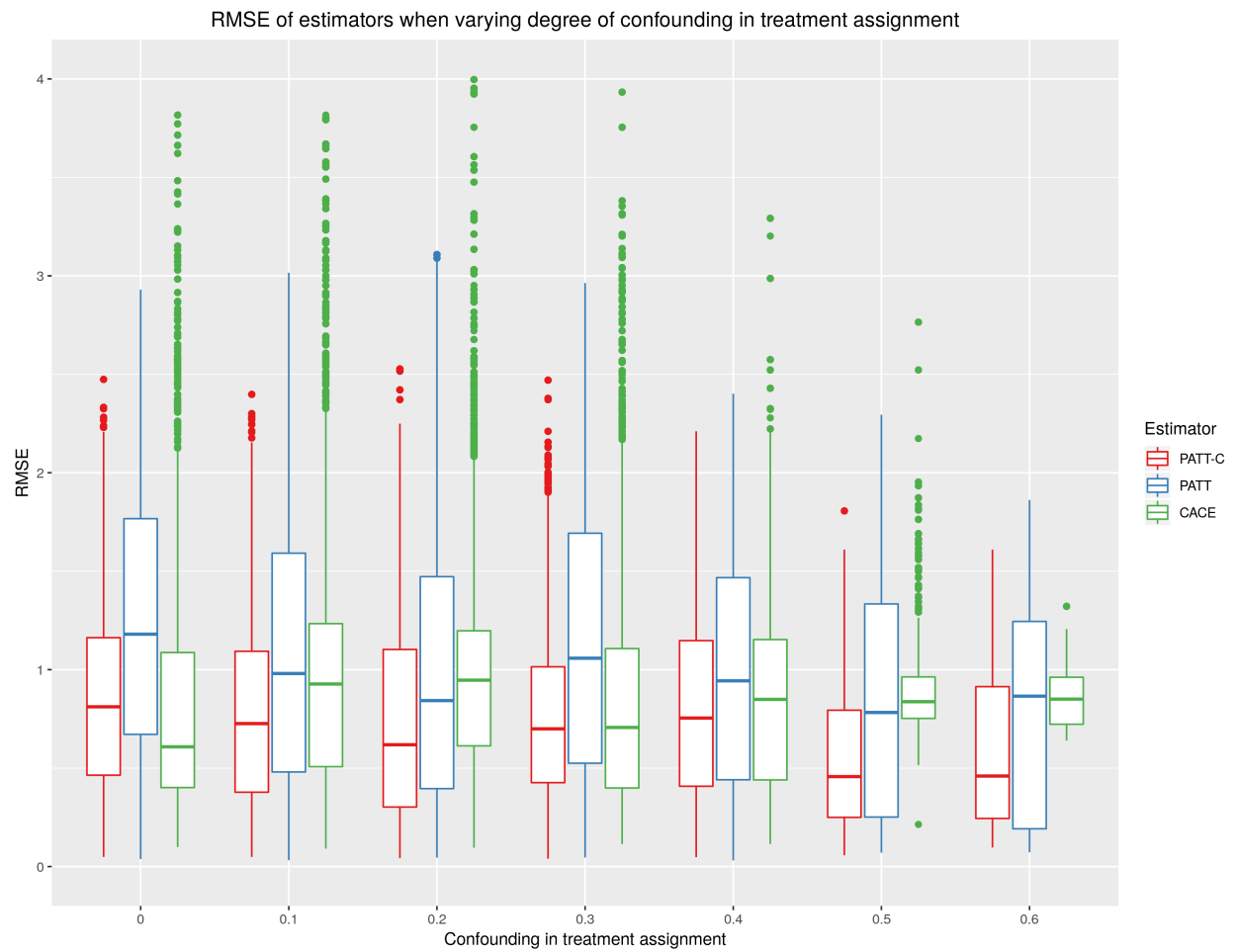


Figure A3: Average RMSE according to degree of confounding in treatment assignment.



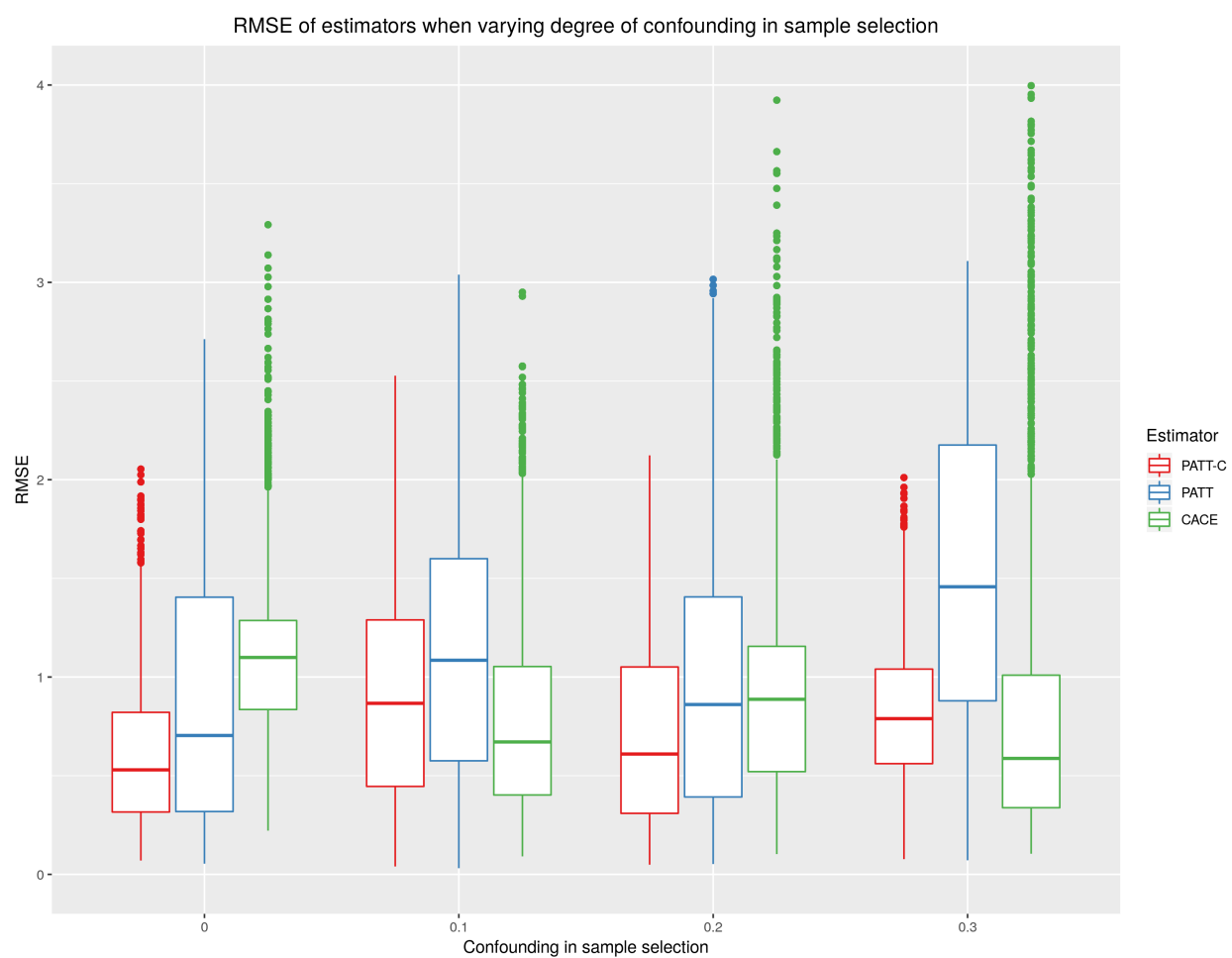


Figure A4: Average RMSE according to degree of confounding in sample selection.

Table A1: Pretreatment covariates and responses for OHIE and NHIS respondents by Medicaid coverage status.

	OHIE no Medicaid <i>n</i> = 4,519		OHIE Medicaid <i>n</i> = 6,100		NHIS Medicaid <i>n</i> = 6,261	
<b>Covariate</b>	<b>n</b>	<b>%</b>	<b>n</b>	<b>%</b>	<b>n</b>	<b>%</b>
<i>Sex:</i>						
Female	2,538	56.2	3506	57.5	4,288	68.5
<i>Age:</i>						
19-49	1,288	28.5	1,625	26.6	4324	69.1
50-64	3,231	71.5	4,475	73.4	1,937	30.9
<i>Race:</i>						
White	3,956	87.5	5,183	85.0	3,902	62.3
Black	193	4.3	247	4.0	1,723	27.5
Hispanic	264	5.8	538	8.8	1,570	25.1
<i>Health status:</i>						
Diabetes	459	10.2	637	10.4	866	13.8
Asthma	823	18.2	1,094	17.9	1272	20.3
High blood pressure	1,362	30.1	1,705	27.9	2,166	34.6
Heart condition	120	2.7	189	3.1	529	8.4
<i>Education:</i>						
Less than high school	858	19.0	1,154	18.9	1,942	31.0
High school diploma or GED	2,589	57.3	3,279	53.8	2,076	33.2
Voc. training / 2-year degree	804	17.8	1,186	19.4	1,810	28.9
4-year college degree or more	268	5.9	481	7.9	433	6.9
<i>Income:</i>						
< \$10k	4,518	100.0	4,111	67.4	2,588	41.3
\$10k-\$25k	1	0.0	1,616	26.5	3,098	49.5
> \$25k	0	0.0	373	6.1	575	9.2
<b>Responses</b>	<b><math>\bar{x}</math></b>	<b>sd</b>	<b><math>\bar{x}</math></b>	<b>sd</b>	<b><math>\bar{x}</math></b>	<b>sd</b>
# ER visits	0.44	0.95	0.44	0.99	0.48	1.0
# outpatient visits	1.9	3.01	1.9	2.8	2.08	2.3

Notes: weighted using OHIE and NHIS survey weights.

Table A2: Distribution of OHIE participants by status of treatment assignment ( $T_i$ ) and treatment received ( $D_i$ ).

	$D_i = 0$	$D_i = 1$	n
$T_i = 0$	10,010	1,556	11,566
$T_i = 1$	6,446	5,193	11,639
n	16,456	6,749	23,205

Notes: weighted using OHIE survey weights.

Table A3: Placebo test results comparing the mean outcomes of RCT compliers and adjusted **population members who received treatment**.

Outcome	RCT complier mean	Adjusted population mean	Difference	Per-comparison $p$ -value	Family-wise $p$ -value
# ER vists	0.45	0.45	0.002	0.85	
# outpatient visits	1.90	1.94	-0.03	0.26	

Notes:  $p$ -values for survey-weighted difference-in-means calculated from two-sided t-test. Family-wise  $p$ -value calculated across the two different measures of health care utilization.

Table A4: Comparison of population and sample estimates.

Outcome Estimator	Any ER visit	# ER visits	# outpatient visits
PATT-C	0.0001 [0.0001, 0.0001]	0.0005 [0.0002, 0.0008]	0.002 [0.002, 0.002]
PATT	0.0006 [0.0005, 0.0007]	-0.004 [-0.005, -0.004]	0.02 [0.02, 0.02]
CACE	-0.001 [-0.02, 0.02]	0.005 [-0.05, 0.06]	-0.02 [-0.19, 0.14]

Notes: Estimates in brackets represent 95% bootstrap confidence intervals constructed with 1,000 bootstrap samples.

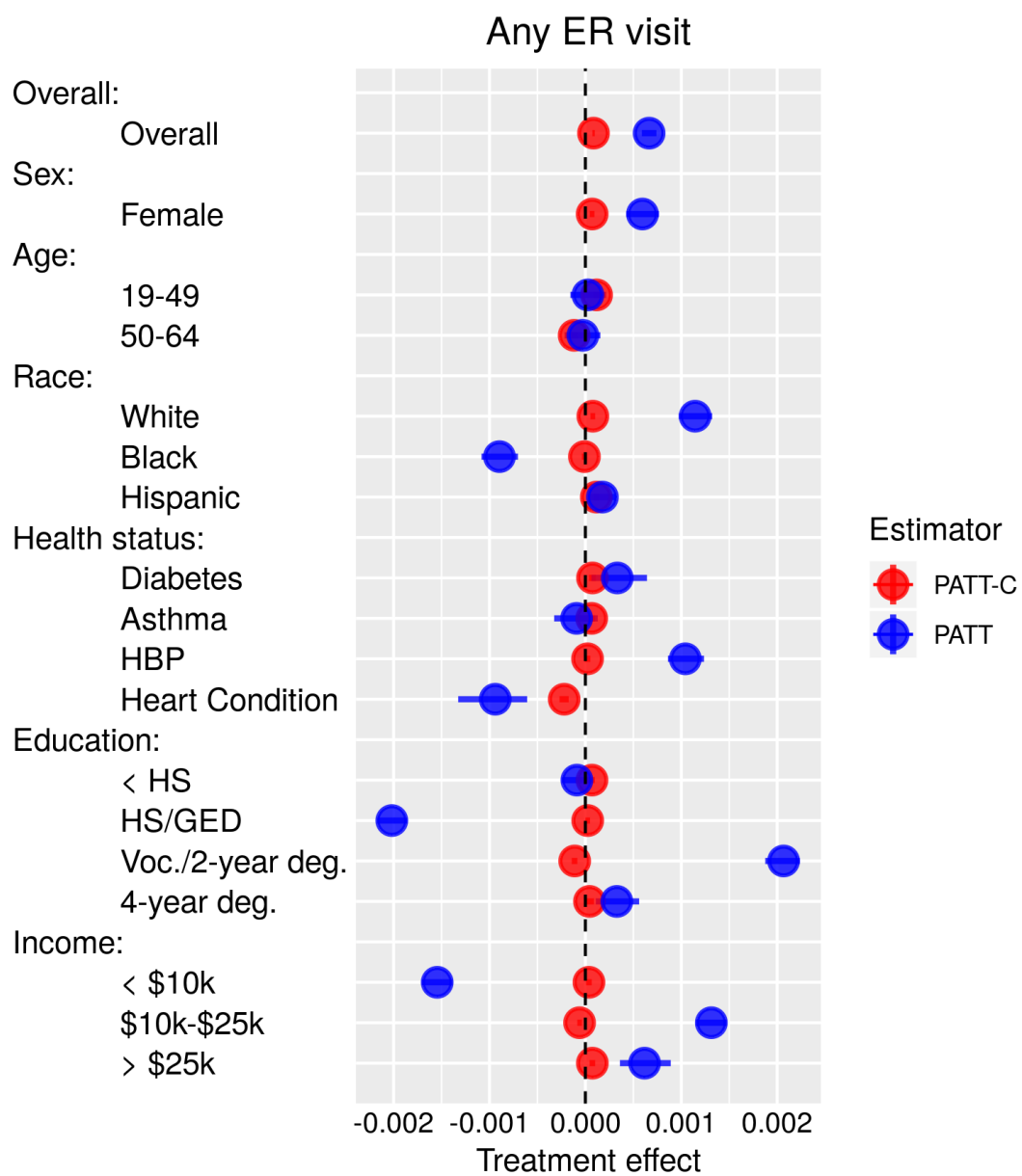


Figure A5: Heterogeneity in sample and population treatment effect estimates: any ER visit. Horizontal lines represent 95% bootstrap confidence intervals constructed with 1,000 bootstrap samples.



Figure A6: Heterogeneity in population treatment effect estimates: # ER visits.

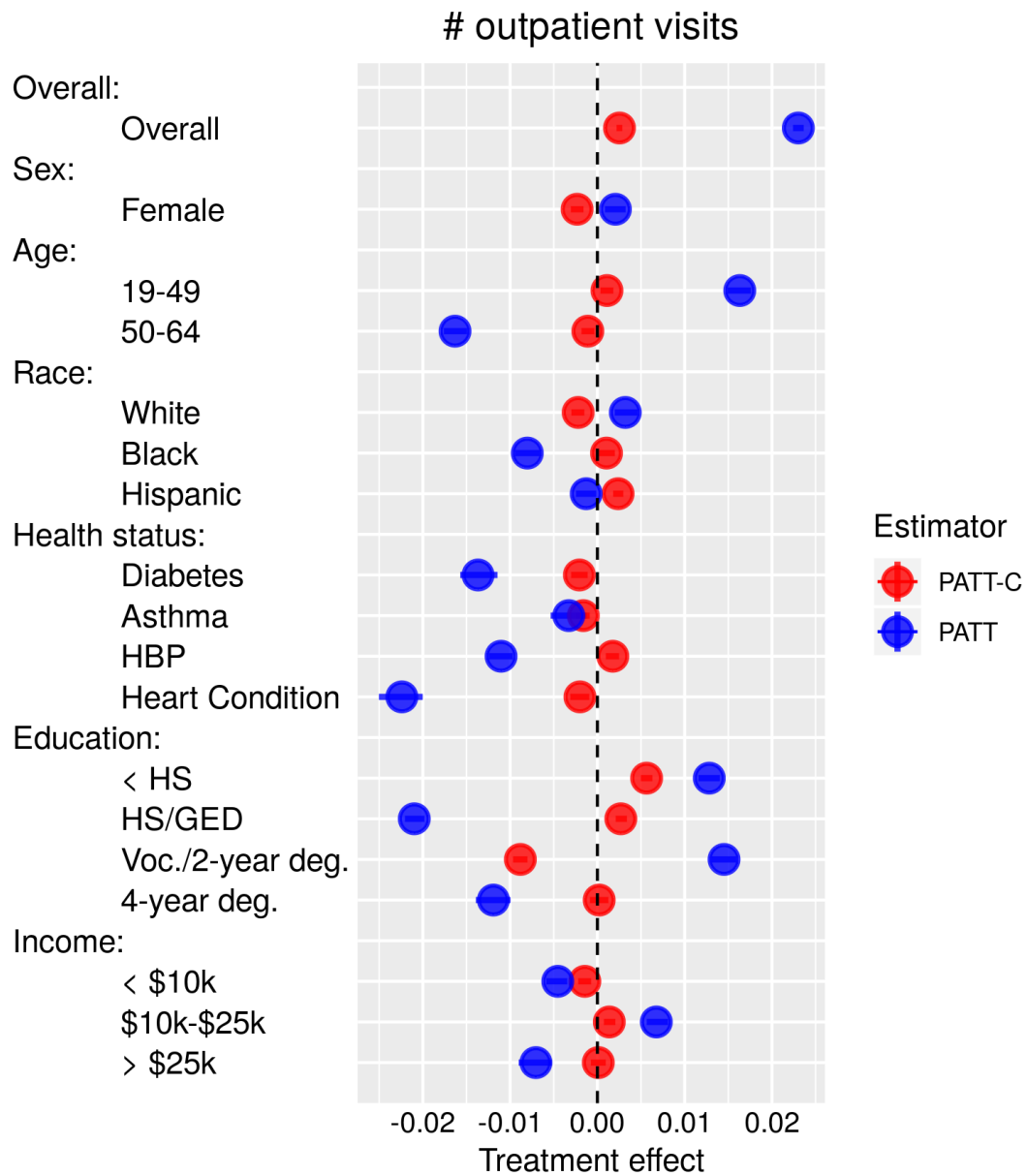


Figure A7: Heterogeneity in population treatment effect estimates: # outpatient visits.

Table A5: Distribution of MSE for compliance ensemble.

Algorithm	Mean	SE	Min.	Max.
Super learner ( <b>SuperLearner</b> )	0.22	0.001	0.21	0.23
Lasso regression ( <b>glmnet</b> )	0.22	0.001	0.21	0.23
Random forests ( <b>randomForest</b> )	0.27	0.002	0.25	0.29
Ridge regression ( <b>glmnet</b> )	0.22	0.001	0.21	0.23

Notes: MSE is 10-fold cross-validated error for super learner ensemble and candidate algorithms. R package used for implementing each algorithm in parentheses.

Table A6: Error and weights for candidate algorithms in response ensemble for RCT compliers.

<b>Any ER visit</b>		
Algorithm	MSE	Weight
Lasso regression ( <code>glmnet</code> )	0.18	1
Random forests, $\#preds. = 1$ ( <code>randomForest</code> )	0.25	0
Random forests, $\#preds. = 10$ ( <code>randomForest</code> )	0.24	0
Regularized logistic regression, $\alpha = 0.25$ ( <code>glmnet</code> )	0.19	0
Regularized logistic regression, $\alpha = 0.5$ ( <code>glmnet</code> )	0.19	0
Regularized logistic regression, $\alpha = 0.75$ ( <code>glmnet</code> )	0.19	0
Ridge regression ( <code>glmnet</code> )	0.18	0
<b># ER visits</b>		
Algorithm	MSE	Weight
Additive regression, degree = 3 ( <code>gam</code> )	0.95	0
Additive regression, degree = 4 ( <code>gam</code> )	0.95	0
Lasso regression ( <code>glmnet</code> )	0.95	0.92
Random forests, $\#preds. = 1$ ( <code>randomForest</code> )	0.95	0
Random forests, $\#preds. = 10$ ( <code>randomForest</code> )	0.99	0.08
Regularized linear regression, $\alpha = 0.25$ ( <code>glmnet</code> )	0.95	0
Regularized linear regression, $\alpha = 0.5$ ( <code>glmnet</code> )	0.95	0
Regularized linear regression, $\alpha = 0.75$ ( <code>glmnet</code> )	0.95	0
Ridge regression ( <code>glmnet</code> )	0.18	0
<b># outpatient visits</b>		
Algorithm	MSE	Weight
Additive regression, degree = 3 ( <code>gam</code> )	8.40	0
Additive regression, degree = 4 ( <code>gam</code> )	8.40	0
Lasso regression ( <code>glmnet</code> )	8.38	0
Random forests, $\#preds. = 1$ ( <code>randomForest</code> )	8.38	0
Random forests, $\#preds. = 10$ ( <code>randomForest</code> )	8.79	0.08
Regularized linear regression, $\alpha = 0.25$ ( <code>glmnet</code> )	8.38	0
Regularized linear regression, $\alpha = 0.5$ ( <code>glmnet</code> )	8.38	0
Regularized linear regression, $\alpha = 0.75$ ( <code>glmnet</code> )	8.38	0.92
Ridge regression ( <code>glmnet</code> )	8.38	0

Notes: cross-validated error and weights used for each algorithm in super learner ensemble. *MSE* is the ten-fold cross-validated mean squared error for each algorithm. *Weight* is the coefficient for the Super Learner, which is estimated using non-negative least squares based on the Lawson-Hanson algorithm. R package used for implementing each algorithm in parentheses.  $\#preds.$  is the number of predictors randomly sampled as candidates in each decision tree in random forests algorithm.  $\alpha$  is a parameter that mixes L1 and L2 norms. degree is the smoothing term for smoothing splines.



Table A7: Error and weights for candidate algorithms in response ensemble for all RCT participants.

<b>Any ER visit</b>		
Algorithm	MSE	Weight
Lasso regression ( <b>glmnet</b> )	0.18	0.96
Random forests, $\#preds. = 1$ ( <b>randomForest</b> )	0.25	0
Random forests, $\#preds. = 10$ ( <b>randomForest</b> )	0.24	0.04
Regularized logistic regression, $\alpha = 0.25$ ( <b>glmnet</b> )	0.18	0
Regularized logistic regression, $\alpha = 0.5$ ( <b>glmnet</b> )	0.18	0
Regularized logistic regression, $\alpha = 0.75$ ( <b>glmnet</b> )	0.18	0
<b># ER visits</b>		
Algorithm	MSE	Weight
Additive regression, degree = 3 ( <b>gam</b> )	0.94	0
Additive regression, degree = 4 ( <b>gam</b> )	0.94	0
Lasso regression ( <b>glmnet</b> )	0.93	0.88
Random forests, $\#preds. = 1$ ( <b>randomForest</b> )	0.93	0
Random forests, $\#preds. = 10$ ( <b>randomForest</b> )	0.97	0.11
Regularized linear regression, $\alpha = 0.25$ ( <b>glmnet</b> )	0.93	0
Regularized linear regression, $\alpha = 0.5$ ( <b>glmnet</b> )	0.93	0
Regularized linear regression, $\alpha = 0.75$ ( <b>glmnet</b> )	0.93	0
Ridge regression ( <b>glmnet</b> )	0.93	0
<b># outpatient visits</b>		
Algorithm	MSE	Weight
Additive regression, degree = 3 ( <b>gam</b> )	8.42	0
Additive regression, degree = 4 ( <b>gam</b> )	8.42	0
Lasso regression ( <b>glmnet</b> )	8.41	0
Random forests, $\#preds. = 1$ ( <b>randomForest</b> )	8.41	0.99
Random forests, $\#preds. = 10$ ( <b>randomForest</b> )	8.79	0.01
Regularized linear regression, $\alpha = 0.25$ ( <b>glmnet</b> )	8.41	0
Regularized linear regression, $\alpha = 0.5$ ( <b>glmnet</b> )	8.41	0
Regularized linear regression, $\alpha = 0.75$ ( <b>glmnet</b> )	8.41	0
Ridge regression ( <b>glmnet</b> )	8.41	0

See notes to Fig.A6.