

# Lasso adjustments of treatment effect estimates in randomized experiments

Adam Bloniarz<sup>a,1</sup>, Hanzhong Liu<sup>a,1</sup>, Cun-Hui Zhang<sup>b</sup>, Jasjeet S. Sekhon<sup>a,c</sup>, and Bin Yu<sup>a,d,2</sup>

<sup>a</sup>Department of Statistics, University of California, Berkeley, CA 94720; <sup>b</sup>Department of Statistics and Biostatistics, Rutgers University, Piscataway, NJ 08854; <sup>c</sup>Department of Political Science, University of California, Berkeley, CA 94720; and <sup>d</sup>Department of Electrical Engineering and Computer Science, University of California, Berkeley, CA 94720

Edited by Richard M. Shiffrin, Indiana University, Bloomington, IN, and approved December 1, 2015 (received for review June 3, 2015)

**We provide a principled way for investigators to analyze randomized experiments when the number of covariates is large. Investigators often use linear multivariate regression to analyze randomized experiments instead of simply reporting the difference of means between treatment and control groups. Their aim is to reduce the variance of the estimated treatment effect by adjusting for covariates. If there are a large number of covariates relative to the number of observations, regression may perform poorly because of overfitting. In such cases, the least absolute shrinkage and selection operator (Lasso) may be helpful. We study the resulting Lasso-based treatment effect estimator under the Neyman–Rubin model of randomized experiments. We present theoretical conditions that guarantee that the estimator is more efficient than the simple difference-of-means estimator, and we provide a conservative estimator of the asymptotic variance, which can yield tighter confidence intervals than the difference-of-means estimator. Simulation and data examples show that Lasso-based adjustment can be advantageous even when the number of covariates is less than the number of observations. Specifically, a variant using Lasso for selection and ordinary least squares (OLS) for estimation performs particularly well, and it chooses a smoothing parameter based on combined performance of Lasso and OLS.**

randomized experiment | Neyman–Rubin model | average treatment effect | high-dimensional statistics | Lasso

Randomized experiments are widely used to measure the efficacy of treatments. Randomization ensures that treatment assignment is not influenced by any potential confounding factors, both observed and unobserved. Experiments are particularly useful when there is no rigorous theory of a system's dynamics, and full identification of confounders would be impossible. This advantage was cast elegantly in mathematical terms in the early 20th century by Jerzy Neyman, who introduced a simple model for randomized experiments, which showed that the difference of average outcomes in the treatment and control groups is statistically unbiased for the average treatment effect (ATE) over the experimental sample (1).

However, no experiment occurs in a vacuum of scientific knowledge. Often, baseline covariate information is collected about individuals in an experiment. Even when treatment assignment is not related to these covariates, analyses of experimental outcomes often take them into account with the goal of improving the accuracy of treatment effect estimates. In modern randomized experiments, the number of covariates can be very large—sometimes even larger than the number of individuals in the study. In clinical trials overseen by regulatory bodies like the Food and Drug Administration and the Medicines and Healthcare products Regulatory Agency, demographic and genetic information may be recorded about each patient. In applications in the tech industry, where randomization is often called A/B testing, there is often a huge amount of behavioral data collected on each user. However, in this “big data” setting, much of these data may be irrelevant to the outcome being studied or there may be more potential covariates than observations, especially once interactions are taken into account. In these cases, selection

of important covariates or some form of regularization is necessary for effective regression adjustment.

To ground our discussion, we examine a randomized trial of the pulmonary artery catheter (PAC) that was carried out in 65 intensive care units in the United Kingdom between 2001 and 2004, called PAC-man (2). The PAC is a monitoring device commonly inserted into critically ill patients after admission to intensive care, and it provides a continuous measurement of several indicators of cardiac activity. However, insertion of PAC is an invasive procedure that carries some risk of complications (including death), and it involves significant expenditure both in equipment costs and personnel (3). Controversy over its use came to a head when an observational study found that PAC had an adverse effect on patient survival and led to increased cost of care (4). This led to several large-scale randomized trials, including PAC-man.

In the PAC-man trial, randomization of treatment was largely successful, and a number of covariates were measured about each patient in the study. If covariate interactions are included, the number of covariates exceeds the number of individuals in the study; however, few of them are predictive of the patient's outcome. As it turned out, the (pretreatment) estimated probability of death was imbalanced between the treatment and control groups ( $P = 0.005$ , Wilcoxon rank sum test). Because the control group had, on average, a slightly higher risk of death, the unadjusted difference-in-means estimator may overestimate the benefits of receiving a PAC. Adjustment for this imbalance seems advantageous in this case, because the pretreatment probability of death is clearly predictive of health outcomes posttreatment.

In this paper, we study regression-based adjustment, using the least absolute shrinkage and selection operator (Lasso) to select relevant covariates. Standard linear regression based on ordinary least squares (OLS) suffers from overfitting if a large number of covariates and interaction terms are included in the model. In such cases, researchers sometimes perform model selection based on observing which covariates are unbalanced given the realized randomization. This generally leads to misleading inferences because of incorrect test levels (5). The Lasso (6) provides researchers with an alternative that can mitigate these problems and still perform model selection. We define an estimator,  $\widehat{ATE}_{\text{Lasso}}$ , which is based on running an  $l_1$ -penalized linear

This paper results from the Arthur M. Sackler Colloquium of the National Academy of Sciences, “Drawing Causal Inference from Big Data,” held March 26–27, 2015, at the National Academies of Sciences in Washington, DC. The complete program and video recordings of most presentations are available on the NAS website at [www.nasonline.org/Big-data](http://www.nasonline.org/Big-data).

Author contributions: A.B., H.L., C.-H.Z., J.S.S., and B.Y. designed research, performed research, analyzed data, and wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

<sup>1</sup>A.B. and H.L. contributed equally to this work.

<sup>2</sup>To whom correspondence should be addressed. Email: [binyu@stat.berkeley.edu](mailto:binyu@stat.berkeley.edu).

This article contains supporting information online at [www.pnas.org/lookup/suppl/doi:10.1073/pnas.1510506113/-DCSupplemental](http://www.pnas.org/lookup/suppl/doi:10.1073/pnas.1510506113/-DCSupplemental).

regression of the outcome on treatment, covariates, and, following the method introduced in ref. 7, treatment by covariate interactions. Because of the geometry of the  $l_1$  penalty, the Lasso will usually set many regression coefficients to 0, and is well defined even if the number of covariates is larger than the number of observations. The Lasso's theoretical properties under the standard linear model have been widely studied in the last decade; consistency properties for coefficient estimation, model selection, and out-of-sample prediction are well understood (see ref. 8 for an overview).

In the theoretical analysis in this paper, instead of assuming that the standard linear model is the true data-generating mechanism, we work under the aforementioned nonparametric model of randomization introduced by Neyman (1) and popularized by Donald Rubin (9). In this model, the outcomes and covariates are fixed quantities, and the treatment group is assumed to be sampled without replacement from a finite population. The treatment indicator, rather than an error term, is the source of randomness, and it determines which of two potential outcomes is revealed to the experimenter. Unlike the standard linear model, the Neyman–Rubin model makes few assumptions not guaranteed by the randomization itself. The setup of the model does rely on the stable unit treatment value assumption, which states that there is only one version of treatment, and that the potential outcome of one unit should be unaffected by the particular assignment of treatments to the other units; however, it makes no assumptions of linearity or exogeneity of error terms. OLS (7, 10, 11), logistic regression (12), and poststratification (13) are among the adjustment methods that have been studied under this model.

To be useful to practitioners, the Lasso-based treatment effect estimator must be consistent and yield a method to construct valid confidence intervals. We outline conditions on the covariates and potential outcomes that will guarantee these properties. We show that an upper bound for the asymptotic variance can be estimated from the model residuals, yielding asymptotically conservative confidence intervals for the ATE, which can be substantially narrower than the unadjusted confidence intervals. Simulation studies are provided to show the advantage of the Lasso-adjusted estimator and to show situations where it breaks down. We apply the estimator to the PAC-man data, and compare the estimates and confidence intervals derived from the unadjusted, OLS-adjusted, and Lasso-adjusted methods. We also compare different methods of selecting the Lasso tuning parameter on these data.

## Framework and Definitions

We give a brief outline of the Neyman–Rubin model for a randomized experiment; the reader is urged to consult refs. 1, 9, and 14 for more details. We follow the notation introduced in refs. 7 and 10. For concreteness, we illustrate the model in the context of the PAC-man trial.

For each individual in the study, the model assumes that there exists a pair of quantities representing his/her health outcomes under the possibilities of receiving and not receiving the catheter. These are called the potential outcomes under treatment and control, and are denoted as  $a_i$  and  $b_i$ , respectively. In the course of the study, the experimenter observes only one of these quantities for each individual, because the catheter is either inserted or not. The causal effect of the treatment on individual  $i$  is defined, in theory, to be  $a_i - b_i$ , but this is unobservable. Instead of trying to infer individual-level effects, we will assume that the intention is to estimate the average causal effect over the whole population, as outlined in the next section.

In the mathematical specification of this model, we consider the potential outcomes to be fixed, nonrandom quantities, even though they are not all observable. The only randomness in the model comes from the assignment of treatment, which is

controlled by the experimenter. We define random treatment indicators  $T_i$ , which take on a value 1 for a treated individual, or 0 for an untreated individual. We will assume that the set of treated individuals is sampled without replacement from the full population, where the size of the treatment group is fixed beforehand; thus, the  $T_i$  are identically distributed but not independent. The model for the observed outcome for individual  $i$ , defined as  $Y_i$ , is thus as follows:

$$Y_i = T_i a_i + (1 - T_i) b_i.$$

This equation simply formalizes the idea that the experimenter observes the potential outcome under treatment for those who receive the treatment, and the potential outcome under control for those who do not.

Note that the model does not incorporate any covariate information about the individuals in the study, such as physiological characteristics or health history. However, we will assume we have measured a vector of baseline, preexperimental covariates for each individual  $i$ . These might include, for example, age, gender, and genetic makeup. We denote the covariates for individual  $i$  as the column vector  $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})^T \in \mathbb{R}^p$  and the full design matrix of the experiment as  $X = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T$ . In *Theoretical Results*, we will assume that there is a correlational relationship between an individual's potential outcomes and covariates, but we will not assume a generative statistical model.

Define the set of treated individuals as  $A = \{i \in \{1, \dots, n\} : T_i = 1\}$ , and similarly define the set of control individuals as  $B$ . Define the number of treated and control individuals as  $n_A = |A|$  and  $n_B = |B|$ , respectively, so that  $n_A + n_B = n$ . We add a line on top of a quantity to indicate its average and a subscript  $A$  or  $B$  to label the treatment or control group. Thus, for example, the average values of the potential outcomes and the covariates in the treatment group are as follows:

$$\bar{a}_A = n_A^{-1} \sum_{i \in A} a_i, \quad \bar{\mathbf{x}}_A = n_A^{-1} \sum_{i \in A} \mathbf{x}_i,$$

respectively. Note that these are random quantities in this model, because the set  $A$  is determined by the random treatment assignment. Averages over the whole population are denoted as

$$\bar{a} = n^{-1} \sum_{i=1}^n a_i, \quad \bar{b} = n^{-1} \sum_{i=1}^n b_i, \quad \bar{\mathbf{x}} = n^{-1} \sum_{i=1}^n \mathbf{x}_i.$$

Note that the averages of potential outcomes over the whole population are not considered random, but are unobservable.

## Treatment Effect Estimation

Our main inferential goal will be average effect of the treatment over the whole population in the study. In a trial such as PAC-man, this represents the difference between the average outcome if everyone had received the catheter, and the average outcome if no one had received it. This is defined as follows:

$$ATE = \bar{a} - \bar{b}.$$

The most natural estimator arises by replacing the population averages with the sample averages:

$$\widehat{ATE}_{\text{unadj}} = \bar{a}_A - \bar{b}_B,$$

The subscript “unadj” indicates an estimator without regression adjustment. The foundational work in ref. 1 points out that, under a randomized assignment of treatment,  $\widehat{ATE}_{\text{unadj}}$  is unbiased for  $ATE$ , and derives a conservative procedure for estimating its variance.

Although  $\widehat{ATE}_{\text{unadj}}$  is an attractive estimator, covariate information can be used to make adjustments in the hope of reducing variance. A commonly used estimator is as follows:

$$\widehat{ATE}_{\text{adj}} = \left[ \bar{a}_A - (\bar{\mathbf{x}}_A - \bar{\mathbf{x}})^T \hat{\boldsymbol{\beta}}^{(a)} \right] - \left[ \bar{b}_B - (\bar{\mathbf{x}}_B - \bar{\mathbf{x}})^T \hat{\boldsymbol{\beta}}^{(b)} \right],$$

where  $\hat{\boldsymbol{\beta}}^{(a)}, \hat{\boldsymbol{\beta}}^{(b)} \in \mathbb{R}^p$  are adjustment vectors for the treatment and control groups, respectively, as indicated by the superscripts. The terms  $\bar{\mathbf{x}}_A - \bar{\mathbf{x}}$  and  $\bar{\mathbf{x}}_B - \bar{\mathbf{x}}$  represent the fluctuation of the covariates in the subsample relative to the full sample, and the adjustment vectors fit the linear relationships between the covariates and potential outcomes under treatment and control. For example, in the PAC-man trial, this would help alleviate the imbalance in the pretreatment estimated probability of death: the corresponding element of  $\bar{\mathbf{x}}_B - \bar{\mathbf{x}}$  would be positive (due to the higher average probability of death in the control group), the corresponding element of  $\hat{\boldsymbol{\beta}}^{(b)}$  would be negative (a higher probability of death correlates with worse health outcomes), so the overall treatment effect estimate would be adjusted downward. This procedure is equivalent to imputing the unobserved potential outcomes; if we define

$$\hat{a}_B = \bar{a}_A + (\bar{\mathbf{x}}_B - \bar{\mathbf{x}})^T \hat{\boldsymbol{\beta}}^{(a)}, \quad \hat{b}_A = \bar{b}_B + (\bar{\mathbf{x}}_A - \bar{\mathbf{x}})^T \hat{\boldsymbol{\beta}}^{(b)},$$

we can form the equivalent estimator:

$$\widehat{ATE}_{\text{adj}} = n^{-1} (n_A \hat{a}_B + n_B \hat{b}_A) - n^{-1} (n_B \bar{b}_B + n_A \bar{b}_A).$$

If we consider these adjustment vectors to be fixed (nonrandom), or if they are derived from an independent data source, then this estimator is still unbiased, and may have substantially smaller asymptotic and finite-sample variance than the unadjusted estimator. This allows for construction of tighter confidence intervals for the true treatment effect.

In practice, the “ideal” linear adjustment vectors, leading to a minimum-variance estimator of the form of  $\widehat{ATE}_{\text{adj}}$ , cannot be computed from the observed data. However, they can be estimated, possibly at the expense of introducing modest finite-sample bias into the treatment effect estimate. In the classical setup, when the number of covariates is relatively small, OLS regression can be used. The asymptotic properties of this kind of estimator are explored under the Neyman–Rubin model in refs. 7, 11, and 12. We will follow a particular scheme that is studied in ref. 7 and shown to have favorable properties: we regress the outcome on treatment indicators, covariates, and treatment by covariate interactions. This is equivalent to running separate regressions in the treatment and control groups of outcome against an intercept and covariates. If we define  $\hat{\boldsymbol{\beta}}_{\text{OLS}}^{(a)}$  and  $\hat{\boldsymbol{\beta}}_{\text{OLS}}^{(b)}$  as the coefficients from the separate regressions, then the estimator is as follows:

$$\widehat{ATE}_{\text{OLS}} = \left[ \bar{a}_A - (\bar{\mathbf{x}}_A - \bar{\mathbf{x}})^T \hat{\boldsymbol{\beta}}_{\text{OLS}}^{(a)} \right] - \left[ \bar{b}_B - (\bar{\mathbf{x}}_B - \bar{\mathbf{x}})^T \hat{\boldsymbol{\beta}}_{\text{OLS}}^{(b)} \right].$$

This has some finite-sample bias, but ref. 7 shows that it vanishes quickly at the rate of  $1/n$  under moment conditions on the potential outcomes and covariates. Moreover, for a fixed  $p$ , under regularity conditions, the inclusion of interaction terms guarantees that it never has higher asymptotic variance than the unadjusted estimator, and asymptotically conservative confidence intervals for the true parameter can be constructed.

In modern randomized trials, where a large number of covariates are recorded for each individual,  $p$  may be comparable to or even larger than  $n$ . In this case, OLS regression can overfit the data badly, or may even be ill posed, leading to estimators with large finite-sample variance. To remedy this, we propose estimating

the adjustment vectors using the Lasso (6). The adjustment vectors would take the following form:

$$\hat{\boldsymbol{\beta}}_{\text{Lasso}}^{(a)} = \arg \min_{\boldsymbol{\beta}} \left[ \frac{1}{2n_A} \sum_{i \in A} \left( a_i - \bar{a}_A - (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T \boldsymbol{\beta} \right)^2 + \lambda_a \sum_{j=1}^p |\beta_j| \right], \quad [1]$$

$$\hat{\boldsymbol{\beta}}_{\text{Lasso}}^{(b)} = \arg \min_{\boldsymbol{\beta}} \left[ \frac{1}{2n_B} \sum_{i \in B} \left( b_i - \bar{b}_B - (\mathbf{x}_i - \bar{\mathbf{x}}_B)^T \boldsymbol{\beta} \right)^2 + \lambda_b \sum_{j=1}^p |\beta_j| \right], \quad [2]$$

and the proposed Lasso-adjusted ATE estimator is as follows:

$$\widehat{ATE}_{\text{Lasso}} = \left[ \bar{a}_A - (\bar{\mathbf{x}}_A - \bar{\mathbf{x}})^T \hat{\boldsymbol{\beta}}_{\text{Lasso}}^{(a)} \right] - \left[ \bar{b}_B - (\bar{\mathbf{x}}_B - \bar{\mathbf{x}})^T \hat{\boldsymbol{\beta}}_{\text{Lasso}}^{(b)} \right].$$

[To simplify the notation, we omit the dependence of  $\hat{\boldsymbol{\beta}}_{\text{Lasso}}^{(a)}, \hat{\boldsymbol{\beta}}_{\text{Lasso}}^{(b)}$ ,  $\lambda_a$ , and  $\lambda_b$  on the population size  $n$ .] Here,  $\lambda_a$  and  $\lambda_b$  are regularization parameters for the Lasso, which must be chosen by the experimenter; simulations show that cross-validation (CV) works well. In the next section, we study this estimator under the Neyman–Rubin model and provide conditions on the potential outcomes, the covariates, and the regularization parameters under which  $\widehat{ATE}_{\text{Lasso}}$  enjoys similar asymptotic and finite-sample advantages as  $\widehat{ATE}_{\text{OLS}}$ .

It is worth noting that, when two different adjustments are made for the treatment and control groups as in ref. 7 and here, the covariates do not have to be the same for the two groups. However, when they are not the same, the Lasso- or OLS-adjusted estimators are no longer guaranteed to have smaller or equal asymptotic variance than the unadjusted one, even in the case of fixed  $p$ . In practice, one may still choose between the adjusted and unadjusted estimators based on the widths of the corresponding confidence intervals.

## Theoretical Results

**Notation.** For a vector  $\boldsymbol{\beta} \in \mathbb{R}^p$  and a subset  $S \subset \{1, \dots, p\}$ , let  $\beta_j$  be the  $j$ th component of  $\boldsymbol{\beta}$ ,  $\boldsymbol{\beta}_S = (\beta_j : j \in S)^T$ ,  $S^c$  be the complement of  $S$ , and  $|S|$  the cardinality of the set  $S$ . For any column vector  $\mathbf{u} = (u_1, \dots, u_m)^T$ , let  $\|\mathbf{u}\|_2^2 = \sum_{i=1}^m u_i^2$ ,  $\|\mathbf{u}\|_1 = \sum_{i=1}^m |u_i|$ ,  $\|\mathbf{u}\|_\infty = \max_{i=1, \dots, m} |u_i|$ , and  $\|\mathbf{u}\|_0 = |\{j : u_j \neq 0\}|$ . For a given  $m \times m$  matrix  $D$ , let  $\lambda_{\min}(D)$  and  $\lambda_{\max}(D)$  be the smallest and largest eigenvalues of  $D$ , respectively, and  $D^{-1}$ , the inverse of the matrix  $D$ . Let  $\xrightarrow{d}$  and  $\xrightarrow{p}$  denote convergence in distribution and in probability, respectively.

**Decomposition of the Potential Outcomes.** The Neyman–Rubin model does not assume a linear relationship between the potential outcomes and the covariates. To study the properties of adjustment under this model, we decompose the potential outcomes into a term linear in the covariates and an error term. Given vectors of coefficients  $\boldsymbol{\beta}^{(a)}, \boldsymbol{\beta}^{(b)} \in \mathbb{R}^p$ , we write for  $i = 1, \dots, n$ ,

$$a_i = \bar{a} + (\mathbf{x}_i - \bar{\mathbf{x}})^T \boldsymbol{\beta}^{(a)} + e_i^{(a)}, \quad [3]$$

$$b_i = \bar{b} + (\mathbf{x}_i - \bar{\mathbf{x}})^T \boldsymbol{\beta}^{(b)} + e_i^{(b)}. \quad [4]$$

[Again, we omit the dependence of  $\boldsymbol{\beta}^{(a)}, \boldsymbol{\beta}^{(b)}$ ,  $\lambda_a$ ,  $\lambda_b$ ,  $e_i^{(a)}$ , and  $e_i^{(b)}$  on  $n$ .]

Note that we have not added any assumptions to the model; we have simply defined unit-level residuals,  $e_i^{(a)}$  and  $e_i^{(b)}$ , given the vectors  $\boldsymbol{\beta}^{(a)}, \boldsymbol{\beta}^{(b)}$ . All of the quantities in [3] and [4] are fixed, deterministic numbers. It is easy to verify that  $\bar{e}^{(a)} = \bar{e}^{(b)} = 0$ . To pursue a theory for the Lasso, we will add assumptions on the populations of  $a_i$ 's,  $b_i$ 's, and  $\mathbf{x}_i$ 's, and we will assume the

existence of  $\beta^{(a)}, \beta^{(b)}$  such that the error terms satisfy certain assumptions.

**Conditions.** We will need the following to hold for both the treatment and control potential outcomes. The first set of assumptions (1–3) are similar to those found in ref. 7.

**Condition 1:** Stability of treatment assignment probability.

$$n_A/n \rightarrow p_A, \text{ as } n \rightarrow \infty, \quad [5]$$

for some  $p_A \in (0, 1)$ .

**Condition 2:** The centered moment conditions. There exists a fixed constant  $L > 0$  such that, for all  $n = 1, 2, \dots$  and  $j = 1, \dots, p$ ,

$$n^{-1} \sum_{i=1}^n (x_{ij} - (\bar{x}_j))^4 \leq L; \quad [6]$$

$$n^{-1} \sum_{i=1}^n (e_i^{(a)})^4 \leq L; \quad n^{-1} \sum_{i=1}^n (e_i^{(b)})^4 \leq L. \quad [7]$$

**Condition 3:** The means  $n^{-1} \sum_{i=1}^n (e_i^{(a)})^2$ ,  $n^{-1} \sum_{i=1}^n (e_i^{(b)})^2$ , and  $n^{-1} \sum_{i=1}^n e_i^{(a)} e_i^{(b)}$  converge to finite limits.

Because we consider the high-dimensional setting where  $p$  is allowed to be much larger than  $n$ , we need additional assumptions to ensure that the Lasso is consistent for estimating  $\beta^{(a)}$  and  $\beta^{(b)}$ . Before stating them, we define several quantities.

**Definition 1:** Given  $\beta^{(a)}$  and  $\beta^{(b)}$ , the sparsity measures for treatment and control groups,  $s^{(a)}$  and  $s^{(b)}$ , are defined as the number of nonzero elements of  $\beta^{(a)}$  and  $\beta^{(b)}$ , i.e.,

$$s^{(a)} = \left| \left\{ j : \beta_j^{(a)} \neq 0 \right\} \right|, \quad s^{(b)} = \left| \left\{ j : \beta_j^{(b)} \neq 0 \right\} \right|, \quad [8]$$

respectively. We will allow  $s^{(a)}$  and  $s^{(b)}$  to grow with  $n$ , although the notation does not explicitly show this.

**Definition 2:** Define  $\delta_n$  to be the maximum covariance between the error terms and the covariates.

$$\delta_n = \max_{\omega=a,b} \left\{ \max_j \left| \frac{1}{n} \sum_{i=1}^n (x_{ij} - (\bar{x}_j)) (e_i^{(\omega)} - \bar{e}^{(\omega)}) \right| \right\}. \quad [9]$$

The following conditions will guarantee that the Lasso consistently estimates the adjustment vectors  $\beta^{(a)}, \beta^{(b)}$  at a fast enough rate to ensure asymptotic normality of  $\widehat{ATE}_{\text{Lasso}}$ . It is an open question whether a weaker form of consistency would be sufficient for our results to hold.

**Condition 4:** Decay and scaling. Let  $s = \max\{s^{(a)}, s^{(b)}\}$ .

$$\delta_n = o\left(\frac{1}{s\sqrt{\log p}}\right). \quad [10]$$

$$(s \log p) / \sqrt{n} = o(1). \quad [11]$$

**Condition 5:** Cone invertibility factor. Define the Gram matrix as  $\Sigma = n^{-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$ . There exist constants  $C > 0$  and  $\xi > 1$  not depending on  $n$ , such that

$$\|\mathbf{h}_S\|_1 \leq Cs \|\Sigma \mathbf{h}\|_\infty, \quad \forall \mathbf{h} \in \mathcal{C}, \quad [12]$$

with  $\mathcal{C} = \{\mathbf{h} : \|\mathbf{h}_{S^c}\|_1 \leq \xi \|\mathbf{h}_S\|_1\}$ , and

$$S = \left\{ j : \beta_j^{(a)} \neq 0 \text{ or } \beta_j^{(b)} \neq 0 \right\}. \quad [13]$$

**Condition 6:** Let  $\tau = \min\{1/70, (3p_A)^2/70, (3-3p_A)^2/70\}$ . For constants  $0 < \eta < (\xi - 1)/(\xi + 1)$  and  $1/\eta < M < \infty$ , assume the

regularization parameters of the Lasso belong to the following sets:

$$\lambda_a \in \left(\frac{1}{\eta}, M\right] \times \left(\frac{2(1+\tau)L^{1/2}}{p_A} \sqrt{\frac{2 \log p}{n}} + \delta_n\right), \quad [14]$$

$$\lambda_b \in \left(\frac{1}{\eta}, M\right] \times \left(\frac{2(1+\tau)L^{1/2}}{p_B} \sqrt{\frac{2 \log p}{n}} + \delta_n\right). \quad [15]$$

Denote, respectively, the population variances of  $e^{(a)}$  and  $e^{(b)}$  and the population covariance between them by the following:

$$\sigma_{e^{(a)}}^2 = n^{-1} \sum_{i=1}^n (e_i^{(a)})^2, \quad \sigma_{e^{(b)}}^2 = n^{-1} \sum_{i=1}^n (e_i^{(b)})^2, \\ \sigma_{e^{(a)}e^{(b)}} = n^{-1} \sum_{i=1}^n e_i^{(a)} e_i^{(b)}.$$

**Theorem 1.** Assume Conditions 1–6 hold for some  $\beta^{(a)}$  and  $\beta^{(b)}$ . Then,

$$\sqrt{n}(\widehat{ATE}_{\text{Lasso}} - ATE) \xrightarrow{d} \mathcal{N}(0, \sigma^2), \quad [16]$$

where

$$\sigma^2 = \lim_{n \rightarrow \infty} \left[ \frac{1-p_A}{p_A} \sigma_{e^{(a)}}^2 + \frac{p_A}{1-p_A} \sigma_{e^{(b)}}^2 + 2\sigma_{e^{(a)}e^{(b)}} \right]. \quad [17]$$

The proof of *Theorem 1* is given in *SI Appendix*. It is easy to show, as in the following corollary of *Theorem 1*, that the asymptotic variance of  $\widehat{ATE}_{\text{Lasso}}$  is no worse than  $\widehat{ATE}_{\text{unadj}}$  when  $\beta^{(a)}$  and  $\beta^{(b)}$  are defined as coefficients of regressing potential outcomes on a subset of covariates. More specifically, suppose there exists a subset  $J \subset \{1, \dots, p\}$ , such that

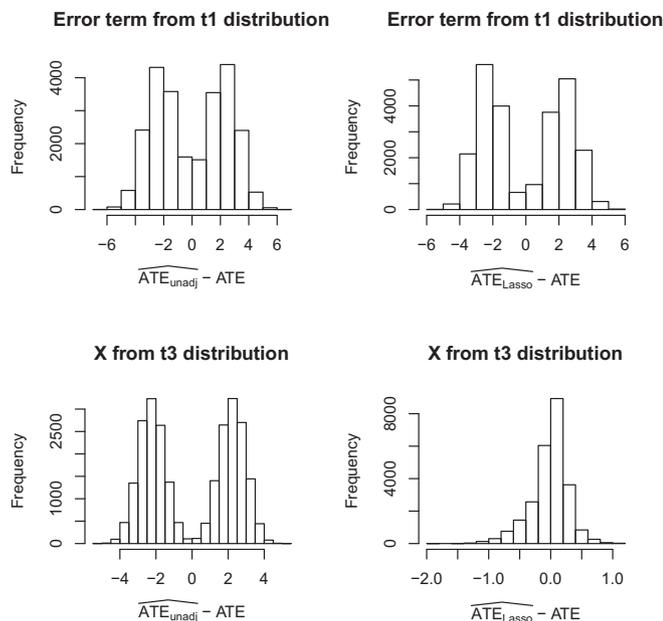
$$\beta^{(a)} = \left( (\beta_J^{(a)})^T, \mathbf{0} \right)^T, \quad \beta^{(b)} = \left( (\beta_J^{(b)})^T, \mathbf{0} \right)^T, \quad [18]$$

where  $\beta_J^{(a)}$  and  $\beta_J^{(b)}$  are the population-level OLS coefficients for regressing the potential outcomes  $a$  and  $b$  on the covariates in the subset  $J$  with intercept, respectively.

**Corollary 1.** For  $\beta^{(a)}$  and  $\beta^{(b)}$  defined in [18] and some  $\lambda_a$  and  $\lambda_b$ , assume Conditions 1–6 hold. Then the asymptotic variance of  $\sqrt{n} \widehat{ATE}_{\text{Lasso}}$  is no greater than that of the  $\sqrt{n} \widehat{ATE}_{\text{unadj}}$ . The difference is  $(1/(p_A(1-p_A)))\Delta$ , where

$$\Delta = - \lim_{n \rightarrow \infty} \|X\beta_E\|_2^2 \leq 0, \quad \beta_E = (1-p_A)\beta^{(a)} + p_A\beta^{(b)}. \quad [19]$$

**Remark 1:** If, instead of *Condition 6*, we assume that the covariates are uniformly bounded, i.e.,  $\max_{i,j} |x_{ij}| \leq L$ , then the fourth moment condition on the error terms, given in [7], can be weakened to a second moment condition. Although we do not prove the necessity of any of our conditions, our simulation studies show that the distributions of the unadjusted and the Lasso-adjusted estimator may be nonnormal when (i) the covariates are generated from Gaussian distributions and the error terms do not satisfy second moment condition, e.g., being generated from a  $t$  distribution with one degree of freedom; or (ii) the covariates do not have bounded fourth moments, e.g., being generated from a  $t$  distribution with three degrees of freedom. See the histograms in Fig. 1, where the corresponding  $p$  values of Kolmogorov–Smirnov testing for normality are less than  $2.2e - 16$ . These findings indicate that our moment conditions cannot be dramatically weakened



**Fig. 1.** Histograms of the unadjusted estimator and the Lasso-adjusted estimator when the moment conditions do not hold. We select the tuning parameters for Lasso using 10-fold CV. The potential outcomes are simulated from linear regression model and then kept fixed (see more details in [SI Appendix](#)). For the upper two subplots, the error terms are generated from  $t$  distribution with one degree of freedom and therefore do not satisfy second moment condition; whereas for the lower two subplots, the covariates are generated from  $t$  distribution with three degrees of freedom and thus violate fourth moment condition.

for asymptotic normality. However, we also find that the Lasso-adjusted estimator still has smaller variance and mean squared error than the unadjusted estimator, even when these moment conditions do not hold. In practice, when the covariates do not have bounded fourth moments, one may perform some transformation—e.g., a logarithm transformation—to ensure that the transformed covariates have bounded fourth moments while having a sufficiently large variance so as to retain useful information. We leave it as future work to explore the properties of different transformations.

**Remark 2:** Statement [11], typically required in debiasing the Lasso (15), is stronger by a factor of  $\sqrt{\log p}$  than the usual requirement for  $l_1$  consistency of the Lasso.

**Remark 3:** *Condition 5* is slightly weaker than the typical restricted eigenvalue condition for analyzing the Lasso.

**Remark 4:** If we assume  $\delta_n = O(\sqrt{\log p/n})$ , which satisfies [10], then *Condition 6* requires that the tuning parameters are proportional to  $\sqrt{\log p/n}$ , which is typically assumed for the Lasso in the high-dimensional linear regression model.

**Remark 5:** For fixed  $p$ ,  $\delta_n = 0$  in [9], *Condition 4* holds automatically, and *Condition 5* holds when the smallest eigenvalue of  $\Sigma$  is uniformly bounded away from 0. In this case, *Corollary 1* reverts to corollary 1.1. in ref. 7. When these conditions are not satisfied, we should set  $\lambda_a$  and  $\lambda_b$  to be large enough to cause the Lasso-adjusted estimator to revert to the unadjusted one.

**Neyman-Type Conservative Variance Estimate.** We note that the asymptotic variance in *Theorem 1* involves the cross-product term  $\sigma_{e^{(a)}e^{(b)}}$ , which is not consistently estimable in the Neyman–Rubin model as  $a_i$  and  $b_i$  are never simultaneously observed. However, we can give a Neyman-type conservative estimate of the variance. Let

$$\hat{\sigma}_{e^{(a)}}^2 = \frac{1}{n_A - df^{(a)}} \sum_{i \in A} \left( a_i - \bar{a}_A - (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T \hat{\boldsymbol{\beta}}_{\text{Lasso}}^{(a)} \right)^2,$$

$$\hat{\sigma}_{e^{(b)}}^2 = \frac{1}{n_B - df^{(b)}} \sum_{i \in B} \left( b_i - \bar{b}_B - (\mathbf{x}_i - \bar{\mathbf{x}}_B)^T \hat{\boldsymbol{\beta}}_{\text{Lasso}}^{(b)} \right)^2,$$

where  $df^{(a)}$  and  $df^{(b)}$  are degrees of freedom defined by the following:

$$df^{(a)} = \hat{s}^{(a)} + 1 = \left\| \hat{\boldsymbol{\beta}}_{\text{Lasso}}^{(a)} \right\|_0 + 1; \quad df^{(b)} = \hat{s}^{(b)} + 1 = \left\| \hat{\boldsymbol{\beta}}_{\text{Lasso}}^{(b)} \right\|_0 + 1.$$

Define the variance estimate of  $\sqrt{n}(\widehat{ATE}_{\text{Lasso}} - ATE)$  as follows:

$$\hat{\sigma}_{\text{Lasso}}^2 = \frac{n}{n_A} \hat{\sigma}_{e^{(a)}}^2 + \frac{n}{n_B} \hat{\sigma}_{e^{(b)}}^2. \quad [20]$$

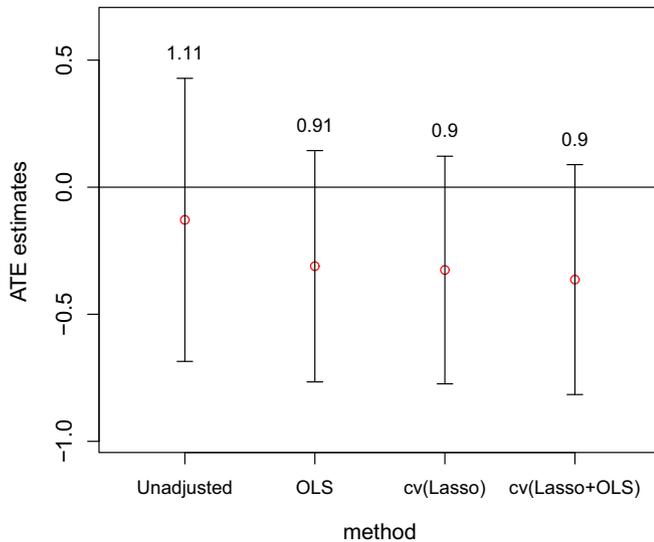
We will show in [SI Appendix, Theorem S1](#), that the limit of  $\hat{\sigma}_{\text{Lasso}}^2$  is greater than or equal to the asymptotic variance of  $\sqrt{n}(\widehat{ATE}_{\text{Lasso}} - ATE)$ , and therefore can be used to construct a conservative confidence interval for the ATE.

**Related Work.** The Lasso has already made several appearances in the literature on treatment effect estimation. In the context of observational studies, ref. 15 constructs confidence intervals for preconceived effects or their contrasts by debiasing the Lasso-adjusted regression, ref. 16 employs the Lasso as a formal method for selecting adjustment variables via a two-stage procedure that concatenates features from models for treatment and outcome, and similarly, ref. 17 gives very general results for estimating a wide range of treatment effect parameters, including the case of instrumental variables estimation. In addition to the Lasso, ref. 18 considers nonparametric adjustments in the estimation of ATE. In works such as these, which deal with observational studies, confounding is the major issue. With confounding, the naive difference-in-means estimator is biased for the true treatment effect, and adjustment is used to form an unbiased estimator. However, in our work, which focuses on a randomized trial, the difference-in-means estimator is already unbiased; adjustment reduces the variance while, in fact, introducing a small amount of finite-sample bias. Another major difference between this prior work and ours is the sampling framework: we operate within the Neyman–Rubin model with fixed potential outcomes for a finite population, where the treatment group is sampled without replacement, whereas these papers assume independent sampling from a probability distribution with random error terms.

Our work is related to the estimation of heterogeneous or subgroup-specific treatment effects, including interaction terms to allow the imputed individual-level treatment effects to vary according to some linear combination of covariates. This is pursued in the high-dimensional setting in ref. 19; this work advocates solving the Lasso on a reduced set of modified covariates, rather than the full set of covariate by treatment interactions, and includes extensions to binary outcomes and survival data. The recent work in ref. 20 considers the problem of designing multiple-testing procedures for detecting subgroup-specific treatment effects; they pose this as an optimization over testing procedures where constraints are added to enforce guarantees on type I error rate and power to detect effects. Again, the sampling framework in these works is distinct from ours; they do not use the Neyman–Rubin model as a basis for designing the methods or investigating their properties.

**PAC Data Illustration and Simulations.** We now return to the PAC-man study introduced earlier. We examine the data in more

## ATE estimates for PAC data



**Fig. 2.** ATE estimates (red circles) and 95% confidence intervals (bars) for the PAC data. The numbers above each bar are the corresponding interval lengths.

detail and explore the results of several adjustment procedures. There were 1,013 patients in the PAC-man study: 506 treated (managed with PAC) and 507 control (managed without PAC, but retaining the option of using alternative devices). The outcome variable is quality-adjusted life years (QALYs). One QALY represents 1 year of life in full health; in-hospital death corresponds to a QALY of zero. We have 59 covariates about each individual in the study; we include all main effects as well as 1,113 two-way interactions, and form a design matrix  $\mathbf{X}$  with 1,172 columns and 1,013 rows. See *SI Appendix* for more details on the design matrix.

The assumptions that underpin the theoretical guarantees of the  $ATE_{Lasso}$  estimator are, in practice, not explicitly checkable, but we attempt to inspect the quantities that are involved in the conditions to help readers make their own judgement. The uniform bounds on the fourth moments refer to a hypothetical sequence of populations; these cannot be verified given that the investigator has a single dataset. However, as an approximation, the fourth moments of the data can be inspected to ensure that they are not too large. In this dataset, the maximum fourth moment of the covariates is 37.3, which is indicative of a heavy-tailed and potentially destabilizing covariate; however, it occurs in an interaction term not selected by the Lasso, and thus does

not influence the estimate. [The fourth moments of the covariates are shown in *SI Appendix*, Fig. S9. The covariates with the largest two fourth moments (37.3 and 34.9, respectively) are quadratic term  $interactnew^2$  and interaction term  $IMscorecrct : systemnew$ . Neither of them is selected by the Lasso to do the adjustment.] Checking the conditions for high-dimensional consistency of the Lasso would require knowledge of the unknown active set  $S$ , and moreover, even if it were known, calculating the cone invertibility factor would involve an infeasible optimization. This is a general issue in the theory of sparse linear high-dimensional estimation. To approximate these conditions, we use the bootstrap to estimate the active set of covariates  $S$  and the error terms  $e^{(a)}$  and  $e^{(b)}$ . See *SI Appendix* for more details. Our estimated  $S$  contains 16 covariates and the estimated second moments of  $e^{(a)}$  and  $e^{(b)}$  are 11.8 and 12.0, respectively. The estimated maximal covariance  $\delta_n$  equals 0.34 and the scaling  $(s \log p)/\sqrt{n}$  is 3.55. Although this is not close to zero, we should mention that the estimation of  $\delta_n$  and  $(s \log p)/\sqrt{n}$  can be unstable and less accurate because it is based on a subsample of the population. As an approximation to *Condition 5*, we examine the largest and smallest eigenvalues of the sub-Gram matrix  $(1/n)\mathbf{X}_S^T \mathbf{X}_S$ , which are 2.09 and 0.18, respectively. Thus, the quantity in *Condition 5* seems reasonably bounded away from zero.

We now estimate the ATE using the unadjusted estimator, the Lasso-adjusted estimator, and the OLS-adjusted estimator, which is computed based on a subdesign matrix containing only the 59 main effects. We also present results for the two-step estimator  $ATE_{Lasso+OLS}$ , which adopts the Lasso to select covariates and then uses OLS to refit the regression coefficients. In the next paragraph and in *SI Appendix, Algorithm 1*, we show how we adapt the CV procedure to select the tuning parameter for  $ATE_{Lasso+OLS}$  based on a combined performance of Lasso and OLS, or  $cv(Lasso+OLS)$ .

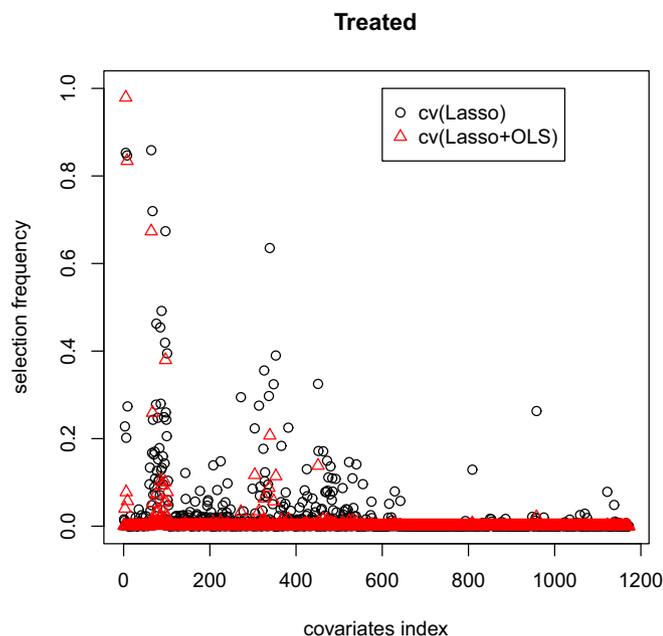
We use the R package “glmnet” to compute the Lasso solution path and select the tuning parameters  $\lambda_a$  and  $\lambda_b$  by 10-fold CV. To indicate the method of selecting tuning parameters, we denote the corresponding estimators as  $cv(Lasso)$  and  $cv(Lasso+OLS)$ , respectively. We should mention that for the  $cv(Lasso+OLS)$ -adjusted estimator, we compute the CV error for a given value of  $\lambda_a$  (or  $\lambda_b$ ) based on the whole Lasso+OLS procedure instead of just the Lasso estimator (*SI Appendix, Algorithm 1*). Therefore, the  $cv(Lasso+OLS)$  and the  $cv(Lasso)$  may select different covariates to do the adjustment. This type of CV requires more computation than the CV based on just the Lasso estimator because it needs to compute the OLS estimator for each fold and each given  $\lambda_a$  (or  $\lambda_b$ ), but it can give better prediction and model selection performance.

Fig. 2 presents the ATE estimates along with 95% confidence intervals. The interval lengths are shown on top of each interval bar. All of the methods give confidence intervals containing 0; hence, this experiment failed to provide sufficient evidence to

**Table 1. Selected covariates for adjustment**

Method	Treatment	Covariates
$cv(Lasso+OLS)$	Treated	Age, p_death, age · age, age:p_death
$cv(Lasso+OLS)$	Control	Age, p_death, age · age, age:p_death, p_death:mch_vent
$cv(Lasso)$	Treated	Pac_rate, age, p_death, age · age, p_death · p_death, region:im_score, region:systemnew, pac_rate:age, pac_rate:p_death, pac_rate:systemnew, im_score:interactnew, age:p_death, age:glasgow, age:systemnew, interactnew:systemnew, pac_rate:creatinine, age:mch_vent, age:respiratory, age:creatinine, interactnew:mch_vent, interactnew:male, glasgow:organ_failure, p_death:mch_vent, systemnew:male
$cv(Lasso)$	Control	Age, p_death, age · age, unitsize:p_death, pac_rate:systemnew, age:p_death, interactnew:mch_vent, p_death:mch_vent*

\*Covariate definitions: age, patient’s age; p\_death, baseline probability of death; mch\_vent, mechanical ventilation at admission; region, geographic region; pac\_rate, PAC rate in unit; creatinine, respiratory, glasgow, interactnew, organ\_failure, systemnew, and im\_score, various physiological indicators.



**Fig. 3.** Selection stability comparison of  $cv(Lasso)$  and  $cv(Lasso+OLS)$  for treatment group.

reject the hypothesis that PAC did not have an effect on patient QALYs (either positive or negative). Because the caretakers of patients managed without PAC retained the option of using less invasive cardiac output monitoring devices, such an effect may have been particularly hard to detect in this experiment.

However, it is interesting to note that, compared with the unadjusted estimator, the OLS-adjusted estimator causes the ATE estimate to decrease (from  $-0.13$  to  $-0.31$ ), and shortens the confidence interval by about 20%. This is due mainly to the imbalance in the pretreatment probability of death, which was highly predictive of the posttreatment QALYs. The  $cv(Lasso)$ -adjusted estimator yields a comparable ATE estimate and confidence interval, but the fitted model is more interpretable and parsimonious than the OLS model: it selects 24 and 8 covariates for treated and control, respectively. The  $cv(Lasso+OLS)$  estimator selects even fewer covariates: 4 and 5 for treated and control, respectively, but performs a similar adjustment as the  $cv(Lasso)$  (see the comparison of fitted values in *SI Appendix, Fig. S8*). We also note that these adjustments agree with the one performed in ref. 13, where the treatment effect was adjusted downward to  $-0.27$  after stratifying into four groups based on predicted probability of death.

The covariates selected by Lasso for adjustment are shown in Table 1, where “A·A” denotes quadratic term of the covariate A, and “A:B” denotes two-way interaction between two covariates A and B. Among them, patient’s age and estimated probability of death ( $p_{death}$ ), together with the quadratic term “age·age” and interactions “age: $p_{death}$ ” and “ $p_{death}$ :mech\_vent” (mechanical ventilation at admission), are the most important covariates for the adjustment. The patients in control group are slightly older and have slightly higher risk of death. These covariates are important predictors of the outcome. Therefore, the unadjusted estimator may overestimate the benefits of receiving PAC.

Because not all of the potential outcomes are observed, we cannot know the true gains of adjustment methods. However, we can estimate the gains via building a simulated set of potential outcomes by matching treated units to control units on observed covariates. We use the matching method described in ref. 21,

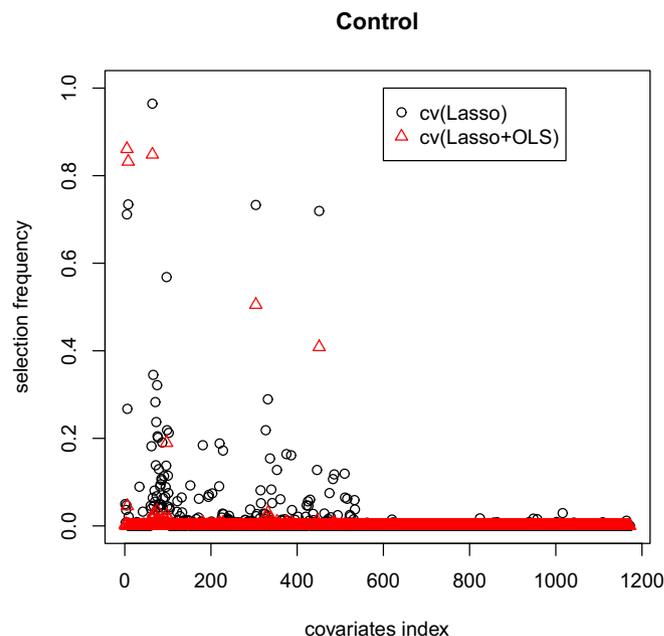
which gives 1,013 observations with all potential outcomes imputed. We match on the 59 main effects only. The ATE is  $-0.29$ . We then use this synthetic dataset to calculate the biases, SDs, and root-mean-square errors ( $\sqrt{MSE}$ ) of different ATE estimators based on 25,000 replicates of a completely randomized experiment, which assigns 506 subjects to the treated group and the remainders to the control group.

*SI Appendix, Table S5*, shows the results. For all of the methods, the bias is substantially smaller (by a factor of 100) than the SD. The SD and  $\sqrt{MSE}$  of the OLS-adjusted estimator are both 10.2% smaller than those of the unadjusted estimator, whereas the  $cv(Lasso)$ - and  $cv(Lasso+OLS)$ -adjusted estimators further improve the SD and  $\sqrt{MSE}$  of the OLS-adjusted estimator by  $\sim 4.7\%$ . Moreover, all these methods provide conservative confidence intervals with coverage probabilities higher than 99%. However, the interval lengths of the OLS-,  $cv(Lasso)$ -, and  $cv(Lasso+OLS)$ -adjusted estimator are comparable and are  $\sim 10\%$  shorter than that of the unadjusted estimator. The  $cv(Lasso+OLS)$ -adjusted estimator is similar to the  $cv(Lasso)$ -adjusted estimator in terms of mean squared error, confidence interval length, and coverage probability, but outperforms the latter with much fewer and more stable covariates in the adjustment (see Figs. 3 and 4 for the selection frequency of each covariate for treatment group and control group, respectively). We also show in *SI Appendix, Fig. S10*, that the sampling distribution of the estimates is very close to Normal.

We conduct additional simulation studies to evaluate the finite sample performance of  $\widehat{ATE}_{Lasso}$  and compare it with that of the OLS-adjusted estimator and the unadjusted estimator. A qualitative analysis of these simulations yields the same conclusions as presented above; however, for the sake of brevity, we defer the simulation details in *SI Appendix*.

## Discussion

We study the Lasso-adjusted ATE estimate under the Neyman–Rubin model for randomization. Our purpose in using the Neyman–Rubin model is to investigate the performance of the Lasso under a realistic sampling framework that does not impose strong assumptions on the data. We provide conditions that ensure



**Fig. 4.** Selection stability comparison of  $cv(Lasso)$  and  $cv(Lasso+OLS)$  for control group.

asymptotic normality, and provide a Neyman-type estimate of the asymptotic variance that can be used to construct a conservative confidence interval for the ATE. Although we do not require an explicit generative linear model to hold, our theoretical analysis requires the existence of latent “adjustment vectors” such that moment conditions of the error terms are satisfied, and that the cone invertibility condition of the sample covariance matrix is satisfied in addition to moment conditions for OLS adjustment as in ref. 7. Both assumptions are difficult to check in practice. In our theory, we do not address whether these assumptions are necessary for our results to hold, although simulations indicate that the moment conditions cannot be substantially weakened. As a by-product of our analysis, we extend Massart’s concentration inequality for sampling without replacement, which is useful for theoretical analysis under the Neyman–Rubin model. Simulation studies and the real-data illustration show the advantage of the Lasso-adjusted estimator in terms of estimation accuracy and model interpretation. In practice, we recommend a variant of Lasso,  $cv(\text{Lasso}+\text{OLS})$ , to select covariates and perform the adjustment, because it gives similar coverage probability and confidence interval length compared with  $cv(\text{Lasso})$ , but with far fewer covariates selected. In future work, we plan to extend our analysis to other popular methods in high-dimensional statistics such as Elastic-Net and ridge regression, which may be more appropriate for estimating adjusted ATE under different assumptions.

The main goal of using Lasso in this paper is to reduce the variance (and overall mean squared error) of ATE estimation. Another important task is to estimate heterogenous treatment effects and provide conditional treatment effect estimates for subpopulations. When the Lasso models of treatment and

control outcomes are different, both in variables selected and coefficient values, this could be interpreted as modeling treatment effect heterogeneity in terms of covariates. However, reducing variance of the ATE estimate and estimating heterogenous treatment effects have completely different targets. Targeting heterogenous treatment effects may result in more variable ATE estimates. Moreover, our simulations show that the set of covariates selected by the Lasso is unstable, and this may cause problems when interpreting them as evidence of heterogenous treatment effects. How best to estimate such effects is an open question that we would like to study in future research.

## Materials and Methods

We did not conduct the PAC-man experiment, and we are analyzing secondary data without any personal identifying information. As such, this study is exempt from human subjects review. The original experiments underwent human subjects review in the United Kingdom (2).

**ACKNOWLEDGMENTS.** We thank David Goldberg for helpful discussions, Rebecca Barter for copyediting and suggestions for clarifying the text, and Winston Lin for comments. We thank Richard Grieve [London School of Hygiene and Tropical Medicine (LSHTM)], Sheila Harvey (LSHTM), David Harrison [Intensive Care National Audit and Research Centre (ICNARC)], and Kathy Rowan (ICNARC) for access to data from the PAC-Man Cost Effectiveness Analysis and the ICNARC Case Mix Programme database. This research was partially supported by NSF Grants DMS-11-06753, DMS-12-09014, DMS-1107000, DMS-1129626, DMS-1209014, Computational and Data-Enabled Science and Engineering in Mathematical and Statistical Sciences 1228246, DMS-1160319 (Focused Research Group); AFOSR Grant FA9550-14-1-0016; NSA Grant H98230-15-1-0040; the Center for Science of Information, a US NSF Science and Technology Center, under Grant Agreement CCF-0939370; Department of Defense for Office of Naval Research Grant N00014-15-1-2367; and the National Defense Science and Engineering Graduate Fellowship Program.

- Splawa-Neyman J, Dabrowska DM, Speed TP (1990) On the application of probability theory to agricultural experiments. *Essay on principles*. Section 9. *Stat Sci* 5(4):465–472.
- Harvey S, et al.; PAC-Man Study Collaboration (2005) Assessment of the clinical effectiveness of pulmonary artery catheters in management of patients in intensive care (PAC-Man): A randomised controlled trial. *Lancet* 366(9484):472–477.
- Dalen JE (2001) The pulmonary artery catheter—friend, foe, or accomplice? *JAMA* 286(3):348–350.
- Connors AF, Jr, et al.; SUPPORT Investigators (1996) The effectiveness of right heart catheterization in the initial care of critically ill patients. *JAMA* 276(11):889–897.
- Permutt T (1990) Testing for imbalance of covariates in controlled experiments. *Stat Med* 9(12):1455–1462.
- Tibshirani R (1994) Regression selection and shrinkage via the Lasso. *J R Stat Soc B* 58(1):267–288.
- Lin W (2013) Agnostic notes on regression adjustments to experimental data: Re-examining Freedman’s critique. *Ann Appl Stat* 7(1):295–318.
- Bühlmann P, Van De Geer S (2011) *Statistics for High-Dimensional Data: Methods, Theory and Applications* (Springer Science and Business Media, Berlin).
- Rubin DB (1974) Estimating causal effects of treatments in randomized and non-randomized studies. *J Educ Psychol* 66(5):688–701.
- Freedman DA (2008) On regression adjustments to experimental data. *Adv Appl Math* 40(2):180–193.
- Freedman DA (2008) On regression adjustments in experiments with several treatments. *Ann Appl Stat* 2(1):176–196.
- Freedman DA (2008) Randomization does not justify logistic regression. *Stat Sci* 23(2): 237–249.
- Miratrix LW, Sekhon JS, Yu B (2013) Adjusting treatment effect estimates by post-stratification in randomized experiments. *J R Stat Soc Series B Stat Methodol* 75(2): 369–396.
- Holland PW (1986) Statistics and causal inference. *J Am Stat Assoc* 81(396):945–960.
- Zhang CH, Zhang SS (2014) Confidence intervals for low dimensional parameters in high dimensional linear models. *J R Stat Soc Series B Stat Methodol* 76(1):217–242.
- Belloni A, Chernozhukov V, Hansen C (2013) Inference on treatment effects after selection among high-dimensional controls. *Rev Econ Stud* 81(2):608–650.
- Belloni A, Chernozhukov V, Fernández-Val I, Hansen C (2013) . Program evaluation with high-dimensional data. arXiv:1311.2645.
- Li L, Tchetgen ET, van der Vaart A, Robins JM (2011) Higher order inference on a treatment effect under low regularity conditions. *Stat Probab Lett* 81(7):821–828.
- Tian L, Alizadeh A, Gentles A, Tibshirani R (2014) A simple method for detecting interactions between a treatment and a large number of covariates. *J Am Stat Assoc* 109(508):1517–1532.
- Rosenblum M, Liu H, Yen EH (2014) Optimal tests of treatment effects for the overall population and two subpopulations in randomized trials, using sparse linear programming. *J Am Stat Assoc* 109(507):1216–1228.
- Diamond A, Sekhon JS (2013) Genetic matching for estimating causal effects: A general multivariate matching method for achieving balance in observational studies. *Rev Econ Stat* 95(3):932–945.

## Supporting Information (Appendix)

### Neyman-type conservative variance estimate

We have given a Neyman-type conservative estimate of the variance in the main text which can be used to construct a conservative confidence interval for the ATE. In this section, we will study the asymptotic behavior of this variance estimate. Recall that

$$\hat{\sigma}_{e^{(a)}}^2 = \frac{1}{n_A - df^{(a)}} \sum_{i \in A} \left( a_i - \bar{a}_A - (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T \hat{\boldsymbol{\beta}}_{\text{Lasso}}^{(a)} \right)^2, \quad (\text{S1})$$

$$\hat{\sigma}_{e^{(b)}}^2 = \frac{1}{n_B - df^{(b)}} \sum_{i \in B} \left( b_i - \bar{b}_B - (\mathbf{x}_i - \bar{\mathbf{x}}_B)^T \hat{\boldsymbol{\beta}}_{\text{Lasso}}^{(b)} \right)^2, \quad (\text{S2})$$

where  $df^{(a)}$  and  $df^{(b)}$  are degrees of freedom defined by

$$df^{(a)} = \hat{s}^{(a)} + 1 = \|\hat{\boldsymbol{\beta}}_{\text{Lasso}}^{(a)}\|_0 + 1; \quad df^{(b)} = \hat{s}^{(b)} + 1 = \|\hat{\boldsymbol{\beta}}_{\text{Lasso}}^{(b)}\|_0 + 1.$$

**Condition S0.1.** For the Gram matrix  $\Sigma$  defined in Condition 5, the largest eigenvalue is bounded away from  $\infty$ , that is, there exists a constant  $\Lambda_{max} < \infty$  such that

$$\lambda_{max}(\Sigma) \leq \Lambda_{max}.$$

Recall that the variance estimate of  $\sqrt{n}(\widehat{ATE}_{\text{Lasso}} - ATE)$  is defined as follows:

$$\hat{\sigma}_{\text{Lasso}}^2 = \frac{n}{n_A} \hat{\sigma}_{e^{(a)}}^2 + \frac{n}{n_B} \hat{\sigma}_{e^{(b)}}^2. \quad (\text{S3})$$

**Theorem S1.** Assume conditions in Theorem 1 and condition S0.1 hold. Then  $\hat{\sigma}_{\text{Lasso}}^2$  converges in probability to

$$\frac{1}{p_A} \lim_{n \rightarrow \infty} \sigma_{e^{(a)}}^2 + \frac{1}{1 - p_A} \lim_{n \rightarrow \infty} \sigma_{e^{(b)}}^2,$$

which is greater than or equal to the asymptotic variance of  $\sqrt{n}(\widehat{ATE}_{\text{Lasso}} - ATE)$ . The difference is

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \left[ a_i - b_i - ATE - (\mathbf{x}_i - \bar{\mathbf{x}})^T (\boldsymbol{\beta}^{(a)} - \boldsymbol{\beta}^{(b)}) \right]^2.$$

**Remark 6.** The Neyman-type conservative variance estimate for the unadjusted estimator is given by

$$\hat{\sigma}_{\text{unadj}}^2 = \frac{n}{n_A} \frac{1}{n_A - 1} \sum_{i \in A} (a_i - \bar{a}_A)^2 + \frac{n}{n_B} \frac{1}{n_B - 1} \sum_{i \in B} (b_i - \bar{b}_B)^2,$$

which, under second moment conditions of potential outcomes  $a$  and  $b$ , converges in probability to

$$\frac{1}{p_A} \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n (a_i - \bar{a})^2 + \frac{1}{1 - p_A} \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n (b_i - \bar{b})^2.$$

Therefore, for the  $\boldsymbol{\beta}^{(a)}$  and  $\boldsymbol{\beta}^{(b)}$  defined in [18], the limit of  $\hat{\sigma}_{\text{Lasso}}^2$  is no greater than that of  $\hat{\sigma}_{\text{unadj}}^2$  and the difference is

$$- \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \frac{1}{p_A} \left[ (\mathbf{x}_i - \bar{\mathbf{x}})^T (\boldsymbol{\beta}^{(a)}) \right]^2 + \frac{1}{1 - p_A} \left[ (\mathbf{x}_i - \bar{\mathbf{x}})^T (\boldsymbol{\beta}^{(b)}) \right]^2.$$

**Remark 7.** With the conservative variance estimate in Theorem S1, the Lasso adjusted confidence interval is also valid for the PATE (Population Average Treatment Effect) if there is a super population of size  $N$  with  $N > n$ .

**Remark 8.** The extra Condition S0.1 is used to obtain the following bounds for the number of selected covariates by the Lasso:  $\max(\hat{s}^{(a)}, \hat{s}^{(b)}) = o_p(\min(n_A, n_B))$ . Condition S0.1 can be removed from Theorem S1 if we redefine  $\hat{\sigma}_{e^{(a)}}^2$  and  $\hat{\sigma}_{e^{(b)}}^2$  without adjusting the degrees of freedom, i.e.,

$$\begin{aligned} (\hat{\sigma}^*)_{e^{(a)}}^2 &= \frac{1}{n_A} \sum_{i \in A} \left( a_i - \bar{a}_A - (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T \hat{\boldsymbol{\beta}}_{\text{Lasso}}^{(a)} \right)^2, \\ (\hat{\sigma}^*)_{e^{(b)}}^2 &= \frac{1}{n_B} \sum_{i \in B} \left( b_i - \bar{b}_B - (\mathbf{x}_i - \bar{\mathbf{x}}_B)^T \hat{\boldsymbol{\beta}}_{\text{Lasso}}^{(b)} \right)^2, \end{aligned}$$

and define  $(\hat{\sigma}^*)_{\text{Lasso}}^2 = \frac{n}{n_A} (\hat{\sigma}^*)_{e^{(a)}}^2 + \frac{n}{n_B} (\hat{\sigma}^*)_{e^{(b)}}^2$ . It follows from the bounds for  $\max(\hat{s}^{(a)}, \hat{s}^{(b)})$  that  $(\hat{\sigma}_{e^{(a)}}^2, \hat{\sigma}_{e^{(b)}}^2)$  and  $((\hat{\sigma}^*)_{e^{(a)}}^2, (\hat{\sigma}^*)_{e^{(b)}}^2)$  have the same asymptotic property.

**Theorem S2.** Assume the conditions in Theorem 1 hold. Then  $(\hat{\sigma}^*)_{\text{Lasso}}^2$  converges in probability to

$$\frac{1}{p_A} \lim_{n \rightarrow \infty} \sigma_{e^{(a)}}^2 + \frac{1}{1 - p_A} \lim_{n \rightarrow \infty} \sigma_{e^{(b)}}^2.$$

**Remark 9.** Though  $(\hat{\sigma}^*)_{\text{Lasso}}^2$  has the same limit as  $\hat{\sigma}_{\text{Lasso}}^2$ , our simulation experience shows that, in finite samples, the confidence intervals based on  $(\hat{\sigma}^*)_{\text{Lasso}}^2$  may yield low coverage probabilities (e.g., the coverage probability for 95% confidence interval can be only 80%). Hence, we recommend readers to use  $\hat{\sigma}_{\text{Lasso}}^2$  in practice.

### Simulation

In this section we carry out simulation studies to evaluate the finite sample performance of  $\widehat{ATE}_{\text{Lasso}}$  estimator. We also present results for the  $\widehat{ATE}_{\text{OLS}}$  estimator when  $p < n$  and the two-step estimator  $\widehat{ATE}_{\text{Lasso+OLS}}$  which adopts Lasso to select covariates and then uses OLS to refit the regression coefficients, see [1, 2, 3, 4] for statistical properties of Lasso+OLS estimator in linear regression model.

Let  $\hat{\beta}^{(a)}$  be the Lasso estimator defined in 1 (we omit the subscript “Lasso” for the sake of simplicity) and let  $\hat{S}^{(a)} = \{j : \hat{\beta}_j^{(a)} \neq 0\}$  be the support of  $\hat{\beta}^{(a)}$ . The Lasso+OLS adjustment vector  $\hat{\beta}_{\text{Lasso+OLS}}^{(a)}$  for treatment group A is defined by

$$\hat{\beta}_{\text{Lasso+OLS}}^{(a)} = \arg \min_{\beta: \beta_j=0, \forall j \notin \hat{S}^{(a)}} \frac{1}{2n_A} \sum_{i \in A} \left( a_i - \bar{a}_A - (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T \beta \right)^2.$$

We can define the Lasso+OLS adjustment vector  $\hat{\beta}_{\text{Lasso+OLS}}^{(b)}$  for control group B similarly. Then  $\widehat{ATE}_{\text{Lasso+OLS}}$  is given by

$$\widehat{ATE}_{\text{Lasso+OLS}} = \left[ \bar{a}_A - (\bar{\mathbf{x}}_A - \bar{\mathbf{x}})^T \hat{\beta}_{\text{Lasso+OLS}}^{(a)} \right] - \left[ \bar{b}_B - (\bar{\mathbf{x}}_B - \bar{\mathbf{x}})^T \hat{\beta}_{\text{Lasso+OLS}}^{(b)} \right].$$

We use the R package “glmnet” to compute the Lasso solution path. We select the tuning parameters  $\lambda_a$  and  $\lambda_b$  by 10-fold Cross Validation (CV) and denote the corresponding adjusted estimators as  $\text{cv}(\text{Lasso})$  and  $\text{cv}(\text{Lasso+OLS})$  respectively. We should mention that for the  $\text{cv}(\text{Lasso+OLS})$  adjusted estimator, we compute the CV error for a given value of the  $\lambda_a$  (or  $\lambda_b$ ) based on the whole Lasso+OLS estimator instead of the Lasso estimator, see Algorithm 1 for details. Therefore, the  $\text{cv}(\text{Lasso+OLS})$  adjusted estimator and the  $\text{cv}(\text{Lasso})$  adjusted estimator may select different covariates to do the adjustment. This type of cross validation for  $\text{cv}(\text{Lasso+OLS})$  requires more computation effort than the cross validation based on just the Lasso estimator since it needs to compute the OLS estimator for each fold and for each  $\lambda_a$  (or  $\lambda_b$ ), but it can give better prediction and covariates selection performance.

The potential outcomes  $a_i$  and  $b_i$  are generated from the following nonlinear model: for  $i = 1, \dots, n$ ,

$$a_i = \sum_{j=1}^s x_{ij} \beta_j^{(a1)} + \exp \left( \sum_{j=1}^s x_{ij} \beta_j^{(a2)} \right) + \epsilon_i^{(a)},$$

$$b_i = \sum_{j=1}^s x_{ij} \beta_j^{(b1)} + \exp \left( \sum_{j=1}^s x_{ij} \beta_j^{(b2)} \right) + \epsilon_i^{(b)},$$

where  $\epsilon_i^{(a)}$  and  $\epsilon_i^{(b)}$  are independent error terms. We set  $n = 250$ ,  $s = 10$ ,  $p = 50$  and 500. For  $p = 50$ , we can compute OLS estimator and compare it with the Lasso. The covariates vector  $\mathbf{x}_i$  is generated from a multivariate normal distribution  $\mathcal{N}(0, \Sigma)$ . We consider two different Toeplitz covariance matrices  $\Sigma$  which control the correlation among the covariates:

$$\Sigma_{ii} = 1; \Sigma_{ij} = \rho^{|i-j|} \quad \forall i \neq j,$$

where  $\rho = 0, 0.6$ . The true coefficients  $\beta_j^{(a1)}, \beta_j^{(a2)}, \beta_j^{(b1)}, \beta_j^{(b2)}$  are generated independently according to

$$\beta_j^{(a1)} \sim t_3; \quad \beta_j^{(a2)} \sim 0.1 * t_3, \quad j = 1, \dots, s,$$

$$\beta_j^{(b1)} \sim \beta_j^{(a1)} + t_3; \quad \beta_j^{(b2)} \sim \beta_j^{(a2)} + 0.1 * t_3, \quad j = 1, \dots, s,$$

where  $t_3$  denotes the  $t$  distribution with three degrees of freedom. This ensures that the treatment effects are not constant across individuals, and that the linear model does not hold in this simulation. The error terms  $\epsilon_i^{(a)}$  and  $\epsilon_i^{(b)}$  are generated according to the following linear model with hidden covariates  $\mathbf{z}_i$ :

$$\epsilon_i^{(a)} = \sum_{j=1}^s z_{ij} \beta_j^{(a1)} + \tilde{\epsilon}_i^{(a)},$$

$$\epsilon_i^{(b)} = \sum_{j=1}^s z_{ij} \beta_j^{(b1)} + \tilde{\epsilon}_i^{(b)},$$

where  $\tilde{\epsilon}_i^{(a)}$  and  $\tilde{\epsilon}_i^{(b)}$  are drawn independently from a standard normal distribution. The vector  $\mathbf{z}_i$  is independent of  $\mathbf{x}_i$  and is also drawn independently from the multivariate normal distribution  $\mathcal{N}(0, \Sigma)$ . The values of  $\mathbf{x}_i$ ,  $\beta^{(a1)}$ ,  $\beta^{(a2)}$ ,  $\beta^{(b1)}$ ,  $\beta^{(b2)}$ ,  $\mathbf{z}_i$ ,  $\tilde{\epsilon}_i^{(a)}$ ,  $\tilde{\epsilon}_i^{(b)}$ ,  $a_i$  and  $b_i$  are generated once and then kept fixed.

After the potential outcomes are generated, a completely randomized experiment is simulated 25000 times, assigning  $n_A = 100, 125, 150$  subjects to treatment A and the remainder to control B. There are 12 different combinations of  $(p, \rho, n_A)$  in total.

Figures S4, S5, S6 show boxplots of different ATE estimators with their standard deviations (computed from 25000 replicates of randomized experiments) presented on top of each box. Regardless of whether the design is balanced ( $n_A = 125$ ) or not ( $n_A = 100, 150$ ), the regression based estimators have much smaller variances and than that of the unadjusted estimator and therefore improve the estimation precision.

To further compare the performance of these estimators, we present the bias, the standard deviation (SD) and the root-mean square error ( $\sqrt{\text{MSE}}$ ) of the estimates in Table S1. Bias is reported as the absolute difference from the true treatment effect. We find that the bias of each method is substantially smaller (more than 10 times smaller) than the SD. The cv(Lasso) and cv(Lasso+OLS) adjusted estimators perform similar in terms of SD and  $\sqrt{\text{MSE}}$ : reducing those of the OLS adjusted estimator and the unadjusted estimator by 10% – 15% and 15% – 31% respectively. We also compare the number of selected covariates by cv(Lasso) and cv(Lasso+OLS) for treatment group and control group separately, see Table S2. It is easy to see that the cv(Lasso+OLS) adjusted estimator uses many fewer (more than 44%) covariates in the adjustment to obtain similar improvement of SD and  $\sqrt{\text{MSE}}$  of ATE estimate as the cv(Lasso) adjusted estimator. Moreover, we find that the covariates selected by the cv(Lasso+OLS) are more stable across different realizations of treatment assignment than the covariates selected by the cv(Lasso). Overall, the cv(Lasso+OLS) adjusted, the cv(Lasso) adjusted, the OLS adjusted and the unadjusted estimators perform from best to worst.

We move now to study the finite sample performance of Neyman-type conservative variance estimates. For each simulation example and each one of the 25000 completely randomized experiments, we calculate the ATE estimates ( $\widehat{ATE}$ ) and the Neyman variance estimates ( $\hat{\sigma}$ ) and then form the 95% confidence intervals  $[\widehat{ATE} - 1.96 \cdot \hat{\sigma}/\sqrt{n}, \widehat{ATE} + 1.96 \cdot \hat{\sigma}/\sqrt{n}]$ . Figures S1, S2, S3 present the boxplot of the interval length with the coverage probability noted on top of each box for the unadjusted, OLS adjusted (only computed when  $p = 50$ ), cv(Lasso) adjusted and cv(Lasso+OLS) adjusted estimators. More results are showed in Table S3. We find that all the confidence intervals for the unadjusted estimator are conservative. The cv(Lasso) adjusted and the cv(Lasso+OLS) adjusted estimators perform very similar: although their coverage probability (at least 92%) may be slightly less than the pre-assigned confidence level (95%), their mean interval length is much shorter (26% – 37%) than that of the unadjusted estimator. The OLS adjusted estimator has comparable interval length with the cv(Lasso) and cv(Lasso+OLS) adjusted estimator, but has slightly worse coverage probability (90% – 93%).

To further investigate how good the Neyman standard deviation (SD) estimate is, we compare them in Figure S7 with the “true” SD presented in Table S1 (the SD of the ATE estimates over 25000 randomized experiments). We find that Neyman SD estimate is very conservative for the unadjusted estimator (its mean is 5% – 14% larger than the “true” SD); while for the OLS adjusted estimator, the mean of Neyman SD estimate can be 6% – 100% smaller than the “true” SD which may be because of over-fitting. For the cv(Lasso) and cv(Lasso+OLS) adjusted estimator, the mean of Neyman SD estimator is within  $1 \pm 7\%$  of the “true” SD. Although the Neyman variance estimate is asymptotically conservative, in small samples the variance estimate may be too small and hence the confidence intervals are too narrow. However, if we increase the sample size  $n$  to 1000, almost all the confidence intervals are conservative.

We conduct more simulation examples to evaluate the conditions assumed for asymptotic normality of the Lasso adjusted estimator. We use the same simulation setup as above, but for simplicity, we generate the potential outcomes from a linear model (set  $\beta^{(a2)} = \beta^{(b2)} = 0$ ) and remove the effects of the hidden covariates  $z_i$  in generating the error terms  $\epsilon_i^{(a)}$  and  $\epsilon_i^{(b)}$  and set  $\rho = 0$ ,  $n_A = 125$ . We find that the distribution of the cv(Lasso) adjusted estimator may be non-normal when:

- (1). The covariates are generated from Gaussian distribution and the error terms do not satisfy second moment condition, e.g., being generated from  $t$  distribution with one degree of freedom, see the upper two subplots of Figure 1 (in the main text) for the histograms of unadjusted the cv(Lasso) adjusted estimators (the corresponding p-values of Kolmogorov–Smirnov testing for normality are less than  $2.2e - 16$ ).
- (2). The covariates do not have bounded fourth moments, e.g., being generated from  $t$  distribution with three degrees of freedom, see the lower two subplots of Figure 1 (in the main text) for the histograms of unadjusted the cv(Lasso) adjusted estimators (again, the corresponding p-values of Kolmogorov–Smirnov testing for normality are less than  $2.2e - 16$ ).

These findings indicate that our moment condition (Condition 2 and Remark 1) cannot be dramatically weakened. However, we also find that the cv(Lasso) adjusted estimator still has smaller SD and  $\sqrt{\text{MSE}}$  than the unadjusted estimator even when these moment conditions do not hold.

### The design matrix of the PAC data

In the PAC data, there are 59 covariates (main effects) including 50 indicators which may be correlated with the outcomes. One of the main effects (called interactnew) has heavy tail, so we perform the transform:  $x \rightarrow \log(|x| + 1)$  to make it more normally distributed. We then centralize and standardize the non-indicator covariates. The quadratic terms (9 in total) of non-indicator covariates and two-way interactions between main effects (1711 in total) may also help predict the potential outcomes, so we included them in the design matrix. The quadratic terms and the interactions between non-indicator covariates and the interactions between indicator covariates and non-indicator covariates are also centered and standardized. Some of the interactions are identical to other effects and we only retain one of them. We also remove the interactions which are highly correlated (with correlation larger than 0.95) with the main effects and remove the indicators with very sparse entries (where the number of 1’s is less than 20). Finally, we form a design matrix  $X$  with 1172 columns (covariates) and 1013 rows (subjects).

### Estimation of constants in the conditions

Let  $S^{(a)} = \{j : \beta_j^{(a)} \neq 0\}$  and  $S^{(b)} = \{j : \beta_j^{(b)} \neq 0\}$  denote the sets of relevant covariates for treatment group and control group respectively. Denote  $S = S^{(a)} \cup S^{(b)} = \{j : \beta_j^{(a)} \neq 0 \text{ or } \beta_j^{(b)} \neq 0\}$ . We use bootstrap to get an estimation of the relevant covariates sets  $S^{(a)}, S^{(b)}$  and then the approximation errors  $e^{(a)}$  and  $e^{(b)}$  are estimated by regressing the observed potential outcomes  $a$  and  $b$  on the covariates in  $S$  respectively. We only present how to estimate  $S^{(a)}$  and  $e^{(a)}$  in detail and the estimation of  $S^{(b)}$  and  $e^{(b)}$  are similar.

Let  $A, B$  be the set of treated subjects (using PAC) and control subjects (without using PAC) respectively. Denote  $a_i, i \in A$  the potential outcomes (quality-adjusted life years (QALYs)) under treatment and  $x_i \in R^{1172}$  the covariates vector of the  $i$ th subject. For each  $d = 1, \dots, 1000$ , we draw a bootstrap sample  $\{(a_i^*(d), x_i^*(d)) : i \in A\}$  with replacement from the data points  $\{(a_i, x_i) : i \in A\}$ . We then compute the LassoOLS(CV) adjusted vector  $\hat{\beta}(d)$  based on each bootstrap sample  $\{(a_i^*(d), x_i^*(d)) : i \in A\}$ . Let  $\tau_j$  be the selection fraction of non-zero  $\hat{\beta}_j(d)$  in the 1000 bootstrap estimators, i.e.,  $\tau_j = (1/1000) \sum_{d=1}^{1000} \mathbb{I}_{\{\hat{\beta}_j(d) \neq 0\}}$ , where  $\mathbb{I}$  is the indicator function. We form the relevant covariates  $S^{(a)}$  by the covariates whose selection fraction are larger than 0.5:  $S^{(a)} = \{j : \tau_j > 0.5\}$ .

To estimate the approximation error  $e^{(a)}$ , we regress  $a_i$  on the relevant covariates  $x_{ij}, j \in S^{(a)}$  and compute OLS estimate and the corresponding residual. That is, let  $T^{(a)}$  denote the complement set of  $S^{(a)}$ ,

$$\beta_{\text{OLS}}^{(a)} = \arg \min_{\beta: \beta_j=0, \forall j \in T^{(a)}} \frac{1}{2n_A} \sum_{i \in A} \left( a_i - \bar{a}_A - (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T \beta \right)^2.$$

$$e_i^{(a)} = a_i - \bar{a}_A - (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T \beta_{\text{OLS}}^{(a)}, \quad i \in A.$$

The maximal covariance  $\delta_n$  is estimated as:

$$\max \left\{ \max_j \left| \frac{1}{n_A} \sum_{i \in A} (x_{ij} - (\bar{\mathbf{x}})_j) (e_i^{(a)} - \bar{e}_A^{(a)}) \right|, \max_j \left| \frac{1}{n_B} \sum_{i \in B} (x_{ij} - (\bar{\mathbf{x}})_j) (e_i^{(b)} - \bar{e}_B^{(b)}) \right| \right\}.$$

### Proofs of Theorems 1, S1, S2 and Corollary 1

In this section, we will prove Theorem 1, S1, S2, and Corollary 1 under weaker sparsity conditions than those given in the main text.

**Definition 1.** We define an approximate sparsity measure. Given the regularization parameter  $\lambda_a, \lambda_b$  and  $\beta^{(a)}$  and  $\beta^{(b)}$ , the sparsity measures for treatment and control groups,  $s_{\lambda_a}^{(a)}$  and  $s_{\lambda_b}^{(b)}$  are defined as

$$s_{\lambda_a}^{(a)} = \sum_{j=1}^p \min \left\{ \frac{|\beta_j^{(a)}|}{\lambda_a}, 1 \right\}, \quad s_{\lambda_b}^{(b)} = \sum_{j=1}^p \min \left\{ \frac{|\beta_j^{(b)}|}{\lambda_b}, 1 \right\}, \quad (\text{S4})$$

respectively. We will allow  $s_{\lambda_a}^{(a)}$  and  $s_{\lambda_b}^{(b)}$  to grow with  $n$ , though the notation does not explicitly show this. Note that this is weaker than strict sparsity, as it allows  $\beta^{(a)}$  and  $\beta^{(b)}$  to have many small non-zero entries.

**Condition (\*).** Suppose there exist  $\beta^{(a)}, \beta^{(b)}, \lambda_a$  and  $\lambda_b$  such that the conditions 1, 2, 3 and the following statements 1, 2, 3 hold simultaneously.

- **Statement 1.** Decay and scaling. Let  $s_\lambda = \max \{s_{\lambda_a}^{(a)}, s_{\lambda_b}^{(b)}\}$ ,

$$\delta_n = o \left( \frac{1}{s_\lambda \sqrt{\log p}} \right), \quad (\text{S5})$$

$$(s_\lambda \log p) / \sqrt{n} = o(1). \quad (\text{S6})$$

- **Statement 2.** Cone invertibility factor. Define the Gram matrix as  $\Sigma = n^{-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$ . There exist constants  $C > 0$  and  $\xi > 1$  not depending on  $n$ , such that

$$\|\mathbf{h}_S\|_1 \leq C s_\lambda \|\Sigma \mathbf{h}\|_\infty, \quad \forall \mathbf{h} \in \mathcal{C}, \quad (\text{S7})$$

with  $\mathcal{C} = \{\mathbf{h} : \|\mathbf{h}_{S^c}\|_1 \leq \xi \|\mathbf{h}_S\|_1\}$ , and

$$S = \{j : |\beta_j^{(a)}| > \lambda_a \text{ or } |\beta_j^{(b)}| > \lambda_b\}. \quad (\text{S8})$$

- **Statement 3.** Let  $\tau = \min \{1/70, (3p_A)^2/70, (3 - 3p_A)^2/70\}$ . For constants  $0 < \eta < \frac{\xi-1}{\xi+1}$  and  $0 < M < \infty$ , assume the regularization parameters of the Lasso belong to the sets

$$\lambda_a \in \left( \frac{1}{\eta}, M \right] \times \left( \frac{2(1+\tau)L^{1/2}}{p_A} \sqrt{\frac{2 \log p}{n}} + \delta_n \right), \quad (\text{S9})$$

$$\lambda_b \in \left(\frac{1}{\eta}, M\right] \times \left(\frac{2(1+\tau)L^{1/2}}{p_B} \sqrt{\frac{2 \log p}{n}} + \delta_n\right). \quad (\text{S10})$$

It is easy to verify that Condition (\*) is implied by Conditions 1 - 6 of the main text. We will prove Theorems 1, S1, S2, and Corollary 1 under the weaker Condition (\*). For ease of notation, we will omit the subscript of  $\hat{\beta}_{\text{Lasso}}^{(a)}$ ,  $\hat{\beta}_{\text{Lasso}}^{(b)}$ ,  $s_\lambda$ ,  $s_{\lambda_a}^{(a)}$  and  $s_{\lambda_b}^{(b)}$ . Note that we can assume, without loss of generality, that

$$\bar{a} = 0, \bar{b} = 0, \bar{\mathbf{x}} = \mathbf{0}. \quad (\text{S11})$$

Otherwise, we can consider  $\check{a}_i = a_i - \bar{a}$ ,  $\check{b}_i = b_i - \bar{b}$  and  $\check{\mathbf{x}}_i = \mathbf{x}_i - \bar{\mathbf{x}}$ . Thus, we assume  $\text{ATE} = \bar{a} - \bar{b} = 0$  and the definition of  $\widehat{\text{ATE}}_{\text{Lasso}}$  becomes

$$\widehat{\text{ATE}}_{\text{Lasso}} = \left[\bar{a}_A - (\bar{\mathbf{x}}_A)^T \hat{\beta}^{(a)}\right] - \left[\bar{b}_B - (\bar{\mathbf{x}}_B)^T \hat{\beta}^{(b)}\right]. \quad (\text{S12})$$

We will rely heavily on the following Massart concentration inequality for sampling without replacement.

**Lemma S1.** *Let  $\{z_i, i = 1, \dots, n\}$  be a finite population of real numbers. Let  $A \subset \{1, \dots, n\}$  be a subset of deterministic size  $|A| = n_A$  that is selected randomly without replacement. Define  $p_A = n_A/n$ ,  $\sigma^2 = n^{-1} \sum_{i=1}^n (z_i - \bar{z})^2$ . Then, for any  $t > 0$ ,*

$$P(\bar{z}_A - \bar{z} \geq t) \leq \exp\left\{-\frac{p_A n_A t^2}{(1+\tau)^2 \sigma^2}\right\}, \quad (\text{S13})$$

with  $\tau = \min\{1/70, (3p_A)^2/70, (3-3p_A)^2/70\}$ .

**Remark.** *Massart showed in his paper [5] that for sampling without replacement, the following concentration inequality holds:*

$$P(\bar{z}_A - \bar{z} \geq t) \leq \exp\left\{-\frac{p_A n_A t^2}{\sigma^2}\right\}.$$

*His proof required that  $n/n_A$  must be an integer. We extend the proof to allow  $n/n_A$  to be a non-integer but with a slightly larger constant factor  $(1+\tau)^2$ .*

*Proof.* Assume  $\bar{z} = 0$  without loss of generality. For  $n_A \leq n/2$ , let  $m \geq 2$  and  $r \geq 0$  be integers satisfying  $n - n_A m = r < n_A$ . Let  $u \geq 0$ , we first prove that

$$E \exp\left(u \sum_{i \in A} z_i\right) \leq E \exp\left(u \delta \sum_{i \in B} z_i / \{m(m+1)\} + u^2 n \sigma^2 / 4\right) \quad (\text{S14})$$

for a random subset  $B \subset \{1, \dots, n\}$  of fixed size  $|B| \leq n/2$  and a certain fixed  $\delta \in \{-1, 1\}$ . Let  $P_1$  be the probability under which  $\{i_1, \dots, i_n\}$  is a random permutation of  $\{1, \dots, n\}$ . Given  $\{i_1, \dots, i_n\}$ , we divide the sequence into consecutive blocks  $B_1, \dots, B_{n_A}$  with  $|B_j| = m+1$  for  $j = 1, \dots, r$  and  $|B_j| = m$  for  $j = r+1, \dots, n_A$ . Let  $\bar{z}_k$  be the mean of  $\{z_i : i \in B_k\}$  and  $P_2$  be a probability conditionally on  $\{i_1, \dots, i_n\}$  under which  $w_k$  is a random element of  $\{z_i : i \in B_k\}$ ,  $k = 1, \dots, n_A$ . Then  $\{w_1, \dots, w_{n_A}\}$  is a random sample from  $\{z_1, \dots, z_n\}$  without replacement under  $P = P_1 P_2$ . Let  $\Delta_k = \max_{i \in B_k} z_i - \min_{i \in B_k} z_i$  and denote  $E_2$  the expectation under  $P_2$ . The Hoeffding inequality gives

$$E_2 \exp\left(u \sum_{k=1}^{n_A} w_k\right) \leq \exp\left(u \sum_{k=1}^{n_A} \bar{z}_k + (u^2/8) \sum_{k=1}^{n_A} \Delta_k^2\right). \quad (\text{S15})$$

As  $\Delta_k^2 \leq 2 \sum_{i \in B_k} (z_i - \bar{z}_k)^2 \leq 2 \sum_{i \in B_k} z_i^2$ ,

$$E_2 \exp\left(u \sum_{k=1}^{n_A} w_k\right) \leq \exp\left(u \sum_{k=1}^{n_A} \bar{z}_k + u^2 n \sigma^2 / 4\right) \quad (\text{S16})$$

Let  $B = \cup_{k=1}^r B_k$ . As  $\bar{z} = 0$ ,

$$\sum_{k=1}^{n_A} \bar{z}_k = \sum_{i \in B} z_i / \{m(m+1)\}. \quad (\text{S17})$$

This yields (S14) with  $\delta = 1$  when  $|B| \leq n/2$ . Otherwise, (S14) holds with  $\delta = -1$  when  $B$  is replaced by  $B^c$ , as  $\sum_{i \in B} z_i = -\sum_{i \in B^c} z_i$  due to  $\bar{z} = 0$ .

Now, as  $m(m+1) \geq 6$ , repeated application of (S14) yields

$$\begin{aligned} E \exp\left(u \sum_{i \in A} z_i\right) &\leq E \exp\left[u \delta' \sum_{i \in B'} z_i / \{m(m+1)m'(m'+1)\} + (1 + \{m(m+1)\}^{-2}) u^2 n \sigma^2 / 4\right] \\ &\leq \exp\left[(1 + \{m(m+1)\}^{-2}(1 + 1/36 + 1/36^2 + \dots)) u^2 n \sigma^2 / 4\right] \\ &= \exp\left[(1 + (36/35)\{m(m+1)\}^{-2}) u^2 n \sigma^2 / 4\right] \\ &\leq \exp\left[(1 + \tau) u^2 n \sigma^2 / 4\right] \end{aligned}$$

with  $\tau = (18/35)\{m(m+1)\}^{-2}$ . The upper bound for  $\tau$  follows from  $2 \leq m < n/n_A < m+1$ . As  $\bar{z} = 0$ , we also have

$$E \exp \left( u \sum_{i \in A} z_i \right) \leq \exp \left[ (1+\tau)^2 u^2 n \sigma^2 / 4 \right] \quad (\text{S19})$$

for  $n_A > n/2$ . This yields (S13) via the usual

$$\begin{aligned} P \{ \bar{z}_A - \bar{z} > t \} &\leq \exp \left[ -ut + (1+\tau)^2 u^2 n \sigma^2 / (4n_A^2) \right] \\ &= \exp \left[ -2 \frac{p_A n_A t^2}{(1+\tau)^2 \sigma^2} + \frac{p_A n_A t^2}{(1+\tau)^2 \sigma^2} \right] \end{aligned} \quad (\text{S20})$$

with  $u = 2p_A n_A t / \{\sigma(1+\tau)\}^2$ . □

**Proof of Theorem 1.**

*Proof.* Recall the decompositions of the potential outcomes:

$$a_i = \bar{a} + (\mathbf{x}_i - \bar{\mathbf{x}})^T \boldsymbol{\beta}^{(a)} + e_i^{(a)} = \mathbf{x}_i^T \boldsymbol{\beta}^{(a)} + e_i^{(a)}, \quad (\text{S21})$$

$$b_i = \bar{b} + (\mathbf{x}_i - \bar{\mathbf{x}})^T \boldsymbol{\beta}^{(b)} + e_i^{(b)} = \mathbf{x}_i^T \boldsymbol{\beta}^{(b)} + e_i^{(b)}. \quad (\text{S22})$$

If we define  $\mathbf{h}^{(a)} = \hat{\boldsymbol{\beta}}^{(a)} - \boldsymbol{\beta}^{(a)}$ ,  $\mathbf{h}^{(b)} = \hat{\boldsymbol{\beta}}^{(b)} - \boldsymbol{\beta}^{(b)}$ , by substitution, we have

$$\sqrt{n}(\widehat{ATE}_{\text{Lasso}} - ATE) = \underbrace{\sqrt{n} \left[ \bar{e}_A^{(a)} - \bar{e}_B^{(b)} \right]}_{*} - \underbrace{\sqrt{n} \left[ (\bar{\mathbf{x}}_A)^T \mathbf{h}^{(a)} - (\bar{\mathbf{x}}_B)^T \mathbf{h}^{(b)} \right]}_{**}.$$

We will analyze these two terms separately, showing that (\*) is asymptotically normal with mean 0 and variance given by (17), and that (\*\*) is  $o_p(1)$ .

Asymptotic normality of (\*) follows from the Theorem 1 in [6] with  $a$  and  $b$  replaced by  $e^{(a)}$  and  $e^{(b)}$  respectively. To bound (\*\*), we will apply Hölder’s inequality to each of the terms. We will focus on the term involving the treatment group  $A$ , but exact same analysis is applied to the control group  $B$ . We have the bound

$$\left| (\bar{\mathbf{x}}_A)^T \mathbf{h}^{(a)} \right| \leq \|\bar{\mathbf{x}}_A\|_\infty \|\mathbf{h}^{(a)}\|_1. \quad (\text{S23})$$

We bound the two terms on the right hand side of (S23) by the following Lemma S2 and Lemma S3, respectively.

**Lemma S2.** Under the moment condition of [6], if we let  $c_n = \frac{(1+\tau)L^{1/4}}{p_A} \sqrt{\frac{2 \log p}{n}}$ , then as  $n \rightarrow \infty$ ,

$$P(\|\bar{\mathbf{x}}_A\|_\infty > c_n) \rightarrow 0$$

Thus,  $\|\bar{\mathbf{x}}_A\|_\infty = O_p\left(\sqrt{\frac{\log p}{n}}\right)$ .

**Lemma S3.** Assume the conditions of Theorem 1 hold. Then  $\|\mathbf{h}^{(a)}\|_1 = o_p\left(\frac{1}{\sqrt{\log p}}\right)$ .

The proofs of these two Lemmas are below. Using these two Lemmas, it is easy to show that (\*\*) =  $\sqrt{n} \cdot O_p\left(\sqrt{\frac{\log p}{n}}\right) \cdot o_p\left(\frac{1}{\sqrt{\log p}}\right) = o_p(1)$ . □

**Proof of Corollary 1.**

*Proof.* By Theorem 1 in [6], the asymptotic variance of  $\sqrt{n} \widehat{ATE}_{\text{unadj}}$  is  $\frac{1-p_A}{p_A} \lim_{n \rightarrow \infty} \sigma_a^2 + \frac{p_A}{1-p_A} \lim_{n \rightarrow \infty} \sigma_b^2 + 2 \lim_{n \rightarrow \infty} \sigma_{ab}$ , so the difference is

$$\frac{1-p_A}{p_A} \lim_{n \rightarrow \infty} (\sigma_{e^{(a)}}^2 - \sigma_a^2) + \frac{p_A}{1-p_A} \lim_{n \rightarrow \infty} (\sigma_{e^{(b)}}^2 - \sigma_b^2) + 2 \lim_{n \rightarrow \infty} (\sigma_{e^{(a)}e^{(b)}} - \sigma_{ab}).$$

We will analyze these three terms separately. Since  $X\boldsymbol{\beta}^{(a)}$  and  $X\boldsymbol{\beta}^{(b)}$  are the orthogonal projections of  $a$  and  $b$  onto the same subspace, we have

$$(X\boldsymbol{\beta}^{(a)})^T e^{(a)} = (X\boldsymbol{\beta}^{(a)})^T e^{(b)} = (X\boldsymbol{\beta}^{(b)})^T e^{(a)} = (X\boldsymbol{\beta}^{(b)})^T e^{(b)} = 0.$$

Simple calculations yield

$$\sigma_{e^{(a)}}^2 - \sigma_a^2 = \|e^{(a)}\|_2^2 - \|a\|_2^2 = -\|X\boldsymbol{\beta}^{(a)}\|_2^2, \quad (\text{S24})$$

$$\sigma_{e^{(b)}}^2 - \sigma_b^2 = \|e^{(b)}\|_2^2 - \|b\|_2^2 = -\|X\boldsymbol{\beta}^{(b)}\|_2^2, \quad (\text{S25})$$

$$\sigma_{e^{(a)}e^{(b)}} - \sigma_{ab} = (e^{(a)})^T(e^{(b)}) - a^T b = -(X\boldsymbol{\beta}^{(a)})^T(X\boldsymbol{\beta}^{(b)}) \quad (\text{S26})$$

Combining (S24), (S25) and (S26), we obtain the corollary.  $\square$

**Proof of Theorem S1.**

*Proof.* To prove Theorem S1, it is enough to show that

$$\hat{\sigma}_{e^{(a)}}^2 \xrightarrow{p} \lim_{n \rightarrow \infty} \sigma_{e^{(a)}}^2, \quad (\text{S27})$$

$$\hat{\sigma}_{e^{(b)}}^2 \xrightarrow{p} \lim_{n \rightarrow \infty} \sigma_{e^{(b)}}^2. \quad (\text{S28})$$

We will only prove the statement (S27) and omit the proof of the statement (S28) since it is identical.

We first state the following two lemmas. Lemma S4 bounds the number of selected covariates (covariates with a nonzero coefficient), while Lemma S5 establishes conditions under which the subsample mean (without replacement) converges in probability to the population mean.

**Lemma S4.** *Under conditions in Theorem S1, there exists a constant  $C$ , such that the following holds with probability going to 1:*

$$\hat{s}^{(a)} \leq Cs; \quad \hat{s}^{(b)} \leq Cs. \quad (\text{S29})$$

The proof of Lemma S4 can be found below.

**Lemma S5.** *Let  $\{z_i, i = 1, \dots, n\}$  be a finite population of real numbers. Let  $A \subset \{i, \dots, n\}$  be a subset of deterministic size  $|A| = n_A$  that is selected randomly without replacement. Suppose that the population mean of the  $z_i$  has a finite limit and that there exist constants  $\epsilon > 0$  and  $L < \infty$  such that*

$$\frac{1}{n} \sum_{i=1}^n |z_i|^{1+\epsilon} \leq L. \quad (\text{S30})$$

If  $\frac{n_A}{n} \rightarrow p_A \in (0, 1)$ , then

$$\bar{z}_A \xrightarrow{p} \lim_{n \rightarrow \infty} \bar{z}. \quad (\text{S31})$$

By definition (S1) and simple calculations,

$$\begin{aligned} \hat{\sigma}_{e^{(a)}}^2 &= \frac{1}{n_A - df^{(a)}} \sum_{i \in A} \left( a_i - \bar{a}_A - (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T \hat{\boldsymbol{\beta}}^{(a)} \right)^2 \\ &= \frac{1}{n_A - df^{(a)}} \sum_{i \in A} \left( a_i - \bar{a}_A - (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T \boldsymbol{\beta}^{(a)} + (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T (\boldsymbol{\beta}^{(a)} - \hat{\boldsymbol{\beta}}^{(a)}) \right)^2 \\ &= \frac{1}{n_A - df^{(a)}} \sum_{i \in A} \left( a_i - \mathbf{x}_i^T \boldsymbol{\beta}^{(a)} - (\bar{a}_A - (\bar{\mathbf{x}}_A)^T \boldsymbol{\beta}^{(a)}) + (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T (\boldsymbol{\beta}^{(a)} - \hat{\boldsymbol{\beta}}^{(a)}) \right)^2 \\ &= \frac{n_A}{n_A - df^{(a)}} \frac{1}{n_A} \sum_{i \in A} \left( e_i^{(a)} - \bar{e}_A^{(a)} + (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T (\boldsymbol{\beta}^{(a)} - \hat{\boldsymbol{\beta}}^{(a)}) \right)^2 \\ &= \frac{n_A}{n_A - df^{(a)}} \left\{ \frac{1}{n_A} \sum_{i \in A} \left( e_i^{(a)} - \bar{e}_A^{(a)} \right)^2 + \frac{1}{n_A} \sum_{i \in A} \left( (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T (\boldsymbol{\beta}^{(a)} - \hat{\boldsymbol{\beta}}^{(a)}) \right)^2 \right\} \\ &\quad + \frac{n_A}{n_A - df^{(a)}} \left\{ \frac{1}{n_A} \sum_{i \in A} (e_i^{(a)} - \bar{e}_A^{(a)}) (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T (\boldsymbol{\beta}^{(a)} - \hat{\boldsymbol{\beta}}^{(a)}) \right\}. \end{aligned}$$

The second to last equality is due to the decomposition of potential outcome  $a$ :

$$a_i = \mathbf{x}_i^T \boldsymbol{\beta}^{(a)} + e_i^{(a)}; \quad \bar{a}_A = (\bar{\mathbf{x}}_A)^T \boldsymbol{\beta}^{(a)} + \bar{e}_A^{(a)}.$$

It is easy to see that

$$\frac{1}{n_A} \sum_{i \in A} \left( e_i^{(a)} - \bar{e}_A^{(a)} \right)^2 = \frac{1}{n_A} \sum_{i \in A} (e_i^{(a)})^2 - (\bar{e}_A^{(a)})^2. \quad (\text{S32})$$

By the 4th moment condition on the approximation error  $e^{(a)}$  (see (7)), and applying Lemma S5 gives

$$\frac{1}{n_A} \sum_{i \in A} (e_i^{(a)})^2 \xrightarrow{p} \lim_{n \rightarrow \infty} \sigma_{e^{(a)}}^2; \quad \bar{e}_A^{(a)} \xrightarrow{p} \lim_{n \rightarrow \infty} \bar{e}^{(a)} = 0.$$

Therefore,

$$\frac{1}{n_A} \sum_{i \in A} \left( e_i^{(a)} - \bar{e}_A^{(a)} \right)^2 \xrightarrow{p} \lim_{n \rightarrow \infty} \sigma_{e^{(a)}}^2. \quad (\text{S33})$$

By simple algebra,

$$\begin{aligned}
 & \frac{1}{n_A} \sum_{i \in A} \left( (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T (\boldsymbol{\beta}^{(a)} - \hat{\boldsymbol{\beta}}^{(a)}) \right)^2 \\
 &= (\boldsymbol{\beta}^{(a)} - \hat{\boldsymbol{\beta}}^{(a)})^T \left[ \frac{1}{n_A} \sum_{i \in A} (\mathbf{x}_i - \bar{\mathbf{x}}_A) (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T \right] (\boldsymbol{\beta}^{(a)} - \hat{\boldsymbol{\beta}}^{(a)}) \\
 &\leq \|\boldsymbol{\beta}^{(a)} - \hat{\boldsymbol{\beta}}^{(a)}\|_1 \cdot \left\| \frac{1}{n_A} \sum_{i \in A} (\mathbf{x}_i - \bar{\mathbf{x}}_A) (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T \right\|_\infty.
 \end{aligned} \tag{S34}$$

We next show that (S34) converges to 0 in probability. By Lemma S3 and Lemma S7, we have

$$\|\boldsymbol{\beta}^{(a)} - \hat{\boldsymbol{\beta}}^{(a)}\|_1 = \|\mathbf{h}^{(a)}\|_1 = o_p \left( \frac{1}{\sqrt{\log p}} \right), \tag{S35}$$

$$\left\| \frac{1}{n_A} \sum_{i \in A} (\mathbf{x}_i - \bar{\mathbf{x}}_A) (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T \right\|_\infty = O_p(1). \tag{S36}$$

Therefore,

$$\frac{1}{n_A} \sum_{i \in A} \left( (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T (\boldsymbol{\beta}^{(a)} - \hat{\boldsymbol{\beta}}^{(a)}) \right)^2 \xrightarrow{p} 0. \tag{S37}$$

By Cauchy-Schwarz inequality,

$$\begin{aligned}
 & \left| \frac{1}{n_A} \sum_{i \in A} (e_i^{(a)} - \bar{e}_A^{(a)}) (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T (\boldsymbol{\beta}^{(a)} - \hat{\boldsymbol{\beta}}^{(a)}) \right| \\
 &\leq \left[ \frac{1}{n_A} \sum_{i \in A} (e_i^{(a)} - \bar{e}_A^{(a)})^2 \right]^{\frac{1}{2}} \left[ \frac{1}{n_A} \sum_{i \in A} \left( (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T (\boldsymbol{\beta}^{(a)} - \hat{\boldsymbol{\beta}}^{(a)}) \right)^2 \right]^{\frac{1}{2}}
 \end{aligned} \tag{S38}$$

which converges to 0 in probability because of (S33) and (S37).

By Lemma S4 and Condition 4, we have

$$\frac{n_A}{n_A - df^{(a)}} = \frac{n_A}{n_A - \hat{s}^{(a)} - 1} \xrightarrow{p} 1. \tag{S39}$$

Combining (S33), (S37), (S38) and (S39), we conclude that

$$\hat{\sigma}_{e^{(a)}}^2 \xrightarrow{p} \lim_{n \rightarrow \infty} \sigma_{e^{(a)}}^2.$$

The remaining part of the proof is to study the difference between the conservative variance estimate and the true asymptotic variance:

$$\begin{aligned}
 & \left( \frac{1}{p_A} \lim_{n \rightarrow \infty} \sigma_{e^{(a)}}^2 + \frac{1}{1-p_A} \lim_{n \rightarrow \infty} \sigma_{e^{(b)}}^2 \right) - \left( \frac{1-p_A}{p_A} \lim_{n \rightarrow \infty} \sigma_{e^{(a)}}^2 + \frac{p_A}{1-p_A} \lim_{n \rightarrow \infty} \sigma_{e^{(b)}}^2 + 2 \lim_{n \rightarrow \infty} \sigma_{e^{(a)}e^{(b)}} \right) \\
 &= \lim_{n \rightarrow \infty} \sigma_{e^{(a)}}^2 + \lim_{n \rightarrow \infty} \sigma_{e^{(b)}}^2 - 2 \lim_{n \rightarrow \infty} \sigma_{e^{(a)}e^{(b)}} \\
 &= \lim_{n \rightarrow \infty} \sigma_{e^{(a)} - e^{(b)}}^2 \\
 &= \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \left( a_i - b_i - \mathbf{x}_i^T (\boldsymbol{\beta}^{(a)} - \boldsymbol{\beta}^{(b)}) \right)^2.
 \end{aligned} \tag{S40}$$

□

### Proof of Theorem S2.

*Proof.* By Lemma S4,  $\max(\hat{s}^{(a)}, \hat{s}^{(b)}) = o_p(\min(n_A, n_B))$ . Therefore,  $(\hat{\sigma}_{e^{(a)}}^2, \hat{\sigma}_{e^{(b)}}^2)$  and  $((\hat{\sigma}^*)_{e^{(a)}}^2, (\hat{\sigma}^*)_{e^{(b)}}^2)$  have the same limits. The conclusion follows from Theorem S1. □

### Proofs of Lemmas

In this section, we will drop the superscript on  $\mathbf{h}$ ,  $e$  and  $\hat{\beta}$  and focus on the proof for treatment group A, as the same analysis can be applied to control group B.

#### Proof of Lemma S2.

*Proof.* Let  $c_n = \frac{(1+\tau)L^{1/4}}{p_A} \sqrt{\frac{2 \log p}{n}}$ . By the union bound,

$$P(\|\bar{\mathbf{x}}_A\|_\infty > c_n) = P\left(\max_{j=1, \dots, p} \left| \frac{1}{n_A} \sum_{i \in A} x_{ij} \right| > c_n\right) \leq \sum_{j=1}^p P\left(\left| \frac{1}{n_A} \sum_{i \in A} x_{ij} \right| > c_n\right). \quad (\text{S41})$$

By Cauchy-Schwarz inequality, we have

$$\frac{1}{n} \sum_{i=1}^n x_{ij}^2 \leq \left(\frac{1}{n} \sum_{i=1}^n x_{ij}^4\right)^{\frac{1}{2}} \left(\frac{1}{n} \sum_{i=1}^n 1^2\right)^{\frac{1}{2}} \leq \sqrt{L}. \quad (\text{S42})$$

Substituting the concentration inequality (S13) into (S41),

$$P(\|\bar{\mathbf{x}}_A\|_\infty > c_n) \leq 2 \exp\left\{\log p - \frac{p_A n_A c_n^2}{(1+\tau)^2 L^{1/2}}\right\} = 2 \exp\{-\log p\} \rightarrow 0.$$

□

#### Proof of Lemma S3.

*Proof.* We start with the KKT condition, which characterizes the solution to the Lasso. Recall the definition of the Lasso estimator  $\hat{\beta}$ :

$$\hat{\beta} = \arg \min_{\beta} \frac{1}{2n_A} \sum_{i \in A} \left(a_i - \bar{a}_A - (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T \beta\right)^2 + \lambda_a \|\beta\|_1.$$

The KKT condition for  $\hat{\beta}$  is

$$\frac{1}{n_A} \sum_{i \in A} (\mathbf{x}_i - \bar{\mathbf{x}}_A) \left(a_i - \bar{a}_A - (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T \hat{\beta}\right) = \lambda_a \kappa, \quad (\text{S43})$$

where  $\kappa$  is the subgradient of  $\|\beta\|_1$  taking value at  $\beta = \hat{\beta}$ , i.e.,

$$\kappa \in \partial \|\beta\|_1 \Big|_{\beta=\hat{\beta}} \quad \text{with} \quad \begin{cases} \kappa_j \in [-1, 1] \text{ for } j \text{ s.t. } \hat{\beta}_j = 0 \\ \kappa_j = \text{sign}(\hat{\beta}_j) \text{ otherwise} \end{cases} \quad (\text{S44})$$

Substituting  $a_i$  by the decomposition (3), (S43) becomes

$$\frac{1}{n_A} \sum_{i \in A} (\mathbf{x}_i - \bar{\mathbf{x}}_A) (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T (\beta - \hat{\beta}) + \frac{1}{n_A} \sum_{i \in A} (\mathbf{x}_i - \bar{\mathbf{x}}_A) (e_i - \bar{e}_A) = \lambda_a \kappa. \quad (\text{S45})$$

Multiplying both sides of (S45) by  $-\mathbf{h}^T = (\beta - \hat{\beta})^T$ , we have

$$\begin{aligned} & \frac{1}{n_A} \sum_{i \in A} \left( (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T \mathbf{h} \right)^2 - \mathbf{h}^T \frac{1}{n_A} \sum_{i \in A} (\mathbf{x}_i - \bar{\mathbf{x}}_A) (e_i - \bar{e}_A) \\ &= \lambda_a (\beta - \hat{\beta})^T \kappa \leq \lambda_a \left( \|\beta\|_1 - \|\hat{\beta}\|_1 \right) \end{aligned}$$

where the last inequality is because

$$\beta^T \kappa \leq \|\beta\|_1 \|\kappa\|_\infty \leq \|\beta\|_1 \quad \text{and} \quad \hat{\beta}^T \kappa = \|\hat{\beta}\|_1.$$

Rearranging and by Hölder’s inequality, we have

$$\begin{aligned} \frac{1}{n_A} \sum_{i \in A} \left( (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T \mathbf{h} \right)^2 &\leq \lambda_a \left( \|\beta\|_1 - \|\hat{\beta}\|_1 \right) + \mathbf{h}^T \frac{1}{n_A} \sum_{i \in A} (\mathbf{x}_i - \bar{\mathbf{x}}_A) (e_i - \bar{e}_A) \\ &\leq \lambda_a \left( \|\beta\|_1 - \|\hat{\beta}\|_1 \right) + \|\mathbf{h}\|_1 \underbrace{\left\| \frac{1}{n_A} \sum_{i \in A} (\mathbf{x}_i - \bar{\mathbf{x}}_A) (e_i - \bar{e}_A) \right\|_\infty}_{*} \end{aligned}$$

To control the term (\*), we define the event  $\mathcal{L} = \{* \leq \eta \lambda_a\}$ . The following Lemma S6 shows that, with  $\lambda_a$  defined appropriately,  $\mathcal{L}$  holds with probability approaching 1. We will prove this Lemma later.

**Lemma S6.** Define  $\mathcal{L} = \left\{ \left\| \frac{1}{n_A} \sum_{i \in A} (\mathbf{x}_i - \bar{\mathbf{x}}_A)(e_i - \bar{e}_A) \right\|_\infty \leq \eta \lambda_a \right\}$ . Then under the conditions of Theorem 1,  $P(\mathcal{L}) \rightarrow 1$ .

On  $\mathcal{L}$

$$\frac{1}{n_A} \sum_{i \in A} \left( (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T \mathbf{h} \right)^2 \leq \lambda_a \left( \|\boldsymbol{\beta}\|_1 - \|\hat{\boldsymbol{\beta}}\|_1 + \eta \|\mathbf{h}\|_1 \right). \quad (\text{S46})$$

By substituting the definition of  $\mathbf{h}$ , and several applications of the triangle inequality, we have

$$\|\boldsymbol{\beta}\|_1 - \|\hat{\boldsymbol{\beta}}\|_1 \leq \|\mathbf{h}_S\|_1 - \|\mathbf{h}_{S^c}\|_1 + 2 \|\boldsymbol{\beta}_{S^c}\|_1.$$

Therefore,

$$\begin{aligned} \frac{1}{n_A} \sum_{i \in A} \left( (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T \mathbf{h} \right)^2 &\leq \lambda_a \left( \|\mathbf{h}_S\|_1 - \|\mathbf{h}_{S^c}\|_1 + 2 \|\boldsymbol{\beta}_{S^c}\|_1 + \eta \|\mathbf{h}\|_1 \right) \\ &\leq \lambda_a \left( (\eta - 1) \|\mathbf{h}_{S^c}\|_1 + (1 + \eta) \|\mathbf{h}_S\|_1 + 2 \|\boldsymbol{\beta}_{S^c}\|_1 \right). \end{aligned}$$

Because  $\frac{1}{n_A} \sum_{i \in A} \left( (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T \mathbf{h} \right)^2 \geq 0$ , we obtain

$$(1 - \eta) \|\mathbf{h}_{S^c}\|_1 \leq (1 + \eta) \|\mathbf{h}_S\|_1 + 2 \|\boldsymbol{\beta}_{S^c}\|_1 \leq (1 + \eta) \|\mathbf{h}_S\|_1 + 2s\lambda_a. \quad (\text{S47})$$

where the last inequality is because of the definition of  $s$  in (S4) and  $S$  in (S8).

Consider the following two cases:

(I) If  $(1 + \eta) \|\mathbf{h}_S\|_1 + 2s\lambda_a \geq (1 - \eta)\xi \|\mathbf{h}_S\|_1$  then by (S47),

$$\|\mathbf{h}\|_1 = \|\mathbf{h}_S\|_1 + \|\mathbf{h}_{S^c}\|_1 \leq \left( \frac{1 + \eta}{1 - \eta} + 1 \right) \|\mathbf{h}_S\|_1 + \frac{2s\lambda_a}{1 - \eta} \leq \frac{2s\lambda_a}{1 - \eta} \left( \frac{2}{(1 - \eta)\xi - (1 + \eta)} + 1 \right).$$

By the definition of  $\lambda_a$  and the scaling assumptions (S5), (S6), we have that  $s\lambda_a = o\left(\frac{1}{\sqrt{\log p}}\right)$ .

(II) If  $(1 + \eta) \|\mathbf{h}_S\|_1 + 2s\lambda_a < (1 - \eta)\xi \|\mathbf{h}_S\|_1$  then by (S47) we have  $\|\mathbf{h}_{S^c}\|_1 \leq \xi \|\mathbf{h}_S\|_1$ . Applying the cone invertibility condition on the design matrix (S7),

$$\|\mathbf{h}\|_1 = \|\mathbf{h}_S\|_1 + \|\mathbf{h}_{S^c}\|_1 \leq (1 + \xi) \|\mathbf{h}_S\|_1 \leq (1 + \xi) C_s \left\| \frac{1}{n} X^T X \mathbf{h} \right\|_\infty \quad (\text{S48})$$

Before applying this inequality we will revisit the KKT condition (S44), but this time we will take the  $l_\infty$ -norm, yielding

$$\left\| \frac{1}{n_A} \sum_{i \in A} (\mathbf{x}_i - \bar{\mathbf{x}}_A)(\mathbf{x}_i - \bar{\mathbf{x}}_A)^T \mathbf{h} \right\|_\infty \leq \lambda_a + \left\| \frac{1}{n_A} \sum_{i \in A} (\mathbf{x}_i - \bar{\mathbf{x}}_A)(e_i - \bar{e}_A) \right\|_\infty \leq (1 + \eta) \lambda_a \quad (\text{S49})$$

where the latter inequality holds on the set  $\mathcal{L}$ . The final step is to control the deviation of the subsampled covariance matrix from the population covariance matrix, so that we can apply (S48). We define another event with constant  $C_1 = \frac{2(1+\tau)L^{1/2}}{p_A}$

$$\mathcal{M} = \left\{ \left\| \frac{1}{n_A} \sum_{i \in A} (\mathbf{x}_i - \bar{\mathbf{x}}_A)(\mathbf{x}_i - \bar{\mathbf{x}}_A)^T - \frac{1}{n} X^T X \right\|_\infty \leq C_1 \sqrt{\frac{\log p}{n}} \right\}$$

**Lemma S7.** Assume stability of treatment assignment probability condition 1 and moment condition 6 hold. Then  $P(\mathcal{M}) \rightarrow 1$ .

We will prove Lemma S7 later. Continuing our inequalities, on the event  $\mathcal{L} \cap \mathcal{M}$ ,

$$\begin{aligned} s \left\| \frac{1}{n} X^T X \mathbf{h} \right\|_\infty &\leq C_1 s \sqrt{\frac{\log p}{n}} \|\mathbf{h}\|_1 + s \left\| \frac{1}{n_A} \sum_{i \in A} (\mathbf{x}_i - \bar{\mathbf{x}}_A)(\mathbf{x}_i - \bar{\mathbf{x}}_A)^T \mathbf{h} \right\|_\infty \\ &\leq o(1) \|\mathbf{h}\|_1 + s(1 + \eta) \lambda_a, \end{aligned}$$

where we have applied the scaling assumption (S6) and (S49) in the second line. Hence, by (S48),

$$\|\mathbf{h}\|_1 \leq (1 + \xi) C \left[ o(1) \|\mathbf{h}\|_1 + s(1 + \eta) \lambda_a \right].$$

Again, applying the scaling assumptions (S5) and (S6), we get  $\|\mathbf{h}\|_1 = o_p\left(\frac{1}{\sqrt{\log p}}\right)$ .  $\square$

**Proof of Lemma S4.**

*Proof.* In the proof of Lemma S3, we have shown that, on  $\mathcal{L}$  defined in Lemma S6,

$$\begin{aligned} \frac{1}{n_A} \sum_{i \in A} \left( (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) \right)^2 &\leq \lambda_a \left( \|\boldsymbol{\beta}\|_1 - \|\hat{\boldsymbol{\beta}}\|_1 + \eta \|\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}\|_1 \right). \\ &\leq \lambda_a (1 + \eta) \|\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}\|_1. \end{aligned} \quad (\text{S50})$$

Let  $\mathbf{x}^j$  be the  $j$ -th column of the design matrix  $X$  and  $\bar{\mathbf{x}}_A^j = n_A^{-1} \sum_{i \in A} x_{ij}$ . Again, by KKT condition, we have

$$\left| \frac{1}{n_A} \sum_{i \in A} (x_{ij} - \bar{x}_A^j) \left( a_i - \bar{a}_A - (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T \hat{\boldsymbol{\beta}} \right) \right| = \lambda_a, \text{ if } \hat{\boldsymbol{\beta}}_j \neq 0.$$

Substituting  $a_i$  by the decomposition (3) yields

$$\left| \frac{1}{n_A} \sum_{i \in A} (x_{ij} - \bar{x}_A^j) (e_i - \bar{e}_A) + \frac{1}{n_A} \sum_{i \in A} (x_{ij} - \bar{x}_A^j) (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) \right| = \lambda_a.$$

Combining with the definition of the event  $\mathcal{L}$ , we have if  $\hat{\boldsymbol{\beta}}_j \neq 0$

$$\Delta_j := \left| \frac{1}{n_A} \sum_{i \in A} (x_{ij} - \bar{x}_A^j) (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) \right| \geq (1 - \eta) \lambda_a. \quad (\text{S51})$$

Let  $Z = (\mathbf{z}_1, \dots, \mathbf{z}_n) \in R^{p \times n}$  with  $\mathbf{z}_i = \mathbf{x}_i - \bar{\mathbf{x}}_A \in R^p$  and denote  $\mathbf{w} = Z^T (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})$ , then

$$\frac{1}{n_A} \|\mathbf{w}_A\|_2^2 = \frac{1}{n_A} \sum_{i \in A} \left( (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) \right)^2 \leq \lambda_a (1 + \eta) \|\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}\|_1.$$

Let  $Z_A = (\mathbf{z}_i : i \in A)$ , since the largest eigenvalues of  $Z_A^T Z_A$  and  $Z_A Z_A^T$  are the same,

$$\begin{aligned} \frac{1}{n_A} \mathbf{w}_A^T Z_A^T Z_A \mathbf{w}_A &\leq \frac{1}{n_A} \lambda_{\max}(Z_A^T Z_A) \|\mathbf{w}_A\|_2^2 \\ &\leq \frac{1}{n_A} \lambda_{\max}(Z_A Z_A^T) \lambda_a (\eta + 1) \|\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}\|_1 \\ &\leq \Lambda_{\max} \frac{n}{n_A} \lambda_a (1 + \eta) \|\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}\|_1. \end{aligned}$$

The last inequality holds because

$$\begin{aligned} \lambda_{\max}(Z_A Z_A^T) &= \max_{\mathbf{u}: \|\mathbf{u}\|_2=1} \mathbf{u}^T Z_A Z_A^T \mathbf{u} \\ &= \max_{\mathbf{u}: \|\mathbf{u}\|_2=1} \mathbf{u}^T \sum_{i \in A} (\mathbf{x}_i - \bar{\mathbf{x}}_A) (\mathbf{x}_i - \bar{\mathbf{x}}_A)^T \mathbf{u} \\ &= \max_{\mathbf{u}: \|\mathbf{u}\|_2=1} \mathbf{u}^T \sum_{i \in A} \mathbf{x}_i \mathbf{x}_i^T \mathbf{u} - n_A \mathbf{u}^T (\bar{\mathbf{x}}_A) (\bar{\mathbf{x}}_A)^T \mathbf{u} \\ &\leq \max_{\mathbf{u}: \|\mathbf{u}\|_2=1} \mathbf{u}^T \sum_{i \in A} \mathbf{x}_i \mathbf{x}_i^T \mathbf{u} \leq n \Lambda_{\max}. \end{aligned} \quad (\text{S52})$$

On the other hand,

$$\frac{1}{n_A} \mathbf{w}_A^T Z_A^T Z_A \mathbf{w}_A = \sum_{j=1}^p \Delta_j^2 \geq \sum_{j: \hat{\boldsymbol{\beta}}_j \neq 0} \Delta_j^2 \geq (1 - \eta)^2 \lambda_a^2 \hat{s}. \quad (\text{S53})$$

Combining (S51), (S53) and the fact that with probability going to 1 (see the proof of Lemma S3)

$$\|\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}\|_1 \leq Cs(1 + \eta)\lambda_a,$$

where  $C$  is a constant. We conclude that with probability going to 1

$$\hat{s} \leq \frac{1}{(1 - \eta)^2} \frac{1}{\lambda_a^2} \Lambda_{\max} \frac{n}{n_A} \lambda_a (1 + \eta) Cs(1 + \eta)\lambda_a \leq \frac{C(1 + \eta)^2}{p_A(1 - \eta)^2} s.$$

□

**Proof of Lemma S5.**

*Proof.* For any  $t > 0$ , we have

$$P(|\bar{z}_A - \lim_{n \rightarrow \infty} \bar{z}| > t) \leq P(|\bar{z}_A - \bar{z}| > t/2) + P(|\bar{z} - \lim_{n \rightarrow \infty} \bar{z}| > t/2). \quad (\text{S54})$$

The second term in the right hand side of (S54) obviously converges to 0 as  $n \rightarrow \infty$ . To bound the first term, we apply the concentration inequality (S13). By (S30), it is easy to show

$$\frac{1}{n} \sum_{i=1}^n z_i^2 = \frac{1}{n} \sum_{i=1}^n |z_i|^{1-\epsilon} |z_i|^{1+\epsilon} \leq (nL)^{\frac{1-\epsilon}{1+\epsilon}} \frac{1}{n} \sum_{i=1}^n |z_i|^{1+\epsilon} \leq L^{\frac{2}{1+\epsilon}} n^{\frac{1-\epsilon}{1+\epsilon}}.$$

Concentration inequality (S13) yields

$$P(|\bar{z}_A - \bar{z}| > t/2) \leq 2 \exp \left\{ -\frac{p_A n_A t^2}{4(1+\tau)^2 L^{\frac{2}{1+\epsilon}} n^{\frac{1-\epsilon}{1+\epsilon}}} \right\} \rightarrow 0.$$

□

**Proof of Lemma S6.**

*Proof.* It is easy to verify that

$$\frac{1}{n_A} \sum_{i \in A} (\mathbf{x}_i - \bar{\mathbf{x}}_A)(e_i - \bar{e}_A) = \frac{1}{n_A} \sum_{i \in A} \mathbf{x}_i e_i - (\bar{\mathbf{x}}_A)(\bar{e}_A).$$

Hence,

$$\left\| \frac{1}{n_A} \sum_{i \in A} (\mathbf{x}_i - \bar{\mathbf{x}}_A)(e_i - \bar{e}_A) \right\|_\infty \leq \left\| \frac{1}{n_A} \sum_{i \in A} \mathbf{x}_i e_i \right\|_\infty + \|(\bar{\mathbf{x}}_A)(\bar{e}_A)\|_\infty. \quad (\text{S55})$$

We analyze these two terms on the right hand side of the inequality separately. For the first term, by triangle inequality and the definition of  $\delta_n$  in (9),

$$\begin{aligned} \left\| \frac{1}{n_A} \sum_{i \in A} \mathbf{x}_i e_i \right\|_\infty &\leq \left\| \frac{1}{n_A} \sum_{i \in A} \mathbf{x}_i e_i - \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i e_i \right\|_\infty + \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i e_i \right\|_\infty \\ &\leq \left\| \frac{1}{n_A} \sum_{i \in A} \mathbf{x}_i e_i - \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i e_i \right\|_\infty + \delta_n. \end{aligned} \quad (\text{S56})$$

We will again bound (S56) by the concentration inequality (S13) in Lemma S1. By Cauchy-Schwarz inequality, we have for any  $j = 1, \dots, p$ ,

$$\frac{1}{n} \sum_{i=1}^n x_{ij}^2 e_i^2 \leq \left( \frac{1}{n} \sum_{i=1}^n x_{ij}^4 \right)^{\frac{1}{2}} \left( \frac{1}{n} \sum_{i=1}^n e_i^4 \right)^{\frac{1}{2}} \leq L.$$

Let  $t_n = \frac{(1+\tau)L^{1/2}}{p_A} \sqrt{\frac{2 \log p}{n}}$ , then by the union bound and the concentration inequality (S13),

$$\begin{aligned} P \left( \left\| \frac{1}{n_A} \sum_{i \in A} \mathbf{x}_i e_i - \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i e_i \right\|_\infty > t_n \right) &\leq 2 \exp \left\{ \log p - \frac{p_A n_A t_n^2}{(1+\tau)^2 L} \right\} \\ &= 2 \exp \{-\log p\} \rightarrow 0. \end{aligned}$$

Taking this back to (S56), we have

$$P \left( \left\| \frac{1}{n_A} \sum_{i \in A} \mathbf{x}_i e_i \right\|_\infty \leq t_n + \delta_n \right) \rightarrow 1. \quad (\text{S57})$$

For the second term, by Lemma S2, we have shown that,

$$P \left( \|\bar{\mathbf{x}}_A\|_\infty \leq \frac{(1+\tau)L^{1/4}}{p_A} \sqrt{\frac{2 \log p}{n}} \right) \rightarrow 1.$$

Similar proof yields

$$P \left( \|\bar{e}_A\|_\infty \leq \frac{(1+\tau)L^{1/4}}{p_A} \sqrt{\frac{2 \log p}{n}} \right) \rightarrow 1.$$

Hence, under the scaling condition (S6),

$$P \left( \|(\bar{\mathbf{x}}_A)(\bar{e}_A)\|_\infty \leq \frac{(1+\tau)L^{1/2}}{p_A} \sqrt{\frac{2 \log p}{n}} \right) \rightarrow 1. \quad (\text{S58})$$

Combining (S57) and (S58) yields

$$P \left( \left\| \frac{1}{n_A} \sum_{i \in A} (\mathbf{x}_i - \bar{\mathbf{x}}_A)(e_i - \bar{e}_A) \right\|_\infty \leq \frac{2(1+\tau)L^{1/2}}{p_A} \sqrt{\frac{2 \log p}{n}} + \delta_n \right) \rightarrow 1.$$

The conclusion follows from the condition  $\lambda_a \in (\frac{1}{\eta}, M] \times \left( \frac{2(1+\tau)L^{1/2}}{p_A} \sqrt{\frac{2 \log p}{n}} + \delta_n \right)$ . □

**Proof of Lemma S7.**

*Proof.* It is easy to see that

$$\frac{1}{n_A} \sum_{i \in A} (\mathbf{x}_i - \bar{\mathbf{x}}_A)(\mathbf{x}_i - \bar{\mathbf{x}}_A)^T = \frac{1}{n_A} \sum_{i \in A} \mathbf{x}_i \mathbf{x}_i^T - (\bar{\mathbf{x}}_A)(\bar{\mathbf{x}}_A)^T.$$

Then, by triangle inequality,

$$\left\| \frac{1}{n_A} \sum_{i \in A} (\mathbf{x}_i - \bar{\mathbf{x}}_A)(\mathbf{x}_i - \bar{\mathbf{x}}_A)^T - \frac{1}{n} X^T X \right\|_\infty \tag{S59}$$

$$\leq \underbrace{\left\| \frac{1}{n_A} \sum_{i \in A} \mathbf{x}_i \mathbf{x}_i^T - \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \right\|_\infty}_* + \underbrace{\|(\bar{\mathbf{x}}_A)(\bar{\mathbf{x}}_A)^T\|_\infty}_{**}. \tag{S60}$$

We control the first term (\*) again using the concentration inequality (S13) and the union bound. By Cauchy-Schwarz inequality, for  $j, k = 1, \dots, p$ ,

$$\frac{1}{n} \sum_{i=1}^n x_{ij}^2 x_{ik}^2 \leq \left( \frac{1}{n} \sum_{i=1}^n x_{ij}^4 \right)^{\frac{1}{2}} \left( \frac{1}{n} \sum_{i=1}^n x_{ik}^4 \right)^{\frac{1}{2}} \leq L.$$

Then,

$$\begin{aligned} & P \left( \left\| \frac{1}{n_A} \sum_{i \in A} \mathbf{x}_i \mathbf{x}_i^T - \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \right\|_\infty \geq \frac{(1+\tau)L^{1/2}}{p_A} \sqrt{\frac{3 \log p}{n}} \right) \\ & \leq 2 \exp \left\{ 2 \log p - \frac{3p_A n_A (1+\tau)^2 L \log p}{(1+\tau)^2 L p_A^2 n} \right\} = 2 \exp \{-\log p\} \rightarrow 0. \end{aligned} \tag{S61}$$

The second term (\*\*) is bounded by again observing that, by Lemma S2 and the scaling condition (S6),

$$(**) \leq \|\bar{\mathbf{x}}_A\|_\infty^2 = o_p \left( \sqrt{\frac{\log p}{n}} \right). \tag{S62}$$

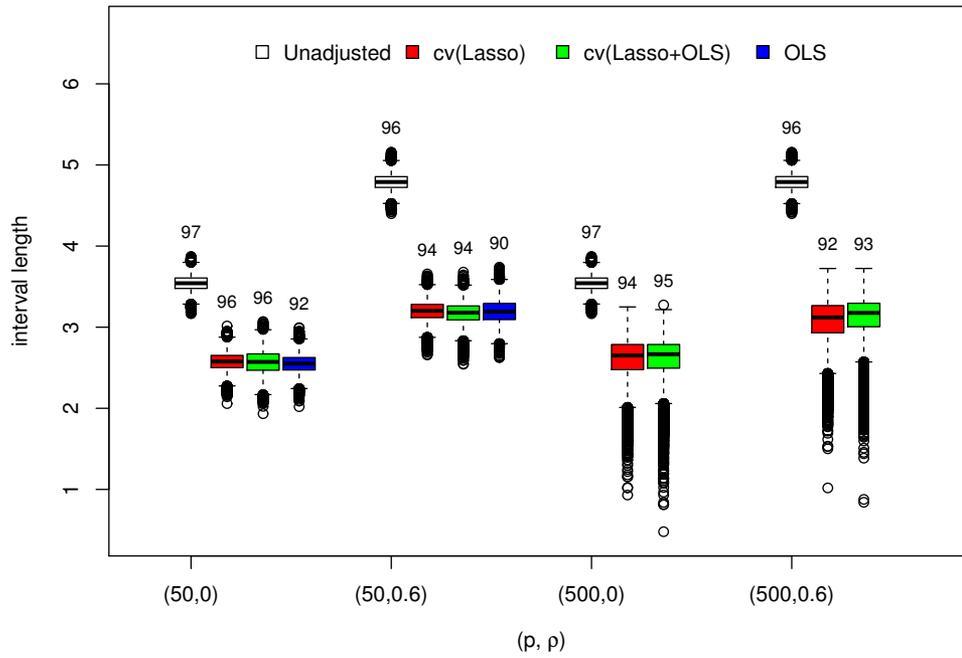
Combining (S61) and (S62) yields the conclusion. □

**Tables and Figures**

**References**

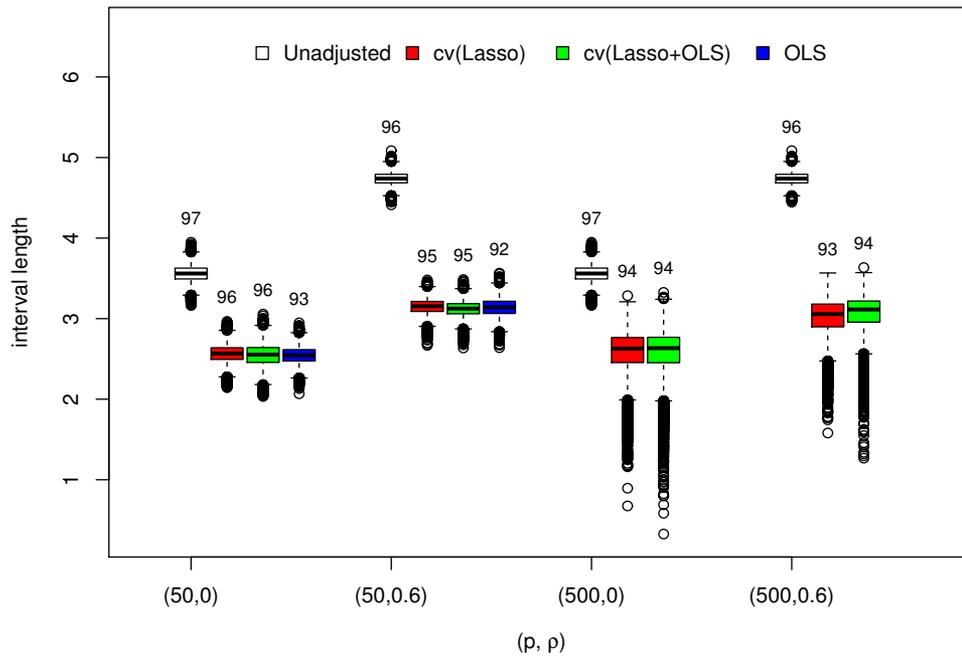
1. Meinshausen N (2007) Relaxed Lasso. *Computational Statistics and Data Analysis* 52(1):374-393.
2. Efron B, Hastie T and Tibshirani R (2004) Least angle regression. *Annals of Statistics* 32(2):407-499.
3. Belloni A, Chernozhukov V (2009) Least Squares After Model Selection in High-dimensional Sparse Models. *Bernoulli* 19(2):521-547.
4. Liu H, Yu B (2013) Asymptotic properties of Lasso+mLS and Lasso+Ridge in sparse high-dimensional linear regression. *Electronic Journal of Statistics* 7(2013):3124-3169.
5. Massart P (1986) in *Geometrical and Statistical Aspects of Probability in Banach Spaces*. (Springer), pp. 73-109.
6. Freedman DA (2008) On regression adjustments in experiments with several treatments. *The Annals of Applied Statistics* 2(1):176-196.

Boxplot of interval length (95% confidence interval) with coverage probability on top ( $n_A=100$ )



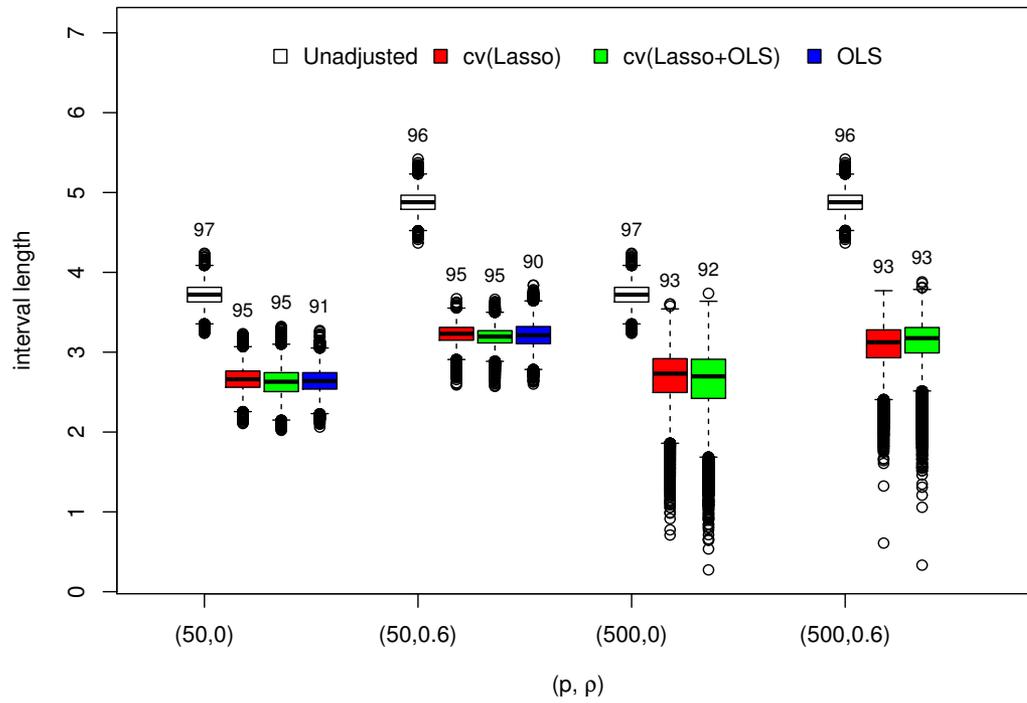
**Fig. S1.** Boxplot of the interval length with coverage probability (%) on top of each box for the unadjusted, OLS adjusted (only computed when  $p = 50$ ), cv(Lasso) adjusted and cv(Lasso+OLS) adjusted estimators with  $n_A = 100$ .

Boxplot of interval length (95% confidence interval) with coverage probability on top ( $n_A=125$ )



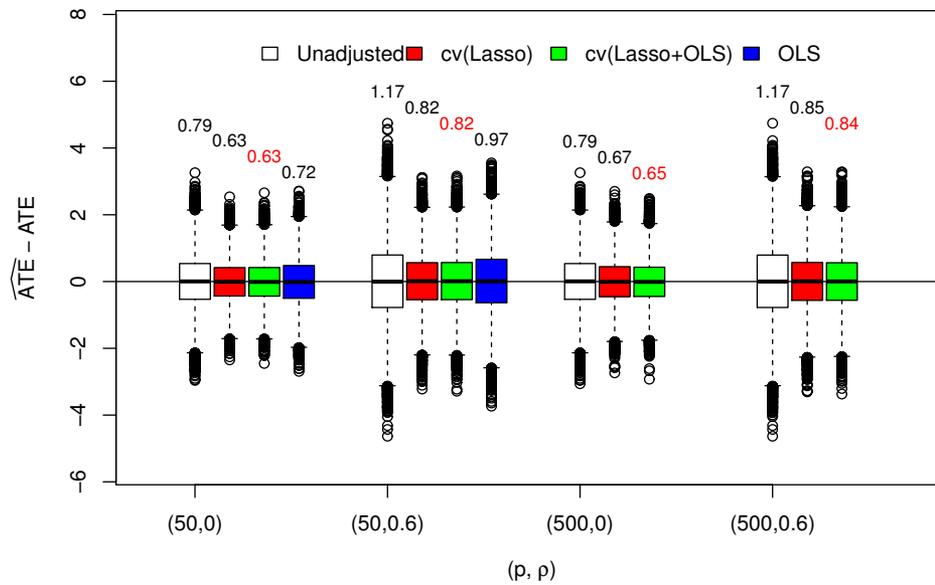
**Fig. S2.** Boxplot of the interval length with coverage probability (%) on top of each box for the unadjusted, OLS adjusted (only computed when  $p = 50$ ), cv(Lasso) adjusted and cv(Lasso+OLS) adjusted estimators with  $n_A = 125$ .

Boxplot of interval length (95% confidence interval) with coverage probability on top ( $n_A=150$ )

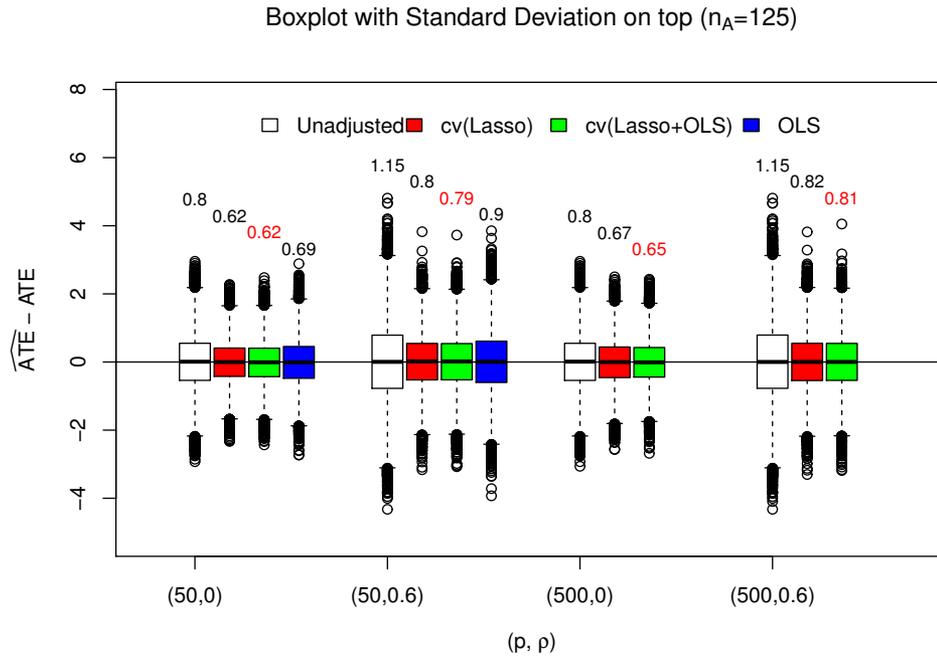


**Fig. S3.** Boxplot of the interval length with coverage probability (%) on top of each box for the unadjusted, OLS adjusted (only computed when  $p = 50$ ), cv(Lasso) adjusted and cv(Lasso+OLS) adjusted estimators with  $n_A = 150$ .

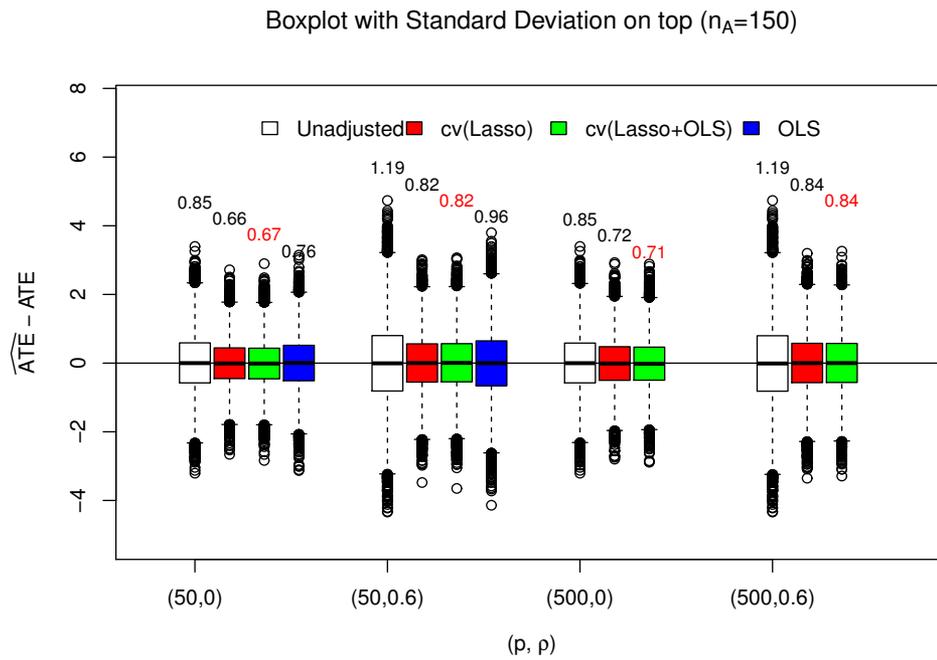
Boxplot with Standard Deviation on top ( $n_A=100$ )



**Fig. S4.** Boxplot of the unadjusted, OLS adjusted (only computed when  $p = 50$ ), cv(Lasso) and cv(Lasso+OLS) adjusted estimators with their standard deviations presented on top of each box for  $n_A = 100$ .



**Fig. S5.** Boxplot of the unadjusted, OLS adjusted (only computed when  $p = 50$ ),  $cv(Lasso)$  and  $cv(Lasso+OLS)$  adjusted estimators with their standard deviations presented on top of each box for  $n_A = 125$ .



**Fig. S6.** Boxplot of the unadjusted, OLS adjusted (only computed when  $p = 50$ ),  $cv(Lasso)$  and  $cv(Lasso+OLS)$  adjusted estimators with their standard deviations presented on top of each box for  $n_A = 150$ .

---

**Algorithm 1**  $K$ -fold Cross Validation (CV) for the Lasso+OLS estimator

---

**Input:** Design matrix  $X$ , response  $Y$  and a sequence of tuning parameter  $\lambda_1, \dots, \lambda_J$ ; Number of folds  $K$ .

**Output:** The optimal tuning parameter selected by CV:  $\lambda_{optimal}$ .

- 1: Divide randomly the data  $z = (X, Y)$  into  $K$  roughly equal parts  $z_k, k = 1, \dots, K$ ;
- 2: For each  $k = 1, \dots, K$ , denote  $\hat{S}^{(k)}(\lambda_0) = \emptyset$  and  $\hat{\beta}_{\text{Lasso+OLS}}^{(k)}(\lambda_0) = 0$ .

- Fit the model with parameters  $\lambda_j, j = 1, \dots, J$  to the other  $K - 1$  parts  $z_{-k} = z \setminus z_k$  of the data, giving the Lasso solution path  $\hat{\beta}^{(k)}(\lambda_j), j = 1, \dots, J$  and compute the selected covariates set  $\hat{S}^{(k)}(\lambda_j) = \{l : \hat{\beta}_l^{(k)}(\lambda_j) \neq 0\}, j = 1, \dots, J$  on the path;
- For each  $j = 1, \dots, J$ , compute the Lasso+OLS estimator:

$$\hat{\beta}_{\text{Lasso+OLS}}^{(k)}(\lambda_j) = \begin{cases} \arg \min_{\beta: \beta_j=0, \forall j \notin \hat{S}^{(k)}(\lambda_j)} \left\{ \frac{1}{2|z_{-k}|} \sum_{i \in z_{-k}} (y_i - x_i^T \beta)^2 \right\}, & \text{if } \hat{S}^{(k)}(\lambda_j) \neq \hat{S}^{(k)}(\lambda_{j-1}), \\ \hat{\beta}_{\text{Lasso+OLS}}^{(k)}(\lambda_{j-1}), & \text{otherwise;} \end{cases} \quad (\text{S63})$$

- Compute the error in predicting the  $k$ th part of the data  $PE^{(k)}$  :

$$PE^{(k)}(\lambda_j) = \frac{1}{|z_k|} \sum_{i \in z_k} \left( y_i - x_i^T \hat{\beta}_{\text{Lasso+OLS}}^{(k)}(\lambda_j) \right)^2 ;$$

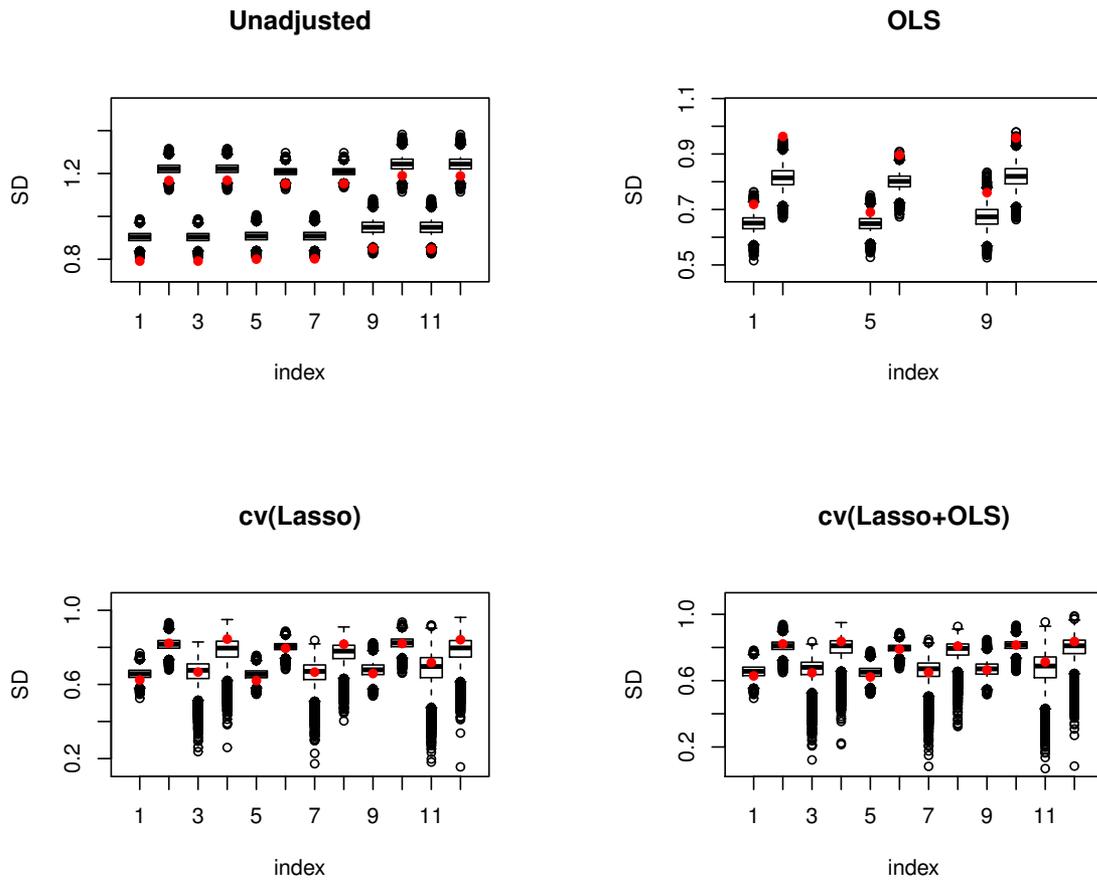
- 3: Compute cross validation error  $CV(\lambda_j), j = 1, \dots, J$ :

$$CV(\lambda_j) = \frac{1}{K} \sum_{k=1}^K PE^{(k)}(\lambda_j);$$

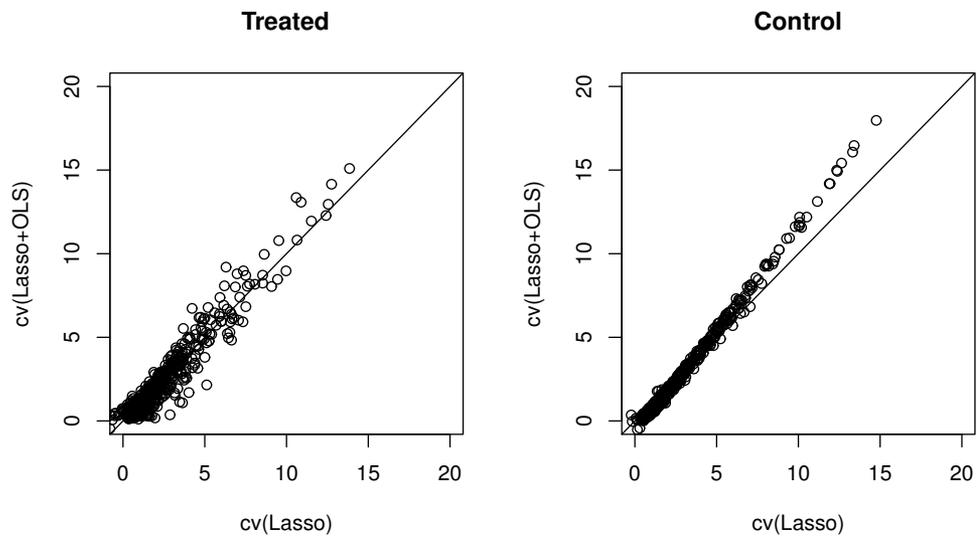
- 4: Compute the optimal  $\lambda$  selected by CV;

$$\lambda_{optimal} = \underset{\lambda_j: j=1, \dots, J}{\operatorname{argmin}} CV(\lambda_j);$$

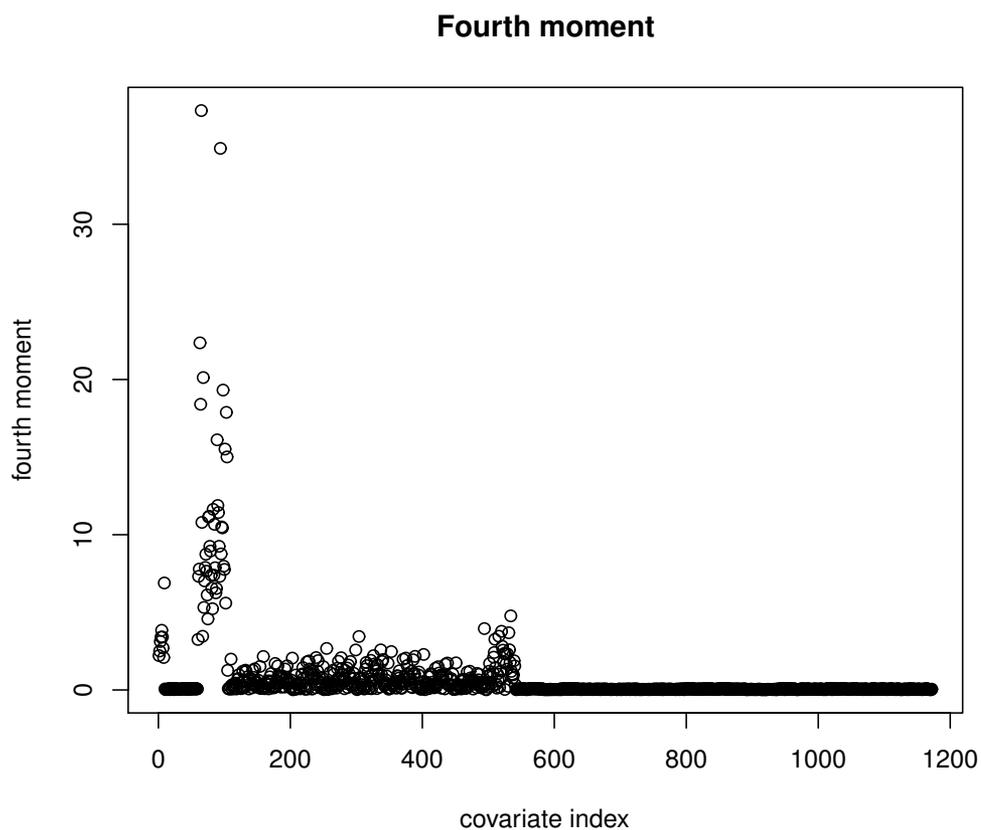
- 5: **return**  $\lambda_{optimal}$ .
-



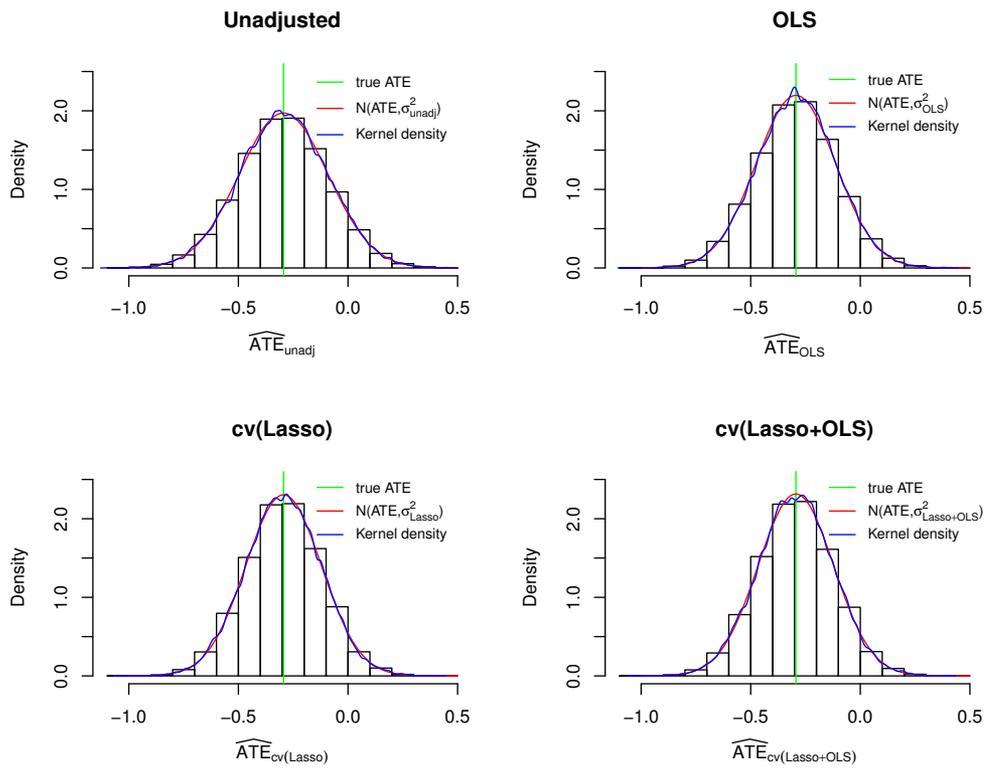
**Fig. S7.** Boxplot of Neyman SD estimate with the “true” SD presented as red dot.



**Fig. S8.** Adjustment (fitted) value comparison for cv(Lasso) and cv(Lasso+OLS).



**Fig. S9.** Fourth moment of each covariate. The covariates with the largest two fourth moments (37.3 and 34.9 respectively) are quadratic term  $interactnew^2$  and interaction term  $IMscorecct : systemnew$  respectively. Neither of them are selected by the Lasso to do the adjustment. All the fourth moments of the main effects are less than 7.



**Fig. S10.** Histograms of ATE estimates. The green vertical lines are the true ATE; the red curves are the density of normal distribution; the blue curves are the kernel density estimate. The blue curves are very close to the red ones meaning that all the ATE estimates follow normal distribution.

**Table S1. Bias, standard deviation (SD) and root-mean square error  $\sqrt{\text{MSE}}$  of ATE estimates**

Statistic	Method	$(p, \rho)$			
		(50,0)	(50,0.6)	(500,0)	(500,0.6)
$n_A = 100$					
bias	Unadjusted	<b>0.003(0.004)*</b>	<b>0.005(0.005)</b>	<b>0.002(0.003)</b>	<b>0.003(0.005)</b>
	OLS	0.014(0.005)	0.013(0.006)	-	-
	cv(Lasso)	0.007(0.004)	0.014(0.005)	0.006(0.004)	0.005(0.004)
	cv(Lasso+OLS)	0.011(0.004)	0.013(0.005)	0.009(0.004)	<b>0.003(0.004)</b>
SD	Unadjusted	0.79(0.08)	1.17(0.11)	0.79(0.07)	1.17(0.11)
	OLS	0.72(0.07)	0.96(0.09)	-	-
	cv(Lasso)	<b>0.62(0.06)</b>	<b>0.82(0.08)</b>	0.67(0.06)	<b>0.84(0.08)</b>
	cv(Lasso+OLS)	0.63(0.06)	<b>0.82(0.08)</b>	<b>0.65(0.06)</b>	<b>0.84(0.08)</b>
$\sqrt{\text{MSE}}$	Unadjusted	0.79(0.08)	1.17(0.11)	0.79(0.07)	1.17(0.11)
	OLS	0.72(0.07)	0.97(0.09)	-	-
	cv(Lasso)	<b>0.63(0.06)</b>	<b>0.82(0.08)</b>	0.67(0.06)	0.85(0.08)
	cv(Lasso+OLS)	<b>0.63(0.06)</b>	<b>0.82(0.08)</b>	<b>0.65(0.06)</b>	<b>0.84(0.08)</b>
$n_A = 125$					
bias	Unadjusted	0.008(0.005)	0.011(0.007)	<b>0.006(0.004)</b>	0.01(0.007)
	OLS	0.008(0.004)	<b>0.005(0.005)</b>	-	-
	cv(Lasso)	<b>0.005(0.003)</b>	0.012(0.005)	0.007(0.004)	0.004(0.004)
	cv(Lasso+OLS)	0.012(0.004)	0.012(0.005)	0.011(0.004)	<b>0.003(0.003)</b>
SD	Unadjusted	0.80(0.08)	1.15(0.11)	0.8(0.08)	1.15(0.11)
	OLS	0.69(0.06)	0.90(0.09)	-	-
	cv(Lasso)	<b>0.62(0.06)</b>	<b>0.79(0.07)</b>	0.67(0.06)	0.82(0.08)
	cv(Lasso+OLS)	<b>0.62(0.06)</b>	<b>0.79(0.07)</b>	<b>0.65(0.06)</b>	<b>0.81(0.08)</b>
$\sqrt{\text{MSE}}$	Unadjusted	0.80(0.07)	1.15(0.11)	0.8(0.07)	1.15(0.11)
	OLS	0.69(0.07)	0.90(0.09)	-	-
	cv(Lasso)	<b>0.62(0.06)</b>	0.80(0.08)	0.67(0.06)	0.82(0.08)
	cv(Lasso+OLS)	<b>0.62(0.06)</b>	<b>0.79(0.07)</b>	<b>0.65(0.06)</b>	<b>0.81(0.08)</b>
$n_A = 150$					
bias	Unadjusted	0.004(0.004)	<b>0.000(0.005)</b>	<b>0.002(0.003)</b>	0.005(0.005)
	OLS	<b>0.002(0.003)</b>	0.006(0.005)	-	-
	cv(Lasso)	0.003(0.003)	0.002(0.004)	0.01(0.005)	0.002(0.003)
	cv(Lasso+OLS)	0.011(0.004)	0.006(0.004)	0.017(0.005)	<b>0.001(0.003)</b>
SD	Unadjusted	0.85(0.08)	1.19(0.11)	0.85(0.08)	1.19(0.11)
	OLS	0.76(0.07)	0.96(0.09)	-	-
	cv(Lasso)	<b>0.66(0.06)</b>	0.82(0.08)	0.72(0.07)	<b>0.84(0.08)</b>
	cv(Lasso+OLS)	0.67(0.06)	<b>0.81(0.07)</b>	<b>0.71(0.07)</b>	<b>0.84(0.08)</b>
$\sqrt{\text{MSE}}$	Unadjusted	0.85(0.08)	1.19(0.11)	0.85(0.08)	1.19(0.11)
	OLS	0.76(0.07)	0.96(0.09)	-	-
	cv(Lasso)	<b>0.66(0.06)</b>	<b>0.82(0.08)</b>	0.72(0.07)	<b>0.84(0.08)</b>
	cv(Lasso+OLS)	0.67(0.06)	<b>0.82(0.08)</b>	<b>0.71(0.07)</b>	<b>0.84(0.08)</b>

\*The numbers in parentheses are the corresponding standard errors estimated by using the bootstrap with  $B = 500$  resamplings of the ATE estimates.

**Table S2. Mean number of selected covariates for treated and control group**

Group	Method	$(p, \rho)$			
		(50,0)	(50,0.6)	(500,0)	(500,0.6)
$n_A = 100$					
treated	cv(Lasso)	16	13	22	22
	cv(Lasso+OLS)	6	6	7	7
control	cv(Lasso)	20	11	32	28
	cv(Lasso+OLS)	8	6	7	7
$n_A = 125$					
treated	cv(Lasso)	17	13	25	24
	cv(Lasso+OLS)	7	6	6	6
control	cv(Lasso)	19	11	32	27
	cv(Lasso+OLS)	8	6	9	8
$n_A = 150$					
treated	cv(Lasso)	18	13	29	26
	cv(Lasso+OLS)	8	7	6	6
control	cv(Lasso)	19	12	30	25
	cv(Lasso+OLS)	8	6	11	8

**Table S3. Coverage probability (%) and mean interval length (in parentheses) for 95% confidence interval**

Methods	$(p, \rho)$			
	(50,0)	(50,0.6)	(500,0)	(500,0.6)
$n_A = 100$				
Unadjusted	97.3(3.54)*	95.8(4.79)	97.3(3.54)	95.8(4.79)
OLS	92.2(2.55)	90.0(3.19)	-	-
cv(Lasso)	95.8(2.58)	94.5(3.20)	94.3(2.61)	92.4(3.07)
cv(Lasso+OLS)	95.6(2.57)	94.4(3.17)	94.8(2.60)	93.0(3.11)
$n_A = 125$				
Unadjusted	97.4(3.56)	96.0(4.74)	97.3(3.56)	95.9(4.74)
OLS	93.3(2.54)	91.6(3.14)	-	-
cv(Lasso)	96.0(2.56)	95.0(3.15)	94.1(2.59)	92.9(3.02)
cv(Lasso+OLS)	95.7(2.55)	94.9(3.12)	94.4(2.58)	93.6(3.06)
$n_A = 150$				
Unadjusted	97.1(3.72)	95.8(4.88)	97.1(3.72)	95.8(4.88)
OLS	91.4(2.64)	90.4(3.21)	-	-
cv(Lasso)	95.4(2.66)	94.9(3.23)	92.9(2.68)	92.6(3.08)
cv(Lasso+OLS)	94.7(2.63)	94.8(3.19)	92.0(2.63)	93.1(3.11)

\*The numbers in parentheses are the corresponding mean interval lengths.

**Table S4. Statistics for the PAC illustration**

Methods	$\widehat{ATE}$	$\hat{\sigma}_{ATE}$	95% confidence interval	No. of selected covariates	
				treated	control
Unadjusted	-0.13	0.081	[-0.69,0.43]	-	-
OLS	-0.31	0.054	[-0.77,0.14]	-	-
cv(Lasso)	-0.33	0.052	[-0.77,0.12]	24	8
cv(Lasso+OLS)	-0.36	0.053	[-0.82,0.09]	4	5

**Table S5. Statistics for the PAC synthetic data set**

	Bias	SD	$\sqrt{\text{MSE}}$	Coverage (%)	Length	No. of selected covariates	
						treated	control
unadjusted	0.001(0)*	0.20(0.02)	0.20(0.02)	99	1.06	-	-
OLS	0.002(0)	0.18(0.02)	0.18(0.02)	99	0.95	-	-
cv(Lasso)	0.001(0)	<b>0.17(0.02)</b>	<b>0.17(0.02)</b>	99	<b>0.94</b>	25(23)	15(14)
cv(Lasso+OLS)	<b>0.000(0)</b>	<b>0.17(0.02)</b>	<b>0.17(0.02)</b>	99	0.95	<b>6(6)</b>	<b>4(3)</b>

\*The numbers in parentheses are the corresponding standard errors estimated by using the bootstrap with  $B = 500$  resamplings of the ATE estimates.