

Treatment Effect Estimation via Differentiated Confounder Balancing and Regression

KUN KUANG, Zhejiang University, Tsinghua University

PENG CUI and BO LI, Tsinghua University

MENG JIANG, University of Notre Dame

YASHEN WANG, China Academy of Electronics and Information Technology

FEI WU, Zhejiang University

SHIQIANG YANG, Tsinghua University

Treatment effect plays an important role on decision making in many fields, such as social marketing, healthcare, and public policy. The key challenge on estimating treatment effect in the wild observational studies is to handle confounding bias induced by imbalance of the confounder distributions between treated and control units. Traditional methods remove confounding bias by re-weighting units with supposedly accurate propensity score estimation under the unconfoundedness assumption. Controlling high-dimensional variables may make the unconfoundedness assumption more plausible, but poses new challenge on accurate propensity score estimation. One strand of recent literature seeks to directly optimize weights to balance confounder distributions, bypassing propensity score estimation. But existing balancing methods fail to do selection and differentiation among the pool of a large number of potential confounders, leading to possible underperformance in many high-dimensional settings. In this article, we propose a data-driven Differentiated Confounder Balancing (DCB) algorithm to jointly select confounders, differentiate weights of confounders and balance confounder distributions for treatment effect estimation in the wild high-dimensional settings. Besides, under some settings with heavy confounding bias, in order to further reduce the bias and variance of estimated treatment effect, we propose a Regression Adjusted Differentiated Confounder Balancing (RA-DCB) algorithm based on our DCB algorithm by incorporating outcome regression adjustment. The synergistic learning algorithms we proposed are more capable of reducing the confounding bias in many observational

This work was supported in part by National Program on Key Basic Research Project No. 2018AAA0102004, 2015CB352300, National Key Research and Development Project No. 2017YFC0820503, National Natural Science Foundation of China Major Project No. U1611461; National Natural Science Foundation of China No. 61772304, 61521002, 61531006, U1611461; Beijing Academy of Artificial Intelligence, BAAI. This work was also supported by the research fund of Tsinghua-Tencent Joint Laboratory for Internet Innovation Technology, the Young Elite Scientist Sponsorship Program by CAST, and the Director Foundation Project of National Engineering Laboratory for Public Safety Risk Perception and Control by Big Data (PSRPC). Bo Li's research was supported by the Tsinghua University Initiative Scientific Research Grant, No. 20165080091; National Natural Science Foundation of China, No. 71490723 and No. 71432004; Science Foundation of Ministry of Education of China, No. 16JJD630006. Fei Wu's research was supported by zhejiang lab.

Authors' addresses: K. Kuang and F. Wu, College of Computer Science and Technology, Zhejiang University, Zhejiang, China; emails: kunkuang@zju.edu.cn, wufei@cs.zju.edu.cn; P. Cui (corresponding author) and S. Yang, Department of Computer Science and Technology, Tsinghua University, Beijing, China; emails: {cui, yangshq}@tsinghua.edu.cn; B. Li (corresponding author), School of Economics and Management, Tsinghua University, Beijing, China; email: libo@sem.tsinghua.edu.cn; M. Jiang, Department of Computer Science and Engineering, University of Notre Dame, Notre Dame, IN; email: mjiang2@nd.edu; Y. Wang, National Engineering Laboratory for Public Safety Risk Perception and Control by Big Data (PSRPC), China Academy of Electronics and Information Technology, Beijing, China; email: yashen_wang@126.com. Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2019 Association for Computing Machinery.

1556-4681/2019/12-ART6 \$15.00

<https://doi.org/10.1145/3365677>

studies. To validate the effectiveness of our DCB and RA-DCB algorithms, we conduct extensive experiments on both synthetic and real-world datasets. The experimental results clearly demonstrate that our algorithms outperform the state-of-the-art methods. By incorporating regression adjustment, our RA-DCB algorithm achieves more precise estimation on treatment effect than DCB algorithm, especially under the settings with heavy confounding bias. Moreover, we show that the top features ranked by our algorithm generate accurate prediction of online advertising effect.

CCS Concepts: • **Computing methodologies** → **Causal reasoning and diagnostics**; *Machine learning*; *Statistical relational learning*;

Additional Key Words and Phrases: Treatment effect estimation, confounding bias, differentiated confounder balancing, regression adjustment

ACM Reference format:

Kun Kuang, Peng Cui, Bo Li, Meng Jiang, Yashen Wang, Fei Wu, and Shiqiang Yang. 2019. Treatment Effect Estimation via Differentiated Confounder Balancing and Regression. *ACM Trans. Knowl. Discov. Data* 14, 1, Article 6 (December 2019), 25 pages.
<https://doi.org/10.1145/3365677>

1 INTRODUCTION

Owing to the popularity of Big Data, abundant data are accumulated in various domains such as healthcare and advertising. At the same time, many machine learning and data mining methods are proposed to exploit these data for prediction, aiming to estimate the future outcome in the application of interest. These methods have been proved to be successful in prediction-oriented applications. However, the lack of interpretability of most predictive algorithms makes them less attractive in many settings, especially those requiring decision making, such as healthcare and policy making. How to improve the explainability of learning algorithms is of paramount importance for both academic research and real applications.

Causal inference, which refers to the process of drawing a conclusion about a causal connection based on the conditions of the occurrence of an effect [16], is a powerful statistical modeling tool for explanatory analysis. One fundamental problem in causal inference is treatment effect estimation, and its key challenge is to remove the confounding bias induced by the different confounder distributions between treated and control units. The gold standard approach for removing confounding bias is to conduct randomized experiments like A/B testing [27], where different treatments are randomly assigned to units.¹ But fully randomized experiments are usually time-consuming, expensive [21] and sometimes infeasible [6]. Therefore, many methods are proposed to estimate treatment effect directly from observational data under the unconfoundedness assumption [38]. Most of them adopt the propensity score to reweight units for removing confounding bias [3, 4, 8, 23]. Although these methods are gaining ground in applied work, they require correct model specification on treatment assignment or accurate propensity score estimation. In big data scenarios, controlling high-dimensional variables may make the unconfoundedness assumption more plausible, but poses new challenge on accurate propensity score estimation. Recently, some researchers proposed to balance confounder distributions by directly optimizing the weights, without modeling or estimating the propensity scores [2, 9, 15, 46]. But they balance all observed variables equally without screening and differentiation of confounders, leading to poor performance in high-dimensional settings. Overall, the previous methods can work well in well-designed experimental settings or observational studies with grounded model assumptions and prior knowledge.

¹Units represent the objects of treatment. For example, the units refer to the users in online advertising campaign.

In the wild big data scenarios, however, there are almost always a large number of additional or mostly uncontrolled confounders and identified variables, and the correlations among them are complex and unknown in the real world [41]. Hence, we face the following challenges in estimating treatment effect in the wild observational studies: (1) *unknown model structure of the interactions among variables*: As stated in [41], pretty much everything in the real world interacts with everything else, to some degree, and their interactions are complicated due to the complex nature of the real world. We hardly know the real model structure among variables in the wild, so we cannot make any model specification *a priori* for removing confounding bias. (2) *High-dimensional and noisy variables*: In big data scenario, there are always a large number of observed variables, but not all these variables are confounders and different confounders contribute unequally to the confounding bias in data. Usually, we do not have sufficient prior knowledge to justify the inclusion of hundreds or even thousands of variables. How to differentiate the confounders and their confounding bias is quite difficult.

To address these challenges, we propose a data-driven method, named Differentiated Confounder Balancing (DCB) algorithm. The method is based on the framework of confounder balancing, but in contrast with previous methods that balance all variables equally, we argue that some variables should not be regarded as confounders and we theoretically prove that the weights of confounders should be differentiated in confounder balancing. Motivated by this, we propose an integrated regularization algorithm to jointly select confounders, differentiate weights of confounders and balance confounder distributions for treatment effect estimation. During the treatment effect estimation, the selected confounders and their weights are used to adjust the weights of units, so that the confounder distributions, approximated by their moments, over all units can be balanced in treated and control groups. We find, however, that our DCB algorithm could still have some substantial bias in the settings with heavy confounding bias. To address this problem, we propose a Regression Adjusted DCB (RA-DCB) model based on our DCB algorithm by incorporating regression adjustment on the outcome, aiming to further reduce the bias and variance of estimated treatment effect. We validate our DCB and RA-DCB algorithms with extensive experiments on both synthetic and real datasets. The results clearly demonstrate that our algorithms outperform the state-of-the-art methods on treatment effect estimation in observational studies. And we find that with considering regression adjustment, our RA-DCB algorithm achieves a better performance on treatment effect estimation than DCB algorithm, especially in high-dimensional settings and high bias selection settings.

The main contributions of this article are as follows:

- We address the new challenges of estimating treatment effect in big data scenarios with high-dimensional noisy variables and insufficient prior knowledge on variable interactions, which is beyond the capability of previous methods.
- We propose a novel DCB algorithm to jointly select confounders, optimize the confounder weights and sample weights for confounder balancing, and simultaneously estimate the treatment effect in observational studies.
- In order to further reduce the bias and variance of estimated treatment effect, we propose a new RA-DCB algorithm based on our DCB algorithm by incorporating regression adjustment on the outcome.
- The advantages of our DCB and RA-DCB algorithms are demonstrated in both synthetic and real datasets. We also show that our method can significantly help to improve the prediction performance with real online advertising dataset.

The rest of this article is organized as follows. Section 2 reviews the related work. Section 3 introduces our DCB estimator. Section 4 proposes the algorithm that accurately infers the treatment

effect, and gives some analysis on our model. Section 5 gives the experimental results. Finally, Section 6 concludes the article.

2 RELATED WORK

In this section, we review related fields including weighting based estimators, confounder selection methods and learning for causal inference.

Weighting based estimators: Existing weighting based treatment effect estimation methods in observational studies either employ propensity score or optimize balance weights directly.

The propensity score was first proposed by Rosenbaum and Rubin [38], where it was estimated via a logistic regression. Then many other machine learning algorithms (e.g., lasso [10, 12], boosting regression [30], bagged CART, and neural network [44]) are employed for propensity score estimation. Various estimators have been proposed based on propensity score, such as propensity score matching, inverse propensity weighting, and double robust estimators [3, 4, 8, 24, 25]. Recently, some novel methods [17, 23] have been proposed to improve the performance of propensity score based methods. Imai et al. [17] introduced a covariate balancing propensity score by modeling treatment assignment while optimizing the covariate balance for treatment effect estimation. Kuang et al. [23] proposed a data driven algorithm by jointly optimize variables separation and treatment effect estimation, where the separated confounders were used for confounding bias removing, and the separated adjustment variables were utilized for variance reduction. These methods are gaining ground in applied work, but they either require correct model specification on treatment assignment or precise estimation of the propensity score, which may not be the case in many applications [2], especially in high-dimensional settings.

Recently, researchers proposed new weighting based estimators by focusing on confounder balancing directly [2, 9, 15, 17, 45, 46], bypassing propensity score estimation. Hainmueller [15] introduced entropy balancing method to directly adjust sample weights to the specified sample moments while moving the sample weights as little as possible. Athey et al. [2] proposed approximate residual balancing algorithm, which, motivated by doubly robust approaches, combines outcome modeling using the LASSO with balancing weights constructed to approximately balance covariates between treatment and control groups. Zubizarreta [46] learned the stable balancing weights via minimizing its variance and adjusting for confounder balancing directly. Chan et al. [9] considered a wide class calibration weights constructed to attain confounder balancing directly. Imai et al. [17] introduced covariate balancing propensity score, which models treatment assignment while optimizing covariates balancing. Most of these methods are nonparametrical and require no propensity score estimation, but they do not differentiate the confounders by treating all observed variables as confounders and balanced all of them equally, leading to possible poor performance on treatment effect estimation in the setting of high-dimensional variables.

Hence, it is very likely to improve the treatment effect estimation efficiency by fine-tuned selection and differentiated methods. To achieve the goal, we propose a DCB algorithm to jointly optimize confounder weights and sample weights for precise treatment effect estimation.

Confounders selection: Recently, researchers had realized that not all observed variables are confounders and proposed some approaches for confounders selection [7, 34, 39, 43]. Most of these methods assumed the causal structure, i.e., whether a variable is the cause of treatment or outcome, is known a prior. But the causal structure cannot be well defined via prior knowledge in the wild, especially in the setting of high-dimensional variables. Here, we propose a data driven approach to learn the confounder weights for treatment effect estimation in the wild.

Learning for causal inference: Due to the big success in machine learning, many learning methods were utilized for causal inference, including deep neural network [19, 40], adversarial learning [20, 32], and variational autoencoder [29, 36]. In [19, 40], they proposed to adopt a deep neural

network to learn a variables representation, which has the same distribution on both treated and control groups, and learn a outcome regression model for counterfactual prediction. With adversarial learning techniques, the authors in [20, 32] proposed to learn a sample weight to adjust the variables' distribution on treated and control groups where the adversarial discriminator cannot distinguish the units from treated or control groups. By leveraging the techniques of variational inference, recent work [29, 36] proposed generative models to learn and capture the latent variables of confounders for better balancing. There are books for causal inference [18, 31, 34], and recently Guo et al. [14] present a survey of learning causality with data.

Comparing to the preliminary version [22], this one comprises a substantial amount of additional theoretical, algorithmic and experimental efforts and contributions. Key points of differences lie in the following aspects: First, as the bias analysis in our conference paper, we extend the theoretical analysis to both bias and variance of estimated treatment effect, and propose to utilize L_2 norm of sample weights to bound the variance. Second, by incorporating regression adjustments on outcome, we propose a new RA-DCB algorithm to further reduce the substantial bias of DCB algorithm in the setting with heavy confounding bias. Third, we report a series of statistical tests that examine the performance of the new RA-DCB algorithm for treatment effect estimation, and find that the method achieves more precise and robust results than DCB algorithm, especially in the settings with heavy confounding bias.

3 PROBLEM AND OUR ESTIMATOR

In this section, we first give the notations and problem formulation, then revisit traditional confounder balancing estimators, and propose a novel estimator via DCB and regression adjustment.

3.1 Notations and Problem Formulation

Our goal is to estimate the treatment effect based on potential outcome framework [18, 38]. With the framework, we define a treatment as a random variable T and a potential outcome as $Y(t)$, which corresponds to a specific treatment $T = t$. In this article, we only focus on binary treatment, that is $t \in \{0, 1\}$. We define the units that received treatment ($T = 1$) as treated units and the other units with $T = 0$ as control units. Then, for each unit indexed by $i = 1, 2, \dots, n$, we observe a treatment T_i , an outcome Y_i^{obs} , and a vector of observed pre-treatment variables $X_i \in \mathbb{R}^{p \times 1}$, where the observed outcome Y_i^{obs} of unit i is denoted by

$$Y_i^{obs} = Y_i(T_i) = T_i \cdot Y_i(1) + (1 - T_i) \cdot Y_i(0). \quad (1)$$

The numbers of treated and control units are equal to n_t and n_c , and the dimension of all observed variables is p . In our article, for any column vector $\mathbf{v} = (v_1, v_2, \dots, v_m)^T$, let $\|\mathbf{v}\|_\infty = \max(|v_1|, \dots, |v_m|)$, $\|\mathbf{v}\|_2^2 = \sum_{i=1}^m v_i^2$, and $\|\mathbf{v}\|_1 = \sum_{i=1}^m |v_i|$.

Throughout this article, we assume the SUTV and Unconfoundedness [38] condition is satisfied.

ASSUMPTION 1 (STABLE UNIT TREATMENT VALUE (SUTV)). *The distribution of potential outcome for one unit is assumed to be unaffected by the particular treatment assignment of another unit, when given the observed variables.*

ASSUMPTION 2 (UNCONFOUNDEDNESS). *The distribution of treatment is independent of potential outcome when given the observed variables. Formally, $T \perp (Y(0), Y(1)) | \mathbf{X}$.*

In this article, we focus on estimating the Average Treatment effect on the Treated (ATT), which represents the mean (average) difference between the potential outcomes under treated and control status among the treated subgroup. Formally, the ATT is defined as

$$ATT = E[Y(1)|T = 1] - E[Y(0)|T = 1], \quad (2)$$

Table 1. Symbols and Definitions

Symbol	Definition
n_t (n_c)	Sample size for treated (control) group
p	Dimension of observed (augmented) variables
$T \in \mathbb{R}^{n \times 1}$	Treatment
$Y \in \mathbb{R}^{n \times 1}$	Outcome
$\mathbf{X} \in \mathbb{R}^{n \times p}$	Observed variables
$\mathbf{X}_t \in \mathbb{R}^{n_t \times p}$	Observed variables of treated units
$\mathbf{M}_t \in \mathbb{R}^{n_t \times p}$	Augmented variables of treated units
$\mathbf{M}_c \in \mathbb{R}^{n_c \times p}$	Augmented variables of control units
$W \in \mathbb{R}^{n_c \times 1}$	Sample weights on control units
$\beta \in \mathbb{R}^{p \times 1}$	Confounder weights

where $Y(1)$ and $Y(0)$ represent the potential outcome of units with treatment status as treated $T = 1$ and control $T = 0$, respectively. Our method proposed in this article can be readily extended to estimate the Average Treatment effect on the Control (ATC) and hence the Average Treatment Effect (ATE) for the whole population.

In Equation (2), $E[Y(1)|T = 1]$ can be straightforwardly estimated by the sample analog $\sum_{i:T_i=1} \frac{1}{n_t} \cdot Y_i^{obs}$. But it is cumbersome to estimate $E[Y(0)|T = 1]$, since we cannot observe the potential outcome $Y(0)$ for the treated units. Under Assumption 1, $E[Y(0)|T = 1]$ is usually estimated by re-weighting techniques for removing the confounding bias. The reweighting methods form the surrogates of the unobserved potential outcome ($Y(0)|T = 1$) by reweighting the control units with sample weights W to make the confounder distributions on control units mimic the distributions on treated units. Then with the sample weights W on control units, we can estimate the *ATT* by

$$\widehat{ATT} = \sum_{i:T_i=1} \frac{1}{n_t} \cdot Y_i^{obs} - \sum_{j:T_j=0} W_j \cdot Y_j^{obs}. \quad (3)$$

3.2 Revisiting on Confounder Balancing

It can be seen from Equation (3) that the *ATT* estimation produces to sample weights learning problem. The classical approaches for sample weights learning are propensity score based methods [3, 4, 8]. The good performance of these methods hinges on the correct model specification for treatment assignment or accurate estimates of the propensity scores. Hence, the performance of these methods is often poor in the wild observational studies, where the model structure among variables is unknown.

To reduce the model dependency for applying on data in the wild, researchers proposed non-parametric methods to optimize the sample weights W by focusing on confounder balancing directly [2, 15]. The motivation behind these methods is that the confounders can be balanced by their moments, which uniquely determine their distributions. Therefore, they learn the sample weights W by

$$W = \arg \min_W \|\bar{\mathbf{X}}_t - \sum_{j:T_j=0} W_j \cdot X_j\|_2^2, \quad (4)$$

or

$$W = \arg \min_W \|\bar{\mathbf{X}}_t - \sum_{j:T_j=0} W_j \cdot X_j\|_\infty^2, \quad (5)$$

where the $\bar{\mathbf{X}}_t = \sum_{i:T_i=1} \frac{1}{n_t} X_i$ represents the mean value of observed variables on treated units. The direct confounder balancing methods based on Equation (4) or (5) can be applied on data in the

wild. But they balance all observed variables equally without differentiating confounders, which results in poor performance in the setting of high-dimensional variables.

3.3 Differentiated Confounder Balancing

To precisely estimate the treatment effect with high-dimensional observational data in the wild, we propose to simultaneously learn confounder weights and sample weights. The confounder weights can determine which variable is included and its share of contribution on confounding bias, and the sample weights are designed for confounder balancing.

To be specific, we jointly optimize the confounder weights and sample weights by learning following optimization under some constraints to be clarified later:

$$W = \arg \min_W \left(\beta^T \cdot \left(\bar{X}_t - \sum_{j:T_j=0} W_j \cdot X_j \right) \right)^2, \quad (6)$$

where $W \in \mathbb{R}^{n_c \times 1}$ is sample weights and $\beta \in \mathbb{R}^{p \times 1}$ is the confounder weights. In Equation (6), the confounder weights β differentiate the roles of each confounder in the balancing process, which helps for better removing the confounding bias in the wild observational studies.

Next, we give theoretical analysis on how to differentiate confounders weights with following proposition.

PROPOSITION 3.1. *In observational studies, different confounders make unequal confounding bias on treatment effect ATT with their own weights, and the weights can be learned by regressing potential outcome $Y(0)$ on observed variables \mathbf{X} .*

The general relationship among observed variables \mathbf{X} , treatment T , and outcome Y can be represented as

$$Y = f(\mathbf{X}) + T \cdot g(\mathbf{X}) + \epsilon, \quad (7)$$

where the *true ATT* is $E(g(\mathbf{X}_t))$, and the potential outcome $Y(0)$ can be represented by

$$Y(0) = f(\mathbf{X}) + \epsilon. \quad (8)$$

We prove Proposition 3.1 with following assumption.

ASSUMPTION 3 (LINEARITY). *The regression of potential outcome $Y(0)$ on observed variables \mathbf{X} is linear, that is $f(\mathbf{X}) = c + \alpha\mathbf{X}$.*

Under Assumption 3, we can rewrite the estimator of \widehat{ATT} as

$$\begin{aligned} \widehat{ATT} &= \sum_{i:T_i=1} \frac{1}{n_t} Y_i^{obs} - \sum_{j:T_j=0} W_j Y_j^{obs} \\ &= \sum_{i:T_i=1} \frac{1}{n_t} (c + \alpha X_i + g(X_i) + \epsilon_i) - \sum_{j:T_j=0} W_j (c + \alpha X_j + \epsilon_j) \\ &= E(g(\mathbf{X}_t)) + \left(\sum_{i:T_i=1} \frac{1}{n_t} \alpha X_i - \sum_{j:T_j=0} W_j \alpha X_j \right) + \phi(\epsilon) \\ &= ATT + \underbrace{\sum_{k=1}^p \alpha_k \left(\sum_{i:T_i=1} \frac{1}{n_t} X_{i,k} - \sum_{j:T_j=0} W_j X_{j,k} \right)}_{\text{Bias}} + \phi(\epsilon), \end{aligned} \quad (9)$$

where $\phi(\epsilon) = \sum_{i:T_i=1} \frac{1}{n_t} \epsilon_i - \sum_{j:T_j=0} W_j \epsilon_j \simeq 0$ refers to the difference of noises between treated and control units. In order to reduce the *Bias* term of estimated ATT, we have to regulate the term

$\sum_{k=1}^p \alpha_k \cdot (\sum_{i:T_i=1} \frac{1}{n_t} X_{i,k} - \sum_{j:T_j=0} W_j X_{j,k})$, where $(\sum_{i:T_i=1} \frac{1}{n_t} X_{i,k} - \sum_{j:T_j=0} W_j X_{j,k})$ means the difference of the k th confounder between treated and control units. The parameter α_k represents the confounding bias weight of the k th confounder, and it is the coefficient of X_k in the function $f(\mathbf{X})$. Hence, we can learn the confounder weights from the regression of potential outcome $Y(0)$ on observed variables \mathbf{X} under *Linearity* assumption.

Actually, the regression of potential outcome $Y(0)$ against on observed variables \mathbf{X} is infeasible, because of the counterfactual problem, that we cannot observe the potential outcome $Y(0)$ for treated units. Here, we utilize the sample weights W again to facilitate the construction of surrogates for the potential outcomes $Y(0)$ of the treated units. We will elaborate on this later.

When the function $f(\mathbf{X})$ is nonlinear, that is $f(\mathbf{X})$ allows for powers and interactions among observed variables. It is conceptually easy to extend above results under *Linearity* assumption to bound the bias of *ATT* with Taylor expansion on $f(\mathbf{X})$ by balancing not only observed variables, but also their powers and interactions. Therefore, when $f(\mathbf{X})$ is nonlinear, we have to balance the augmented variables $\mathbf{M} = (\mathbf{X}, \mathbf{X}^2, X_i X_j, \mathbf{X}^3, X_i X_j X_k, \dots)$, and learn the confounder weights by regressing the potential outcome $Y(0)$ on augmented variables \mathbf{M} .

Besides considering the bias of estimated treatment effect, we also give theoretical analysis on its variance with following proposition.

PROPOSITION 3.2. *If we assume the homogeneity of the variance of outcome when given observed variables and treatment, that is $\forall i, \text{Var}(Y_j|X_j, T_j) = \sigma^2$, then we can bound the variance of estimated *ATT* as a L_2 norm regularizer on the sample weights W .*

Under the linearity assumption in 3, we can write the Mean Squared Error (MSE) between estimated *ATT* (\widehat{ATT}) and real *ATT* (ATT) as

$$E\left(\left(\widehat{ATT} - ATT\right)^2 \mid \{X_i, T_i\}_{i=1}^n\right) = \left(\text{Bias} + \left(\sum_{i:T_i=1} \frac{1}{n_t} \epsilon_i - \sum_{j:T_j=0} W_j \epsilon_j \right) \right)^2 \quad (10)$$

$$= \text{Bias}^2 \quad (11)$$

$$= 2 \cdot \text{Bias} \cdot \left(\sum_{i:T_i=1} \frac{1}{n_t} E(\epsilon_i | X_i, T_i) - \sum_{j:T_j=0} W_j E(\epsilon_j | X_j, T_j) \right) \quad (12)$$

$$+ \underbrace{E\left(\sum_{i:T_i=1} \frac{1}{n_t} \epsilon_i - \sum_{j:T_j=0} W_j \epsilon_j\right)^2 \mid \{X_i, T_i\}_{i=1}^n}_{\text{Variance}}, \quad (13)$$

where $\text{Bias} = \sum_{k=1}^p \alpha_k (\sum_{i:T_i=1} \frac{1}{n_t} X_{i,k} - \sum_{j:T_j=0} W_j X_{j,k})$, which refers to the bias of estimated *ATT*. Equation (11) refers to the bias term of estimated *ATT*. Equation (12) equals 0, as $E(\epsilon_i | X_i, T_i) = 0$. By using the assumption that ϵ_i are independent and therefore $E(\epsilon_i \epsilon_j) = 0$ for $i \neq j$, then we can simply rewrite the variance term in Equation (13) as

$$\text{Variance} = \frac{1}{n_t^2} \sum_{i:T_i=1} E(\epsilon_i^2 | X_i, T_i) + \sum_{i:T_i=0} W_j^2 E(\epsilon_j^2 | X_j, T_j). \quad (14)$$

$$= \frac{1}{n_t^2} \sum_{i:T_i=1} \text{Var}(Y_i | X_i, T_i) + \sum_{i:T_i=0} W_j^2 \text{Var}(Y_j | X_j, T_j). \quad (15)$$

If we consider the homogeneity of $\text{Var}(Y_j | X_j, T_j)$, denoted by σ^2 . Then, we can bound above variance term as

$$\text{Variance} = \frac{1}{n_t^2} \sum_{i:T_i=1} \sigma^2 + \sum_{i:T_i=0} W_j^2 \sigma^2.$$

That is we can minimize the variance term by a L_2 norm regularizer on sample weights W .

Consequently, we may write the MSE as

$$E \left((\widehat{ATT} - ATT)^2 \mid \{X_i, T_i\}_{i=1}^n \right) = \underbrace{\left(\sum_{k=1}^p \alpha_k \left(\sum_{i:T_i=1} \frac{1}{n_t} X_{i,k} - \sum_{j:T_j=0} W_j X_{j,k} \right) \right)^2}_{\text{Bias}} \quad (16)$$

$$+ \underbrace{\frac{1}{n_t^2} \sum_{i:T_i=1} \sigma^2 + \sigma^2 \sum_{i:T_i=0} W_j^2}_{\text{Variance}}.$$

Therefore, to minimize the MSE and precisely estimate the ATT, we need to consider the tradeoff between bias term and variance term. In the next section, we will introduce our algorithms for treatment effect estimation by simultaneously minimizing both bias and variance term.

4 MODEL AND OPTIMIZATION

In this section, we give details of our models for treatment effect estimation, including DCB model and Regression Adjusted DCB model.

4.1 Differentiated Confounder Balancing Model

With Proposition 3.1, we know the ATT estimator is affected by the unbalance of the observed variables, and their high order terms. That is the augmented variables \mathbf{M} :

$$\mathbf{M} = (\mathbf{X}, \mathbf{X}^2, X_i X_j, \mathbf{X}^3, X_i X_j X_k, \dots). \quad (17)$$

Combining Equation (6) and (17) and Proposition 3.1, we give our objective function to jointly optimize sample weights and confounder weights for ATT estimation in observational studies as

$$\begin{aligned} \min \quad & \left(\beta^T \cdot (\bar{\mathbf{M}}_t - \mathbf{M}_c^T W) \right)^2, \quad (18) \\ \text{s.t.} \quad & \sum_{j:T_j=0} (1/n_t + W_j) \cdot (Y_j - M_j \cdot \beta)^2 \leq \lambda, \\ & \|W\|_2^2 \leq \delta, \quad \|\beta\|_2^2 \leq \mu, \quad \|\beta\|_1 \leq \nu, \\ & \mathbf{1}^T W = 1 \quad \text{and} \quad W \geq 0, \end{aligned}$$

where W is the sample weights and β is the confounder weights. $\bar{\mathbf{M}}_t$ represents the mean value of augmented variables on treated units. $\sum_{j:T_j=0} (1/n_t + W_j) \cdot (Y_j - M_j \cdot \beta)^2$ refers to the loss function of potential outcome $Y(0)$ when learning the confounder weights, including potential outcome loss on both control units $\sum_{j:T_j=0} 1/n_t (Y_j - M_j \cdot \beta)^2$ and treated units $\sum_{j:T_j=0} W_j \cdot (Y_j - M_j \cdot \beta)^2$, which is again a surrogate by weighting. With the constraints $\|\beta\|_2^2 \leq \mu$ and $\|\beta\|_1 \leq \nu$, we can remove the nonconfounders and smooth the confounder weights. The formula $\mathbf{1}^T W = 1$ normalizes the sample weights on control units to add up to one, with the sample weights on treated units. The terms $W \geq 0$ constraint each of sample weights is nonnegative. With norm $\|W\|_2^2 \leq \delta$, we can reduce the variance of estimated ATT to achieve stability with theoretical guarantee.

These lead to the following optimization problem, which is to minimize $\mathcal{J}(W, \beta)$ with constraints on parameters W .

$$\begin{aligned} \mathcal{J}(W, \beta) \quad & = \quad \left(\beta^T \cdot (\bar{\mathbf{M}}_t - \mathbf{M}_c^T W) \right)^2 + \lambda \sum_{j:T_j=0} 1/n_t + W_j) \cdot (Y_j - M_j \cdot \beta)^2 \quad (19) \\ & + \delta \|W\|_2^2 + \mu \|\beta\|_2^2 + \nu \|\beta\|_1, \\ \text{s.t.} \quad & \mathbf{1}^T W = 1 \quad \text{and} \quad W \geq 0. \end{aligned}$$

Here, we propose an iterative method to minimize the above objective function (19).

ALGORITHM 1: Differentiated Confounder Balancing (DCB)

Input: Tradeoff parameters $\lambda > 0, \delta > 0, \mu > 0, \nu > 0$, Augmented Variables Matrix on treat units \mathbf{M}_t , Augmented Variables Matrix on control units \mathbf{M}_c and Outcome Y .

Output: Confounder Weights β and Sample Weights W

- 1: Initialize Confounder Weights $\beta^{(0)}$ and Sample Weights $W^{(0)}$
- 2: Calculate the current value of $\mathcal{J}(W, \beta)^{(0)} = \mathcal{J}(W^{(0)}, \beta^{(0)})$ with Equation (19)
- 3: Initialize the iteration variable $t \leftarrow 0$
- 4: **repeat**
- 5: $t \leftarrow t + 1$
- 6: Update $\beta^{(t)}$ by solving $\mathcal{J}(\beta^{(t-1)})$ in Equation (20)
- 7: Update $W^{(t)}$ by solving $\mathcal{J}(W^{(t-1)})$ in Equation (21)
- 8: Calculate $\mathcal{J}(W, \beta)^{(t)} = \mathcal{J}(W^{(t)}, \beta^{(t)})$
- 9: **until** $\mathcal{J}(W, \beta)^{(t)}$ converges or max iteration is reached
- 10: **return** β, W .

Firstly, we initialize sample weights $W = \{1/n_c, \dots, 1/n_c\}^T$ and confounder weights $\beta = \{1/p, \dots, 1/p\}^T$. Once the initial values are given, in each iteration, we first update β by fixing W , and then update W by fixing β . These steps are described as follows.

Update β : When fixing W , the problem (19) is equivalent to optimize following objective function:

$$\mathcal{J}(\beta) = (\beta^T \cdot (\bar{\mathbf{M}}_t - \mathbf{M}_c^T W))^2 + \lambda \sum_{j:T_j=0} (1/n_t + W_j) \cdot (Y_j - M_j \cdot \beta)^2 + \mu \|\beta\|_2^2 + \nu \|\beta\|_1, \quad (20)$$

which is a standard ℓ_1 norm regularized least squares problem and can be solved with any LASSO (or elastic net) solver. Here, we use the proximal gradient algorithm [33] with proximal operator to solve the objective function in (20).

Update W : By fixing β , we can obtain W by optimizing (19). It is equivalent to optimize following objective function:

$$\mathcal{J}(W) = (\beta^T \cdot (\bar{\mathbf{M}}_t - \mathbf{M}_c^T W))^2 + \lambda \sum_{j:T_j=0} (1/n_t + W_j) \cdot (Y_j - M_j \cdot \beta)^2 + \delta \|W\|_2^2, \quad (21)$$

$$s.t. \mathbf{1}^T W = 1 \quad \text{and} \quad W \geq 0.$$

For ensuring nonnegative of W with constraint $W \geq 0$, we let $W = \omega \odot \omega$, where $\omega \in \mathbb{R}^{p \times 1}$ and \odot refers to the Hadamard product. Then, the problem (21) can be reformulated as

$$\mathcal{J}(\omega) = (\beta^T \cdot (\bar{\mathbf{M}}_t - \mathbf{M}_c^T (\omega \odot \omega)))^2 + \lambda \sum_{j:T_j=0} (1/n_t + \omega_j \odot \omega_j) \cdot (Y_j - M_j \cdot \beta)^2 + \delta \|\omega \odot \omega\|_2^2, \quad (22)$$

$$s.t. \mathbf{1}^T (\omega \odot \omega) = 1.$$

The partial gradient of term $\mathcal{J}(\omega)$ with respect to ω is

$$\begin{aligned} \frac{\partial \mathcal{J}(\omega)}{\partial \omega} &= -4(\beta^T \cdot (\bar{\mathbf{M}}_t - \mathbf{M}_c^T (\omega \odot \omega))) \cdot \mathbf{M}_c \cdot \beta \odot \omega \\ &\quad + 4\delta \omega \odot \omega \odot \omega + 2\lambda \omega \odot (Y_c - \mathbf{M}_c \cdot \beta)^2. \end{aligned}$$

Then, we determine the step size a with line search, and update ω at t th iteration as

$$\omega^{(t)} = \omega^{(t-1)} - a \cdot \frac{\partial \mathcal{J}(\omega^{(t-1)})}{\partial \omega^{(t-1)}}.$$

With constraint $\mathbf{1}^T(\omega \odot \omega) = 1$, we normalize $\omega^{(t)}$ as

$$\omega^{(t)} = \frac{\omega^{(t)}}{\sqrt{\mathbf{1}^T(\omega^{(t)} \odot \omega^{(t)})}}.$$

Then, we update $W^{(t)}$ at t th iteration with

$$W^{(t)} = \omega^{(t)} \odot \omega^{(t)}.$$

We update β and W iteratively until the objective function (19) converges. The whole algorithm is summarized in Algorithm 1.

Finally, with the optimized sample weights W by our DCB algorithm, we can estimate the ATT with Equation (3).

4.2 Regression Adjusted DCB Model

Recently, regression adjustment [13] has been used to experimental data for treatment effect estimation. And in [5, 28], it has been proved that regression adjustment could help to reduce the variance of estimated treatment effect on experimental data. Also, in the literature of causal inference with observational data [2, 23], regression adjustment also has been applied to reduce the bias and variance of estimated treatment effect, therefore, achieve a more precise estimation of causal effect.

Inspired by regression adjustment in these works, we propose a regression adjusted estimator based on our DCB model for treatment effect estimation in observational studies, named as Regression Adjusted Differentiated Confounder Balancing (RA-DCB) estimator, where we utilize the confounder weights β learned in our DCB algorithm or outcome regression adjustment, aiming to further reduce the substantial bias and variance from DCB algorithm under some settings with heavy confounding bias in observational data.

Therefore, we estimate ATT in our new proposed RA-DCB estimation with regression adjustment as following:

$$\widehat{ATT} = \sum_{i:T_i=1} \frac{1}{n_t} \cdot (Y_i^{obs} - M_i\beta) - \sum_{j:T_j=0} W_j \cdot (Y_j^{obs} - M_j\beta), \quad (23)$$

where the parameter β can be learned by regression augmented variables M on outcome Y . Comparing with previous estimator in (3), our new RA-DCB estimator in Equation (23) can remove the substantial bias from augmented variables M by regression adjustment, which could help to reduce the variance of estimated ATT.

4.3 Complexity Analysis

During the procedure of optimization, the main cost of DCB Algorithm 1 is to calculate the loss $\mathcal{J}(W, \beta)$, update confounder weights β and sample weights W . We analyze the time complexity of each of them, respectively. For the calculation of the loss, its complexity is $O(np)$, where n is the sample size and p is the dimension of (augmented) variables. For updating β , this is standard LASSO problem and its complexity is $O(np)$. For updating W , the complexity is dominated by the step of calculating the partial gradients of function $\mathcal{J}(\omega)$ with respect to variable ω . The complexity of $\frac{\partial \mathcal{J}(\omega)}{\partial \omega}$ is $O(np)$.

In total, the complexity of each iteration in DCB Algorithm 1 is $O(np)$. Similarly, we can obtain the complexity of each iteration in RA-DCB estimator is also $O(np)$ totally.

4.4 Parameters Tuning

No ground truth for parameters tuning is the main challenge of causal inference in observational studies. To address this challenge, we apply matching method to estimate the ATT and set it as the “approximate ground truth” as [1, 23, 37] did. Specially, for each treated unit i , we find its closet match among control units as follow:

$$match(i) = \arg \min_{j:T_j=0} \|X_i - X_j\|_2^2. \quad (24)$$

To make the matching approximate to exactly matching, we drop unit i if $match(i) > \epsilon$. Then, we can obtain the “approximate ground truth” by comparing the average outcome between the matched treated and control units.

With the “approximate ground truth,” we can tune parameters for our algorithm and baselines with cross validation by grid searching.

5 EXPERIMENTS

In this section, we evaluate our algorithm on both synthetic and real-world datasets, comparing with the state-of-the-art methods.

5.1 Baseline Estimators

We implement following baseline estimators to evaluate the ATT for comparison.

- *Unadjusted estimator* \widehat{ATT}_{UNA} : It evaluates the ATT by directly comparing the average outcome between the treated and control units without adjusting data. It ignores the confounding bias in data.
- *IPW estimator* \widehat{ATT}_{IPW} [38]: It evaluates the ATT via reweighting units with inverse of propensity score. It relies on correct model specification for propensity score estimation.
- *Doubly robust estimator* \widehat{ATT}_{DR} [4]: It evaluates the ATT with combination of IPW and regression method. It relies on correct specification of propensity score or outcome regression models.
- *Entropy balancing estimator* \widehat{ATT}_{ENT} [15]: It evaluates the ATT by directly balancing on confounders and entropy loss on sample weights. It ignores the confounder weights.
- *Approximate residual balancing estimator* \widehat{ATT}_{ARB} [2]: It evaluates the ATT by combining weighting adjustment via directly balancing on confounders and regression adjustment on outcome. It ignores the confounder weights.

In this article, we implemented \widehat{ATT}_{IPW} and \widehat{ATT}_{DR} with *lasso regression* for variables selection.

5.2 Experiments on Synthetic Data

In this section, we introduce how to generate the synthetic datasets and demonstrate the effectiveness of our DCB algorithm with extensive experiments.

5.2.1 Dataset. To generate the synthetic datasets, we consider two sample sizes $n = \{2,000, 5,000\}$ and also vary the dimension of observed variables $p = \{50, 100\}$. We first generate the observed variables $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p)$ with independent Gaussian distributions as

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p \stackrel{iid}{\sim} \mathcal{N}(0, 1),$$

where \mathbf{x}_i represents value of the i th variable in \mathbf{X} .

To test the robustness of all estimators, we generate the binary treatment variable T from a logistic function (T_{logit}) and a misspecified function (T_{missp}) as

$$T_{logit} \sim \text{Bernoulli} \left(\frac{1}{1 + \exp \left(- \sum_{i=1}^{p \cdot r_c} s_c \cdot x_i + \mathcal{N}(0, 1) \right)} \right), \text{ and}$$

$$T_{missp} = 1 \text{ if } \sum_{i=1}^{p \cdot r_c} s_c \cdot x_i + \mathcal{N}(0, 1) > 0, T_{missp} = 0 \text{ otherwise.}$$

where we vary both *confounding rate* r_c and *confounding strength* s_c from 0 to 1. The confounding rate represents the ration of confounders to all observed variables, and the confounding strength refers to the bias strength of confounders on treatment. As increasing of confounding rate r_c and confounding strength s_c , the selection bias between treated and control groups become more and more serious.

We generate the outcome Y from a linear function (Y_{linear}) and a nonlinear function (Y_{nonlin}) as

$$Y_{linear} = T + \sum_{j=1}^p \left\{ I(\text{mod}(j, 2) \equiv 0) \cdot \left(\frac{j}{2} + T \right) \cdot \mathbf{x}_j \right\} + \mathcal{N}(0, 3),$$

$$Y_{nonlin} = T + \sum_{j=1}^p \left\{ I(\text{mod}(j, 2) \equiv 0) \cdot \left(\frac{j}{2} + T \right) \cdot \mathbf{x}_j \right\} + \mathcal{N}(0, 3)$$

$$+ \sum_{j=1}^{p-1} \left\{ I(\text{mod}(j, 10) \equiv 1) \cdot \frac{p}{2} \cdot (x_j^2 + x_j \cdot x_{j+1}) \right\},$$

where the $I(\cdot)$ is the indicator function and function $\text{mod}(x, y)$ returns the modulus after division of x by y .

Under different settings on treatment T and outcome Y , we know the *true ATT* in simulation. We evaluate the *ATT* with our algorithm, comparing with baselines.

5.2.2 Results. To evaluate the performance of our proposed method, we carry out the experiments for 100 times independently. Based on the estimated *ATT* (\widehat{ATT}), we calculate its *Bias*, standard deviations (*SD*), mean absolute errors (*MAE*), and root mean square errors (*RMSE*) with following definitions:

$$\text{Bias} = \left| \frac{1}{K} \sum_{k=1}^K \widehat{ATT}_k - ATT \right|$$

$$\text{SD} = \sqrt{\frac{1}{K} \sum_{k=1}^K (\widehat{ATT}_k - \frac{1}{K} \sum_{k=1}^K \widehat{ATT}_k)^2}$$

$$\text{MAE} = \frac{1}{K} \sum_{k=1}^K |\widehat{ATT}_k - ATT|$$

$$\text{RMSE} = \sqrt{\frac{1}{K} \sum_{k=1}^K (\widehat{ATT}_k - ATT)^2},$$

where K is the experimental times, \widehat{ATT}_k is the estimated *ATT* in k th experiment and ATT represents the *true treatment effect*.

By varying sample size n , variables' dimension p , function of treatment T , function of outcome Y , confounding rate r_c , and confounding strength s_c , we obtain the experimental results of our DCB algorithm in total eight different settings by comparing with all baselines. The experimental results are demonstrated in Tables 2 and 3.

From Tables 2 and 3, we have following observations and analyses:

- Unadjusted estimator fails when confounders are associated with both treatment and outcome. From our results, we find the unadjusted estimator makes huge error on ATT estimation, because it ignores the confounding bias in data.
- IPW and DR estimators have poor performance in the setting of high-dimensional variables or when the model specifications are incorrect. IPW and DR estimators make huge error under setting 3 and setting 4, where $T = T_{missp}$ and $Y = Y_{nonlin}$.
- ENT estimator has good performance only when the parameters $s_c = 0.2$ under setting 2, where the confounding bias is small in data, but its performance deteriorates as the confounding bias increasing. Since it ignores the confounder weights, which makes it unable to effectively remove the confounding bias in data.
- ARB estimator achieves better performance than other baselines in most of time, since it is nonparametric method with regression adjustment. However, it is far inferior to our proposed estimator. The key reason is that it balances all observed variables equally.
- Our proposed DCB estimator, by jointly optimizing both sample weights and confounder weights, achieves significant improvements over the baselines in all settings, when varying sample size n , dimension of variables p , confounding rate r_c , and confounding strength s_c .

Robustness test. We also show the robustness of our DCB estimator in Figure 1 by varying the sample size n , dimension of variables p , confounding rate r_c , and confounding strength s_c . From Figure 1, we find that as we decrease n or increase p , r_c , and s_c , the MAE of our DCB estimator is consistent stable and small, while the MAE of baseline estimators increases continuously. This demonstrates that our proposed estimator is more precise and robust than the baselines.

RA-DCB VS. DCB. In Table 4, we report the experimental results by comparing RA-DCB with DCB in different settings. From the results, we find our proposed RA-DCB algorithm achieve a comparable results with previous DCB algorithm in the settings with mild confounding bias. But when the bias become severe (severe bias could be induced by high dimension, high confounding rate and confounding strength), RA-DCB algorithm could have a better performance than DCB on treatment effect estimation. For example, in setting 3, when $n = 5,000$, $p = 100$, $s_c = 1$, $r_c = 0.8$, our RA-DCB algorithm make an obvious improvement on treatment effect estimation than DCB algorithm. To clearly demonstrate the effective of our RA-DCB algorithm, we report the results in Figure 3 by comparing with DCB algorithm in a severe bias setting, where $T = T_{missp}$, $Y = Y_{nonlin}$, $n = 5,000$, $p = 200$, $s_c = 1.0$, $r_c = 0.8$. From the results, we conclude that with considering the regression adjustment, our RA-DCB algorithm can have a better performance on treatment effect estimation than DCB algorithm in settings with severe bias.

5.2.3 Parameter Analysis. In our DCB algorithm, we have hype-parameters λ , δ , μ , and ν . As mentioned before, we tuned these parameters in our experiments with cross validation by grid searching, and each parameter is uniformly varied from $\{0.001, 0.01, 0.1, 1, 10, 100, 1,000\}$. We displayed the *Bias* of treatment effect estimation with respect to λ , δ , μ , and ν , respectively. As seen from Figure 2, the *Bias* do not change too much and the performance are relatively stable when parameters $\lambda \geq 1$ and $\delta, \mu, \nu \leq 1$. From Figure 2(a), we can see the *Bias* is huge when parameter λ is too small. The main reason is that small value of λ would slack the constrain on confounder weights learning, resulting in imprecise confounder weights, even the trivial solution $\beta = 0$. From

Table 2. Results on Synthetic Dataset in Setting 1 to 4

Setting 1: $T = T_{logit}$, $Y = Y_{linear}$, and $s_c = 1$													
	n/p	$n = 2,000, p = 50$			$n = 2,000, p = 100$			$n = 5,000, p = 50$			$n = 5,000, p = 100$		
r_c	Estimator	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE
$r_c = 0.2$	\widehat{ATT}_{UNa}	6.483 (3.460)	6.682	7.349	18.60 (8.859)	18.67	20.61	6.420 (2.050)	6.420	6.739	18.53 (5.148)	18.53	19.23
	\widehat{ATT}_{IPW}	2.220 (6.224)	4.866	6.609	8.365 (15.40)	14.47	17.52	1.907 (4.092)	3.648	4.514	8.033 (9.852)	10.52	12.71
	\widehat{ATT}_{DR}	0.118 (0.307)	0.253	0.329	1.591 (0.512)	1.591	1.672	0.059 (0.174)	0.145	0.183	1.446 (0.337)	1.446	1.485
	\widehat{ATT}_{ENT}	0.371 (0.477)	0.453	0.605	4.924 (3.167)	5.052	5.855	0.046 (0.254)	0.210	0.258	2.425 (1.229)	2.429	2.719
	\widehat{ATT}_{ARB}	0.074 (0.472)	0.376	0.477	0.868 (0.435)	0.881	0.971	0.027 (0.269)	0.217	0.270	0.365 (0.371)	0.447	0.520
	\widehat{ATT}_{DCB}	0.014 (0.121)	0.099	0.122	0.006 (0.119)	0.101	0.119	0.001 (0.073)	0.053	0.073	0.001 (0.085)	0.067	0.085
$r_c = 0.8$	\widehat{ATT}_{UNa}	51.06 (3.725)	51.06	51.19	143.0 (9.389)	143.0	143.3	50.45 (1.900)	50.45	50.48	142.1 (5.647)	142.1	142.2
	\widehat{ATT}_{IPW}	29.99 (4.048)	29.99	30.26	98.24 (8.462)	98.24	98.60	29.38 (2.216)	29.38	29.46	96.86 (5.899)	96.86	97.04
	\widehat{ATT}_{DR}	0.345 (0.253)	0.367	0.428	4.492 (0.333)	4.492	4.504	0.338 (0.136)	0.338	0.365	4.306 (0.227)	4.306	4.312
	\widehat{ATT}_{ENT}	15.06 (1.745)	15.06	15.16	63.02 (4.551)	63.02	63.19	10.09 (1.473)	10.09	10.19	51.99 (3.206)	51.99	52.09
	\widehat{ATT}_{ARB}	0.231 (0.645)	0.553	0.685	2.909 (0.491)	2.909	2.951	0.189 (0.504)	0.428	0.538	2.259 (0.468)	2.259	2.307
	\widehat{ATT}_{DCB}	0.003 (0.127)	0.102	0.127	0.020 (0.135)	0.114	0.136	0.003 (0.088)	0.072	0.088	0.012 (0.088)	0.073	0.089
Setting 2: $T = T_{logit}$, $Y = Y_{linear}$, and $r_c = 0.5$													
	n/p	$n = 2,000, p = 50$			$n = 2,000, p = 100$			$n = 5,000, p = 50$			$n = 5,000, p = 100$		
s_c	Estimator	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE
$s_c = 0.2$	\widehat{ATT}_{UNa}	11.80 (3.243)	11.80	12.24	43.38 (9.170)	43.38	44.34	11.53 (2.142)	11.53	11.73	42.64 (6.103)	42.64	43.07
	\widehat{ATT}_{IPW}	3.897 (2.759)	4.144	4.775	18.37 (8.317)	18.38	20.17	3.873 (2.055)	3.875	4.384	17.13 (5.971)	17.13	18.14
	\widehat{ATT}_{DR}	0.053 (0.150)	0.124	0.159	1.255 (0.265)	1.255	1.283	0.056 (0.104)	0.090	0.118	1.148 (0.180)	1.148	1.162
	\widehat{ATT}_{ENT}	0.023 (0.168)	0.128	0.170	0.174 (0.193)	0.208	0.260	0.001 (0.116)	0.090	0.116	0.089 (0.119)	0.120	0.149
	\widehat{ATT}_{ARB}	0.002 (0.170)	0.129	0.170	0.011 (0.184)	0.151	0.185	0.004 (0.119)	0.094	0.120	0.006 (0.121)	0.093	0.122
	\widehat{ATT}_{DCB}	0.011 (0.107)	0.086	0.107	0.013 (0.098)	0.080	0.099	0.003 (0.065)	0.053	0.065	0.004 (0.073)	0.060	0.073
$s_c = 0.8$	\widehat{ATT}_{UNa}	22.81 (3.610)	22.81	23.09	69.28 (9.608)	69.28	69.94	21.91 (1.908)	21.91	21.99	68.72 (5.410)	68.72	68.93
	\widehat{ATT}_{IPW}	9.984 (4.878)	10.15	11.11	40.64 (12.48)	40.64	42.51	9.263 (3.615)	9.263	9.943	40.31 (7.185)	40.31	40.94
	\widehat{ATT}_{DR}	0.185 (0.256)	0.256	0.316	3.234 (0.449)	3.234	3.265	0.177 (0.166)	0.205	0.243	3.051 (0.245)	3.051	3.061
	\widehat{ATT}_{ENT}	2.805 (1.153)	2.805	3.033	23.53 (4.432)	23.53	23.94	0.742 (0.447)	0.759	0.866	15.97 (2.519)	15.97	16.16
	\widehat{ATT}_{ARB}	0.059 (0.564)	0.455	0.567	1.861 (0.491)	1.861	1.924	0.005 (0.408)	0.327	0.408	1.133 (0.451)	1.133	1.219
	\widehat{ATT}_{DCB}	0.007 (0.124)	0.102	0.124	0.015 (0.123)	0.102	0.124	0.001 (0.083)	0.067	0.083	0.017 (0.076)	0.063	0.078
Setting 3: $T = T_{missp}$, $Y = Y_{nonlin}$, and $s_c = 1$													
	n/p	$n = 2,000, p = 50$			$n = 2,000, p = 100$			$n = 5,000, p = 50$			$n = 5,000, p = 100$		
r_c	Estimator	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE
$r_c = 0.2$	\widehat{ATT}_{UNa}	6.527 (5.367)	7.041	8.450	18.67 (14.04)	20.01	23.36	7.340 (3.425)	7.366	8.099	20.54 (9.992)	20.54	22.84
	\widehat{ATT}_{IPW}	5.061 (8.998)	8.542	10.32	17.31 (19.22)	21.90	25.86	6.707 (6.494)	7.934	9.336	19.81 (15.04)	21.79	24.87
	\widehat{ATT}_{DR}	6.334 (8.628)	8.562	10.70	23.65 (26.32)	29.16	35.38	6.493 (6.698)	7.637	9.329	23.44 (16.62)	24.77	28.74
	\widehat{ATT}_{ENT}	3.770 (2.166)	3.842	4.348	13.46 (5.854)	13.58	14.68	3.096 (1.285)	3.102	3.352	12.16 (3.585)	12.16	12.68
	\widehat{ATT}_{ARB}	0.643 (0.292)	0.647	0.706	3.757 (0.483)	3.757	3.788	0.512 (0.247)	0.517	0.569	3.288 (0.262)	3.288	3.299
	\widehat{ATT}_{DCB}	0.016 (0.316)	0.263	0.317	0.021 (0.364)	0.294	0.365	0.017 (0.169)	0.139	0.169	0.082 (0.214)	0.183	0.230
$r_c = 0.8$	\widehat{ATT}_{UNa}	53.26 (5.308)	53.26	53.53	145.2 (13.47)	145.2	145.9	53.12 (3.673)	53.12	53.24	145.2 (9.247)	145.2	145.4
	\widehat{ATT}_{IPW}	39.46 (6.404)	39.46	39.97	113.0 (16.91)	113.0	114.3	39.04 (4.424)	39.04	39.29	111.7 (10.19)	111.7	112.1
	\widehat{ATT}_{DR}	15.12 (8.433)	15.40	17.31	34.07 (28.29)	37.09	44.28	14.26 (5.613)	14.28	15.33	30.92 (15.90)	31.70	34.77
	\widehat{ATT}_{ENT}	29.83 (1.795)	29.83	29.89	97.32 (6.507)	97.32	97.54	25.73 (1.155)	25.73	25.76	85.63 (3.114)	85.63	85.68
	\widehat{ATT}_{ARB}	1.342 (0.337)	1.342	1.384	7.440 (0.566)	7.440	7.462	1.102 (0.230)	1.102	1.126	6.526 (0.325)	6.526	6.535
	\widehat{ATT}_{DCB}	0.076 (0.321)	0.255	0.330	0.024 (0.388)	0.298	0.389	0.003 (0.207)	0.171	0.207	0.021 (0.304)	0.248	0.305

(Continued)

Table 2. Continued

Setting 4: $T = T_{missp}$, $Y = Y_{nonlin}$, and $r_c = 0.5$													
	n/p	$n = 2,000, p = 50$			$n = 2,000, p = 100$			$n = 5,000, p = 50$			$n = 5,000, p = 100$		
s_c	Estimator	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE
$s_c = 0.2$	\widehat{ATT}_{UNa}	18.01 (5.556)	18.01	18.84	59.49 (14.13)	59.49	61.15	18.01 (3.178)	18.01	18.29	60.34 (8.923)	60.34	60.99
	\widehat{ATT}_{IPW}	7.288 (6.605)	8.429	9.836	32.24 (19.66)	33.23	37.76	7.372 (4.505)	7.516	8.639	33.39 (12.87)	33.39	35.78
	\widehat{ATT}_{DR}	3.408 (5.953)	5.735	6.859	13.87 (21.90)	21.33	25.92	3.130 (4.146)	4.360	5.194	13.87 (12.53)	15.54	18.69
	\widehat{ATT}_{ENT}	1.812 (0.818)	1.812	1.988	25.54 (6.241)	25.54	26.29	0.273 (0.160)	0.282	0.317	14.49 (2.800)	14.49	14.76
	\widehat{ATT}_{ARB}	0.159 (0.254)	0.244	0.300	2.960 (0.385)	2.960	2.985	0.055 (0.150)	0.131	0.160	1.899 (0.241)	1.899	1.915
	\widehat{ATT}_{DCB}	0.005 (0.223)	0.178	0.223	0.011 (0.288)	0.228	0.288	0.012 (0.120)	0.095	0.120	0.025 (0.158)	0.125	0.160
$s_c = 0.8$	\widehat{ATT}_{dir}	24.58 (5.276)	24.58	25.14	72.30 (13.95)	72.30	73.63	24.10 (3.219)	24.10	24.31	71.20 (8.771)	71.20	71.74
	\widehat{ATT}_{IPW}	18.34 (6.819)	18.34	19.56	57.07 (18.02)	57.07	59.85	17.65 (4.755)	17.65	18.28	54.95 (9.861)	54.95	55.83
	\widehat{ATT}_{DR}	11.23 (8.757)	12.46	14.24	32.35 (26.22)	35.39	41.65	11.17 (5.492)	11.17	12.44	28.06 (14.24)	28.29	31.46
	\widehat{ATT}_{ENT}	12.88 (1.956)	12.88	13.03	48.40 (5.818)	48.40	48.75	10.46 (1.315)	10.46	10.55	40.79 (2.773)	40.79	40.88
	\widehat{ATT}_{ARB}	0.993 (0.343)	0.993	1.050	6.052 (0.525)	6.052	6.075	0.807 (0.255)	0.807	0.846	5.176 (0.279)	5.176	5.183
	\widehat{ATT}_{DCB}	0.042 (0.310)	0.246	0.313	0.023 (0.364)	0.306	0.365	0.006 (0.211)	0.167	0.211	0.013 (0.237)	0.194	0.238

The *Bias* refers to the absolute error between the true and estimated ATT. The *SD*, *MAE*, and *RMSE* represent the standard deviations, mean absolute errors, and root mean square errors of estimated ATT (\widehat{ATT}) after 100 times independently experiments, respectively. The smaller *Bias*, *SD*, *MAE*, and *RMSE*, the better.

Figure 2(c) and (d), we find that the *Bias* increased as the increasing of μ and ν . This is because that large value of μ and ν would enforce the confounder weights close to *zero*. The Figure 2(b) demonstrates that the performance is insensitive to the parameter δ . To sum up, we can easily obtain the best hyper-parameters for our DCB algorithm.

5.3 Experiments on Real Data

In this section, we apply our DCB algorithm on two real datasets for ATT estimation and application, including the LaLonde dataset and an online advertising dataset.

5.3.1 LaLonde Dataset. First, we apply our DCB algorithm on the LaLonde [26] dataset,² a canonical benchmark in the causal inference literature [11, 15]. The LaLonde dataset used in our article consists of two parts. The first part comes from a randomized experiment on a large scale job training program, the National Support Work Demonstration (NSW).³ In the second part data, as [15] did, we replace the control group in randomized experiment with another control group drawn from the Current Population Survey-Social Security Administration file (CPS-1) where the measured covariates are the same with the experimental data. The treatment in this data is whether the participant attend the particular job training program or not, and the outcome is the earning in the year 1978. The data contains 10 raw observed variables, including earnings and employment status for year 1974 and 1975, education status (years of schooling and an indicator for completed high school degree), age, ethnicity (indicators for black and hispanic), and the married status.

Overall, there are 185 program participants (the treated units) and 260 nonparticipants (the control units) in the experimental data NSW. In the observational data CPS-1, we have 185 program participants and 15,992 nonparticipants. The randomized experimental data NSW provide the ground truth for estimating the ATT of the program. We estimate the ATT with the observational data CPS-1, comparing our proposed algorithm with the baselines.

²The dataset is available at <http://users.nber.org/~rdehejia/data/nswdata2.html>.

³Notice that we focus on the Dehejia and Wahba sampled dataset of the LaLonde.

Table 3. Results on Synthetic Dataset in Setting 5 to 8

Setting 5: $T = T_{missp}$, $Y = Y_{linear}$, and $s_c = 1$													
	n/p	$n = 2,000, p = 50$			$n = 2,000, p = 100$			$n = 5,000, p = 50$			$n = 5,000, p = 100$		
r_c	Estimator	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE
$r_c = 0.2$	\widehat{ATT}_{UNA}	7.366 (2.857)	7.370	7.901	19.31 (10.17)	19.58	21.83	7.388 (2.182)	7.388	7.704	20.43 (5.003)	20.43	21.04
	\widehat{ATT}_{IPW}	4.335 (4.960)	5.350	6.587	11.45 (13.30)	14.54	17.55	3.132 (3.484)	3.764	4.685	12.76 (8.194)	13.11	15.17
	\widehat{ATT}_{DR}	0.163 (0.290)	0.266	0.333	2.309 (0.446)	2.309	2.352	0.163 (0.168)	0.185	0.234	2.265 (0.302)	2.265	2.285
	\widehat{ATT}_{ENT}	2.284 (1.232)	2.287	2.595	8.716 (3.705)	8.729	9.470	1.556 (0.789)	1.565	1.745	7.311 (2.202)	7.311	7.636
	\widehat{ATT}_{ARB}	0.077 (0.643)	0.538	0.648	1.724 (0.447)	1.724	1.781	0.094 (0.497)	0.404	0.506	1.437 (0.435)	1.437	1.501
	\widehat{ATT}_{DCB}	0.005 (0.134)	0.106	0.134	0.025 (0.117)	0.092	0.120	0.003 (0.084)	0.067	0.084	0.000 (0.067)	0.052	0.067
$r_c = 0.8$	\widehat{ATT}_{UNA}	52.46 (3.347)	52.46	52.56	145.9 (8.598)	145.9	146.1	52.06 (1.963)	52.06	52.10	145.7 (5.380)	145.7	145.8
	\widehat{ATT}_{IPW}	35.31 (3.548)	35.31	35.49	105.3 (8.115)	105.3	105.6	34.51 (2.012)	34.51	34.57	104.5 (5.467)	104.5	104.6
	\widehat{ATT}_{DR}	0.437 (0.251)	0.442	0.504	4.885 (0.348)	4.885	4.897	0.396 (0.132)	0.396	0.417	4.649 (0.252)	4.649	4.656
	\widehat{ATT}_{ENT}	23.72 (1.416)	23.72	23.77	76.10 (3.331)	76.10	76.17	20.70 (1.059)	20.70	20.72	68.32 (2.304)	68.32	68.36
	\widehat{ATT}_{ARB}	0.357 (0.528)	0.514	0.637	3.534 (0.488)	3.534	3.567	0.276 (0.539)	0.457	0.605	3.034 (0.421)	3.034	3.063
	\widehat{ATT}_{DCB}	0.005 (0.128)	0.106	0.128	0.034 (0.124)	0.105	0.129	0.004 (0.084)	0.066	0.084	0.002 (0.086)	0.068	0.086
Setting 6: $T = T_{missp}$, $Y = Y_{linear}$, and $r_c = 0.5$													
	n/p	$n = 2,000, p = 50$			$n = 2,000, p = 100$			$n = 5,000, p = 50$			$n = 5,000, p = 100$		
s_c	Estimator	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE
$s_c = 0.2$	\widehat{ATT}_{UNA}	18.00 (3.090)	18.00	18.26	58.80 (9.332)	58.80	59.54	17.70 (1.734)	17.70	17.79	59.63 (5.402)	59.63	59.88
	\widehat{ATT}_{IPW}	5.868 (3.710)	5.905	6.943	27.90 (10.52)	27.91	29.82	5.612 (2.314)	5.612	6.071	27.45 (6.236)	27.45	28.15
	\widehat{ATT}_{DR}	0.093 (0.197)	0.182	0.218	2.191 (0.347)	2.191	2.218	0.099 (0.125)	0.129	0.160	1.972 (0.206)	1.972	1.983
	\widehat{ATT}_{ENT}	0.106 (0.237)	0.215	0.260	5.948 (1.987)	5.948	6.271	0.041 (0.147)	0.122	0.153	0.540 (0.279)	0.540	0.607
	\widehat{ATT}_{ARB}	0.007 (0.237)	0.190	0.237	0.444 (0.383)	0.484	0.586	0.002 (0.148)	0.118	0.148	0.017 (0.230)	0.190	0.231
	\widehat{ATT}_{DCB}	0.003 (0.099)	0.080	0.099	0.007 (0.124)	0.098	0.124	0.002 (0.070)	0.057	0.070	0.002 (0.075)	0.063	0.075
$s_c = 0.8$	\widehat{ATT}_{UNA}	23.99 (3.322)	23.99	24.22	71.72 (8.267)	71.72	72.19	24.25 (1.828)	24.25	24.32	72.19 (5.520)	72.19	72.40
	\widehat{ATT}_{IPW}	14.18 (3.898)	14.18	14.71	47.86 (9.081)	47.86	48.72	14.00 (2.514)	14.00	14.23	47.90 (6.710)	47.90	48.37
	\widehat{ATT}_{DR}	0.356 (0.244)	0.367	0.431	3.910 (0.466)	3.910	3.937	0.280 (0.141)	0.282	0.314	3.830 (0.268)	3.830	3.839
	\widehat{ATT}_{ENT}	9.040 (1.216)	9.040	9.122	35.08 (3.207)	35.08	35.22	6.990 (0.981)	6.990	7.058	30.22 (2.387)	30.22	30.32
	\widehat{ATT}_{ARB}	0.214 (0.579)	0.494	0.617	2.756 (0.528)	2.756	2.806	0.110 (0.530)	0.439	0.542	2.417 (0.420)	2.417	2.454
	\widehat{ATT}_{DCB}	0.003 (0.123)	0.099	0.123	0.013 (0.123)	0.098	0.123	0.000 (0.073)	0.057	0.073	0.003 (0.077)	0.065	0.077
Setting 7: $T = T_{logit}$, $Y = Y_{nonlin}$, and $s_c = 1$													
	n/p	$n = 2,000, p = 50$			$n = 2,000, p = 100$			$n = 5,000, p = 50$			$n = 5,000, p = 100$		
r_c	Estimator	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE
$r_c = 0.2$	\widehat{ATT}_{UNA}	6.639 (5.061)	7.128	8.348	18.60 (14.50)	19.95	23.59	6.301 (3.417)	6.403	7.168	16.44 (9.478)	16.80	18.98
	\widehat{ATT}_{IPW}	3.132 (10.15)	8.573	10.62	13.00 (26.69)	22.22	29.69	2.793 (6.610)	5.872	7.176	11.82 (15.47)	15.43	19.47
	\widehat{ATT}_{DR}	1.646 (8.908)	7.640	9.059	12.98 (25.93)	22.17	29.00	2.516 (6.266)	5.410	6.752	12.90 (15.89)	15.99	20.47
	\widehat{ATT}_{ENT}	1.908 (1.659)	2.062	2.529	10.88 (6.590)	11.19	12.72	0.780 (0.672)	0.835	1.029	7.509 (3.069)	7.552	8.112
	\widehat{ATT}_{ARB}	0.310 (0.305)	0.371	0.435	2.854 (0.464)	2.854	2.892	0.150 (0.228)	0.215	0.273	2.100 (0.267)	2.100	2.117
	\widehat{ATT}_{DCB}	0.000 (0.251)	0.204	0.251	0.004 (0.314)	0.257	0.314	0.015 (0.160)	0.129	0.160	0.023 (0.175)	0.139	0.176
$r_c = 0.8$	\widehat{ATT}_{UNA}	49.87 (5.283)	49.87	50.15	143.6 (15.26)	143.6	144.4	50.13 (3.167)	50.13	50.23	143.5 (10.02)	143.5	143.9
	\widehat{ATT}_{IPW}	31.81 (6.563)	31.81	32.48	105.5 (16.47)	105.5	106.8	32.58 (4.659)	32.58	32.91	104.6 (11.58)	104.6	105.2
	\widehat{ATT}_{DR}	10.86 (8.339)	11.57	13.69	24.45 (22.87)	28.56	33.48	11.63 (5.477)	11.67	12.86	28.23 (15.73)	29.17	32.32
	\widehat{ATT}_{ENT}	23.27 (2.175)	23.27	23.37	89.07 (5.759)	89.07	89.26	17.79 (1.395)	17.79	17.85	73.90 (3.717)	73.90	74.00
	\widehat{ATT}_{ARB}	1.032 (0.367)	1.032	1.096	6.783 (0.529)	6.783	6.804	0.774 (0.274)	0.774	0.821	5.697 (0.342)	5.697	5.707
	\widehat{ATT}_{DCB}	0.033 (0.308)	0.246	0.310	0.040 (0.395)	0.324	0.397	0.026 (0.185)	0.147	0.186	0.156 (0.251)	0.246	0.295

(Continued)

Table 3. Continued

Setting 8: $T = T_{logit}$, $Y = Y_{nonlin}$, and $r_c = 0.5$													
	n/p	$n = 2,000, p = 50$			$n = 2,000, p = 100$			$n = 5,000, p = 50$			$n = 5,000, p = 100$		
s_c	Estimator	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE
$s_c = 0.2$	\widehat{ATT}_{UnA}	10.85 (5.138)	10.87	12.01	41.68 (13.38)	41.68	43.77	11.53 (3.348)	11.53	12.01	41.02 (9.734)	41.02	42.16
	\widehat{ATT}_{IPW}	3.970 (4.980)	5.027	6.369	18.16 (14.94)	19.72	23.51	4.525 (3.891)	4.988	5.968	17.96 (10.84)	18.03	20.98
	\widehat{ATT}_{DR}	1.175 (4.624)	3.740	4.771	3.810 (15.86)	12.57	16.31	1.482 (3.303)	2.970	3.620	4.847 (9.978)	8.739	11.09
	\widehat{ATT}_{ENT}	0.154 (0.188)	0.203	0.243	9.315 (3.602)	9.315	9.987	0.101 (0.124)	0.133	0.160	2.035 (0.733)	2.035	2.163
	\widehat{ATT}_{ARB}	0.009 (0.172)	0.139	0.172	1.035 (0.302)	1.035	1.078	0.002 (0.113)	0.092	0.113	0.406 (0.151)	0.406	0.433
	\widehat{ATT}_{DCB}	0.004 (0.159)	0.121	0.159	0.006 (0.192)	0.152	0.192	0.008 (0.112)	0.089	0.112	0.015 (0.142)	0.116	0.143
$s_c = 0.8$	\widehat{ATT}_{UnA}	21.48 (5.483)	21.48	22.17	71.82 (14.83)	71.82	73.33	21.98 (3.225)	21.98	22.21	69.41 (9.158)	69.41	70.01
	\widehat{ATT}_{IPW}	10.64 (9.112)	11.94	14.01	49.78 (22.10)	50.77	54.47	12.22 (5.276)	12.23	13.31	47.33 (11.83)	47.33	48.79
	\widehat{ATT}_{DR}	5.907 (8.284)	8.353	10.17	21.41 (22.76)	25.53	31.24	6.601 (6.164)	7.591	9.031	21.09 (14.99)	22.45	25.87
	\widehat{ATT}_{ENT}	7.549 (2.105)	7.549	7.837	40.91 (7.127)	40.91	41.53	4.448 (1.021)	4.448	4.564	30.59 (3.313)	30.59	30.76
	\widehat{ATT}_{ARB}	0.596 (0.312)	0.607	0.672	4.843 (0.449)	4.843	4.864	0.356 (0.235)	0.369	0.426	3.884 (0.288)	3.884	3.894
	\widehat{ATT}_{DCB}	0.012 (0.282)	0.221	0.282	0.042 (0.325)	0.253	0.327	0.007 (0.153)	0.121	0.154	0.014 (0.188)	0.145	0.188

The *Bias* refers to the absolute error between the true and estimated ATT. The SD, MAE, and RMSE represent the standard deviations, mean absolute errors, and root mean square errors of estimated ATT (\widehat{ATT}) after 100 times independently experiments, respectively. The smaller *Bias*, *SD*, *MAE*, and *RMSE*, the better.

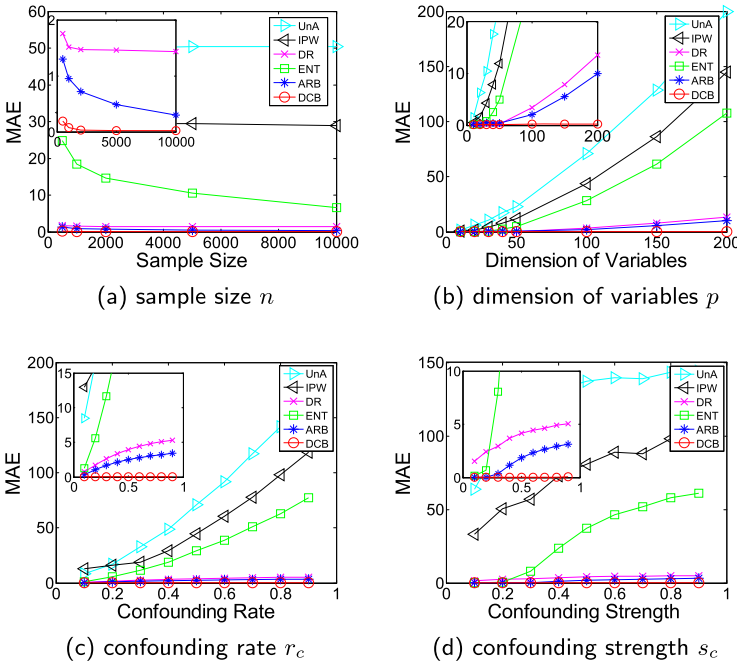


Fig. 1. MAE on ATT estimation when varying different parameters, with setting $T = T_{logit}$, $Y = Y_{linear}$. The subfigure on the top left corner of each main figure is plot by freezing MAE on Y-axis with a limit. The results show our proposed DCB estimator is more precise and robust than the baselines.

Table 4. Results on Synthetic Dataset by Comparing RA-DCB with DCB in Different Settings

Setting 1: $T = T_{logit}$, $Y = Y_{linear}$, and $s_c = 1$													
	n/p	$n = 2,000, p = 50$			$n = 2,000, p = 100$			$n = 5,000, p = 50$			$n = 5,000, p = 100$		
r_c	Estimator	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE
$r_c = 0.2$	\widehat{ATT}_{DCB}	0.024 (0.102)	0.083	0.104	0.024 (0.124)	0.103	0.127	0.004 (0.071)	0.057	0.071	0.003 (0.090)	0.071	0.090
	\widehat{ATT}_{RA-DCB}	0.024 (0.108)	0.088	0.110	0.029 (0.119)	0.099	0.123	0.005 (0.069)	0.056	0.070	0.002 (0.090)	0.074	0.090
$r_c = 0.8$	\widehat{ATT}_{DCB}	0.022 (0.129)	0.102	0.131	0.011 (0.142)	0.119	0.142	0.011 (0.089)	0.068	0.090	0.003 (0.106)	0.082	0.106
	\widehat{ATT}_{RA-DCB}	0.025 (0.129)	0.104	0.131	0.016 (0.146)	0.120	0.146	0.014 (0.094)	0.075	0.095	0.008 (0.101)	0.081	0.102
Setting 2: $T = T_{logit}$, $Y = Y_{linear}$, and $r_c = 0.5$													
	n/p	$n = 2,000, p = 50$			$n = 2,000, p = 100$			$n = 5,000, p = 50$			$n = 5,000, p = 100$		
s_c	Estimator	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE
$s_c = 0.2$	\widehat{ATT}_{DCB}	0.003 (0.098)	0.078	0.098	0.015 (0.124)	0.102	0.125	0.007 (0.075)	0.057	0.075	0.010 (0.060)	0.049	0.061
	\widehat{ATT}_{RA-DCB}	0.002 (0.104)	0.085	0.104	0.011 (0.129)	0.107	0.129	0.009 (0.074)	0.056	0.074	0.009 (0.066)	0.052	0.067
$s_c = 0.8$	\widehat{ATT}_{DCB}	0.043 (0.130)	0.107	0.137	0.018 (0.122)	0.103	0.123	0.003 (0.081)	0.065	0.082	0.004 (0.113)	0.097	0.113
	\widehat{ATT}_{RA-DCB}	0.042 (0.133)	0.109	0.139	0.020 (0.129)	0.113	0.130	0.005 (0.082)	0.065	0.082	0.004 (0.111)	0.094	0.111
Setting 3: $T = T_{missp}$, $Y = Y_{nonlin}$, and $s_c = 1$													
	n/p	$n = 2,000, p = 50$			$n = 2,000, p = 100$			$n = 5,000, p = 50$			$n = 5,000, p = 100$		
r_c	Estimator	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE
$r_c = 0.2$	\widehat{ATT}_{DCB}	0.018 (0.334)	0.270	0.335	0.020 (0.322)	0.249	0.323	0.010 (0.186)	0.147	0.186	0.101 (0.212)	0.185	0.235
	\widehat{ATT}_{RA-DCB}	0.018 (0.336)	0.271	0.337	0.002 (0.327)	0.257	0.327	0.013 (0.187)	0.148	0.188	0.063 (0.202)	0.170	0.211
$r_c = 0.8$	\widehat{ATT}_{DCB}	0.022 (0.269)	0.211	0.270	0.026 (0.337)	0.260	0.337	0.010 (0.172)	0.133	0.173	0.080 (0.238)	0.211	0.251
	\widehat{ATT}_{RA-DCB}	0.007 (0.262)	0.204	0.262	0.124 (0.385)	0.339	0.404	0.022 (0.173)	0.137	0.174	0.120 (0.181)	0.169	0.217
Setting 4: $T = T_{missp}$, $Y = Y_{nonlin}$, and $r_c = 0.5$													
	n/p	$n = 2,000, p = 50$			$n = 2,000, p = 100$			$n = 5,000, p = 50$			$n = 5,000, p = 100$		
s_c	Estimator	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE
$s_c = 0.2$	\widehat{ATT}_{DCB}	0.021 (0.199)	0.157	0.200	0.024 (0.284)	0.228	0.285	0.013 (0.122)	0.097	0.123	0.003 (0.183)	0.140	0.183
	\widehat{ATT}_{RA-DCB}	0.024 (0.200)	0.158	0.202	0.060 (0.278)	0.241	0.284	0.009 (0.122)	0.097	0.122	0.028 (0.181)	0.144	0.183
$s_c = 0.8$	\widehat{ATT}_{DCB}	0.024 (0.320)	0.270	0.321	0.031 (0.338)	0.278	0.339	0.005 (0.205)	0.163	0.205	0.037 (0.214)	0.178	0.218
	\widehat{ATT}_{RA-DCB}	0.029 (0.327)	0.276	0.329	0.132 (0.359)	0.291	0.382	0.000 (0.204)	0.163	0.204	0.080 (0.198)	0.177	0.213
Setting 5: $T = T_{missp}$, $Y = Y_{linear}$, and $s_c = 1$													
	n/p	$n = 2,000, p = 50$			$n = 2,000, p = 100$			$n = 5,000, p = 50$			$n = 5,000, p = 100$		
r_c	Estimator	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE
$r_c = 0.2$	\widehat{ATT}_{DCB}	0.046 (0.114)	0.102	0.123	0.014 (0.121)	0.096	0.122	0.018 (0.076)	0.062	0.078	0.012 (0.083)	0.066	0.084
	\widehat{ATT}_{RA-DCB}	0.045 (0.119)	0.106	0.127	0.015 (0.124)	0.101	0.125	0.020 (0.076)	0.063	0.079	0.014 (0.082)	0.065	0.083
$r_c = 0.8$	\widehat{ATT}_{DCB}	0.022 (0.134)	0.111	0.136	0.015 (0.146)	0.114	0.147	0.007 (0.089)	0.070	0.089	0.018 (0.078)	0.063	0.080
	\widehat{ATT}_{RA-DCB}	0.023 (0.138)	0.115	0.140	0.013 (0.163)	0.126	0.164	0.017 (0.090)	0.073	0.091	0.007 (0.082)	0.063	0.083
Setting 6: $T = T_{missp}$, $Y = Y_{linear}$, and $r_c = 0.5$													
	n/p	$n = 2,000, p = 50$			$n = 2,000, p = 100$			$n = 5,000, p = 50$			$n = 5,000, p = 100$		
r_c	Estimator	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE
$s_c = 0.2$	\widehat{ATT}_{DCB}	0.039 (0.106)	0.086	0.113	0.031 (0.148)	0.121	0.152	0.000 (0.065)	0.054	0.065	0.005 (0.075)	0.060	0.076
	\widehat{ATT}_{RA-DCB}	0.036 (0.104)	0.086	0.110	0.034 (0.151)	0.122	0.155	0.000 (0.065)	0.053	0.065	0.003 (0.077)	0.063	0.077
$s_c = 0.8$	\widehat{ATT}_{DCB}	0.002 (0.119)	0.094	0.119	0.021 (0.117)	0.089	0.119	0.001 (0.091)	0.073	0.091	0.001 (0.078)	0.065	0.078
	\widehat{ATT}_{RA-DCB}	0.003 (0.124)	0.097	0.124	0.015 (0.118)	0.090	0.119	0.002 (0.092)	0.072	0.092	0.002 (0.082)	0.069	0.082

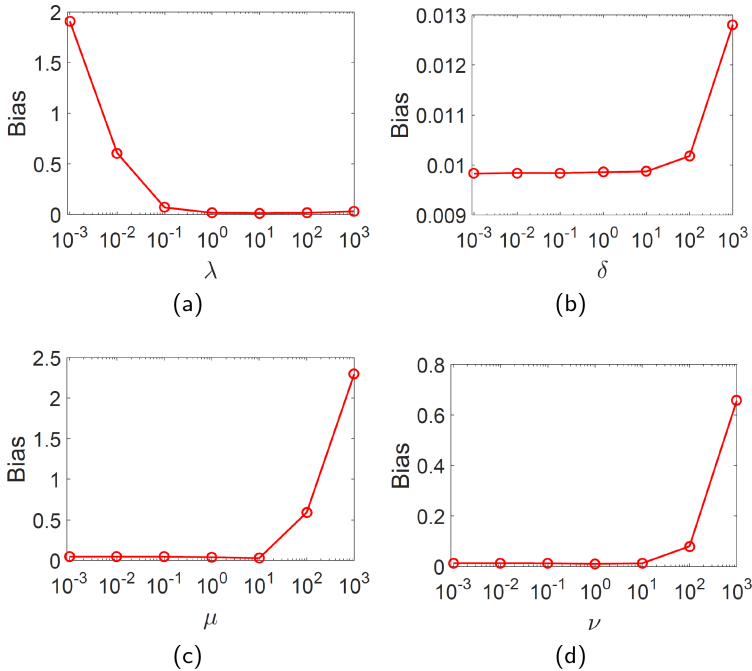
(Continued)

Table 4. Continued

Setting 7: $T = T_{logit}$, $Y = Y_{nonlinear}$, and $s_c = 1$													
r_c	n/p	$n = 2,000, p = 50$			$n = 2,000, p = 100$			$n = 5,000, p = 50$			$n = 5,000, p = 100$		
		Estimator	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE
$r_c = 0.2$	\widehat{ATT}_{DCB}	0.028 (0.260)	0.212	0.261	0.035 (0.272)	0.227	0.274	0.010 (0.150)	0.127	0.150	0.034 (0.206)	0.178	0.209
	\widehat{ATT}_{RA-DCB}	0.029 (0.258)	0.209	0.260	0.058 (0.283)	0.238	0.289	0.013 (0.149)	0.126	0.150	0.015 (0.205)	0.175	0.205
$r_c = 0.8$	\widehat{ATT}_{DCB}	0.001 (0.296)	0.239	0.296	0.144 (0.291)	0.261	0.325	0.021 (0.169)	0.135	0.170	0.136 (0.237)	0.215	0.274
	\widehat{ATT}_{RA-DCB}	0.011 (0.299)	0.242	0.299	0.159 (0.397)	0.335	0.427	0.019 (0.168)	0.135	0.169	0.042 (0.238)	0.201	0.241

Setting 8: $T = T_{logit}$, $Y = Y_{nonlinear}$, and $r_c = 0.5$													
r_c	n/p	$n = 2,000, p = 50$			$n = 2,000, p = 100$			$n = 5,000, p = 50$			$n = 5,000, p = 100$		
		Estimator	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE	RMSE	Bias (SD)	MAE
$s_c = 0.2$	\widehat{ATT}_{DCB}	0.003 (0.098)	0.078	0.098	0.015 (0.124)	0.102	0.125	0.007 (0.075)	0.057	0.075	0.010 (0.060)	0.049	0.061
	\widehat{ATT}_{RA-DCB}	0.002 (0.104)	0.085	0.104	0.011 (0.129)	0.107	0.129	0.009 (0.074)	0.056	0.074	0.009 (0.066)	0.052	0.067
$s_c = 0.8$	\widehat{ATT}_{DCB}	0.043 (0.130)	0.107	0.137	0.018 (0.122)	0.103	0.123	0.003 (0.081)	0.065	0.082	0.004 (0.113)	0.097	0.113
	\widehat{ATT}_{RA-DCB}	0.042 (0.133)	0.109	0.139	0.020 (0.129)	0.113	0.130	0.005 (0.082)	0.065	0.082	0.004 (0.111)	0.094	0.111

The *Bias* refers to the absolute error between the true and estimated ATT. The *SD*, *MAE*, and *RMSE* represent the standard deviations, mean absolute errors, and root mean square errors of estimated ATT (\widehat{ATT}) after 100 times independently experiments, respectively. The smaller *Bias*, *SD*, *MAE*, and *RMSE*, the better.

Fig. 2. The effect of hyper-parameters λ , δ , μ , and ν .

Experimental settings. In our experiments, we randomly split the observational data CPS-1 as six partitions, with the first three partitions, we train our model and baseline models for parameters tuning with cross validation by grid searching, and test model performance and robustness with the last three partitions. We conduct our DCB algorithm and baselines on two variables sets, V-RAW and V-INTERACTION. The V-RAW refers to the 10 raw observed variables, and the V-

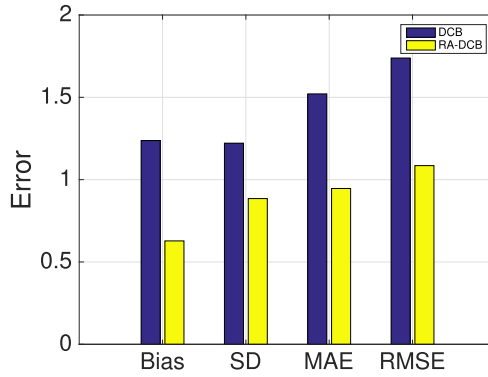


Fig. 3. Experimental results by comparing RA-DCB algorithm with DCB algorithm under setting $T = T_{missp}$, $Y = Y_{nonlin}$, $n = 5,000$, $p = 200$, $s_c = 1.0$, $r_c = 0.8$.

Table 5. ATT Estimation Results on LaLonde Dataset, Where the True ATT from Randomized Experiment is 1,794

Variables Set	V-RAW		V-INTERACTION	
	\widehat{ATT}	<i>Bias</i> (SD)	\widehat{ATT}	<i>Bias</i> (SD)
\widehat{ATT}_{dir}	-8,471	10,265 (374)	-8,471	10,265 (374)
\widehat{ATT}_{IPW}	-4,481	6,275 (971)	-4,365	6,159 (1024)
\widehat{ATT}_{DR}	1,154	639 (491)	1,590	204 (812)
\widehat{ATT}_{ENT}	1,535	259 (995)	1,405	388 (787)
\widehat{ATT}_{ARB}	1,537	257 (996)	1,627	167 (957)
\widehat{ATT}_{DCB}	1,958	164 (728)	1,836	43 (716)
\widehat{ATT}_{RA-DCB}	1,731	63 (523)	1,877	83 (520)

The smaller *Bias* and SD, the better.

INTERACTION refers to the set of raw variables, their pairwise one-way interaction, and their squared terms.

Results. We report the results in Table 5, where the smaller *Bias* and *SD*, the better. From the results, we have following observations. (1) Unadjusted estimator failed due to the existing of confounding bias in the LaLonde data. (2) IPW generates a big error on ATT estimation in both V-RAW and V-INTERACTION settings. The main reason is that the specification model of IPW is incorrect and the sample size between treated and control units is unbalanced. (3) Our proposed DCB and RA-DCB estimators outperform than all the baselines, since our estimators simultaneously optimizes sample weights and confounder weights, and requires no any model specification on treatment assignment. (4) With considering the regression adjustment, our estimator RA-DCB obtain a more accurate and robust result for ATT estimation than DCB algorithm under both V-RAW and V-INTERACTION settings, since the regression adjustment can help to further remove the confounding bias and reduce the variance of estimated ATT. (5) Under V-INTERACTION setting, our DCB and RA-DCB and also obtain a more robust (smaller SD) result than V-RAW setting. This demonstrates that our estimators can achieve a better confounder balancing and bias removing with including the high-order terms of observed variables in augmented variables.

In Table 6, we show the confounder weights optimized by our DCB algorithm with V-RAW variables set. From this table, we know that the confounders of Earnings 1975 & 1974 and Education

Table 6. Confounder Weights Learned from Our DCB Algorithm with V-RAW Variables Set

Rank	Confounder	Weight
1	Earnings 1975	0.335
2	Earnings 1974	0.241
3	Employed 1975	0.141
4	Education	0.138
5	Employed 1974	0.050
6	Married	0.039
7	High School Degree 1975	0.017
8	Age	-0.013
9	Black	-0.003
10	Hispanic	-0.001

are very important for the outcome (Earning 1978), but the Black and Hispanic have few effects on the outcome. That is the confounders of Earnings 1975 & 1974 and Education are more important, and should be balanced first.

5.3.2 Online Advertising Dataset. The real online advertising dataset we used is collected from Tencencent WeChat App⁴ during September 2015. In WeChat, each user can share (receive) posts to (from) his/her friends as like the Twitter and Facebook. Then the advertisers could push their advertisements to users, by merging them into the list of the user’s wallposts. For each advertisement, there are two types of feedbacks: “Like” and “Dislike.” When the user clicks the “Like” button, his/her friends will receive the advertisements with this action.

The online advertising campaign used in our article is about the LONGCHAMP handbags for young ladies.⁵ This campaign contains 14,891 user feedbacks with Like and 93,108 Dislikes. For each user, we have 56 features including (1) demographic attributes, such as age, gender, (2) number of friends, (3) device (iOS or Android), and (4) the user settings on WeChat, for example, whether allowing strangers to see his/her album and whether installing the online payment service.

Experimental settings. In our experiments, we set the feedback of users on the advertisement as outcome Y . Specifically, we set the outcome $Y_i = 1$ when user i likes the advertisement and $Y_i = 0$ when user i dislikes it. And we alternatively set one of the user features as the treatment T and others as the observed variables X . Therefore, we can estimate the ATT for each user feature. We tuned the parameters in our algorithm and baseline methods with the “approximal ground truth” via cross validation by grid searching.

Evaluation and baselines. In this dataset, we have no ground truth about the treatment effect of each user feature, but we are interesting in whether the top k features ranked by our proposed DCB estimator is able to get good performance in predicting the Like and Dislike behaviors of users, comparing with all above ATT baseline estimators and two commonly used methods for correlation-based feature selection, including MRel (Maximum Relevance) [42] and mRMR (Maximum Relevance Minimum Redundancy) [35]. Our estimator and other ATT baseline estimator rank the user features by their absolute causal effect. We use MAE as the evaluation metric, which

⁴<http://www.wechat.com/en/>.

⁵<http://en.longchamp.com/en/womens-bags>.

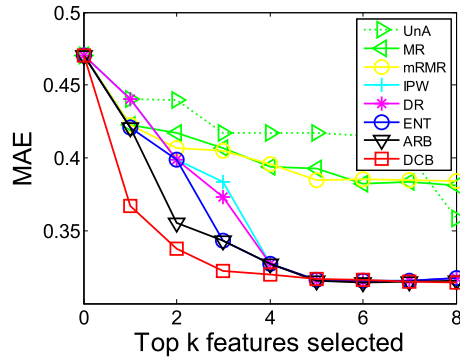


Fig. 4. Our proposed DCB estimator outperforms the baselines when selecting the top k significant causal features to predict whether user will like or dislike an advertisement.

is defined as

$$MAE = \frac{1}{m} \sum_{i=1}^m |\hat{Y}_i - Y_i|,$$

where m is the number of users in test data, \hat{Y}_i and Y_i represent the predict and actual feedback of user i on the advertisement, respectively.

Results. We plot the results in Figure 4. From the results, we can find that our proposed DCB estimator achieves the best prediction accuracy with different number of features. Also, our method can get almost the optimal prediction performance with much less features than other baselines. The main reason is that with differentiating the confounders, our estimator can estimate the causal effect of each user feature more precise by better confounding bias removing. Another important observation is that the two commonly used correlation-based feature selection methods perform worse than our method and even the other causal estimators. This is because of the sample selection bias between the training and testing datasets, the correlation-based methods cannot handle this issue, while the causal estimators can solve the problem to a certain extent by balancing treated and control units and removing the confounding bias.

The results demonstrate that treatment effect estimation can significantly help to improve the prediction performance, as long as the confounding problems are appropriately addressed.

6 CONCLUSION

In this article, we focus on how to estimate the treatment effect more precisely with high-dimensional observational data in the wild. We argued that most previous weighting based estimators do not take confounder differentiation into account or require model specification, leading to poor performance in the setting of high-dimensional variables or in the wild. Therefore, we proposed the concept of confounder weights for confounders differentiation with theoretical analysis. We proposed a DCB algorithm to jointly optimize the confounder weights and sample weights for treatment effect estimation. Then, with considering regression adjustment, we propose a Regression Adjusted DCB algorithm based on DCB algorithm for further removing confounding bias and improve the robustness of treatment effect estimation under some settings with severe confounding bias. Extensive experiments on both synthetic and real datasets demonstrated that our proposed algorithms can significantly and consistently outperforms the start-of-the-art methods, and RA-DCB algorithm can obtain more precise and robust estimation of treatment effect than DCB algorithm under the settings with severe bias. We also demonstrated that the top

ranked features by our algorithm have the best prediction performance on an online advertising dataset.

Our future will focus on causal inference with unobserved confounders in observational studies by data driven approaches.

ACKNOWLEDGMENTS

All opinions, findings, and conclusions in this article are those of the authors and do not necessarily reflect the views of the funding agencies.

REFERENCES

- [1] Susan Athey and Guido W. Imbens. 2015. Machine learning methods for estimating heterogeneous causal effects. *Stat* 1050, 5 (2015), 1–26.
- [2] Susan Athey, Guido W. Imbens, and Stefan Wager. 2018. Approximate residual balancing: De-biased inference of average treatment effects in high dimensions. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 80, 4 (2018), 597–623.
- [3] Peter C. Austin. 2011. An introduction to propensity score methods for reducing the effects of confounding in observational studies. *Multivariate Behavioral Research* 46, 3 (2011), 399–424.
- [4] Heejung Bang and James M. Robins. 2005. Doubly robust estimation in missing data and causal inference models. *Biometrics* 61, 4 (2005), 962–973.
- [5] Adam Bloniarz, Hanzhong Liu, Cun-Hui Zhang, Jasjeet S. Sekhon, and Bin Yu. 2016. Lasso adjustments of treatment effect estimates in randomized experiments. *Proceedings of the National Academy of Sciences* 113, 27 (2016), 7383–7390.
- [6] Léon Bottou, Jonas Peters, Joaquin Quinonero Candela, Denis Xavier Charles, Max Chickering, Elon Portugaly, Dipankar Ray, Patrice Y. Simard, and Ed Snelson. 2013. Counterfactual reasoning and learning systems: The example of computational advertising. *Journal of Machine Learning Research* 14, 1 (2013), 3207–3260.
- [7] M. Alan Brookhart, Sebastian Schneeweiss, Kenneth J. Rothman, Robert J. Glynn, Jerry Avorn, and Til Stürmer. 2006. Variable selection for propensity score models. *American Journal of Epidemiology* 163, 12 (2006), 1149–1156.
- [8] David Chan, Rong Ge, Ori Gershony, Tim Hesterberg, and Diane Lambert. 2010. Evaluating online ad campaigns in a pipeline: Causal models at scale. In *Proceedings of the 16th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 7–16.
- [9] Kwun Chuen Gary Chan, Sheung Chi Phillip Yam, and Zheng Zhang. 2016. Globally efficient non-parametric inference of average treatment effects by empirical balancing calibration weighting. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 78, 3 (2016), 673–700.
- [10] V. Chernozhukov, D. Chetverikov, M. Demirer, E. Duflo, C. Hansen and W. K. Newey. 2016. Double machine learning for treatment and causal parameters. Technical report, cemmap working paper.
- [11] Alexis Diamond and Jasjeet S. Sekhon. 2013. Genetic matching for estimating causal effects: A general multivariate matching method for achieving balance in observational studies. *Review of Economics and Statistics* 95, 3 (2013), 932–945.
- [12] Max H. Farrell. 2015. Robust inference on average treatment effects with possibly more covariates than observations. *Journal of Econometrics* 189, 1 (2015), 1–23.
- [13] David A. Freedman. 2008. On regression adjustments to experimental data. *Advances in Applied Mathematics* 40, 2 (2008), 180–193.
- [14] Ruo Cheng Guo, Lu Cheng, Jundong Li, P Richard Hahn, and Huan Liu. 2018. A survey of learning causality with data: Problems and methods. *arXiv:1809.09337*.
- [15] Jens Hainmueller. 2012. Entropy balancing for causal effects: A multivariate reweighting method to produce balanced samples in observational studies. *Political Analysis* 20, 1 (2012), 25–46.
- [16] Paul W. Holland. 1986. Statistics and causal inference. *Journal of the American statistical Association* 81, 396 (1986), 945–960.
- [17] Kosuke Imai and Marc Ratkovic. 2014. Covariate balancing propensity score. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 76, 1 (2014), 243–263.
- [18] Guido W. Imbens and Donald B. Rubin. 2015. *Causal Inference in Statistics, Social, and Biomedical Sciences*. Cambridge University Press.
- [19] Fredrik Johansson, Uri Shalit, and David Sontag. 2016. Learning representations for counterfactual inference. In *Proceedings of the International Conference on Machine Learning*. 3020–3029.
- [20] Nathan Kallus. 2018. Deepmatch: Balancing deep covariate representations for causal inference using adversarial training. *arXiv:1802.05664*.

- [21] Ron Kohavi and Roger Longbotham. 2011. Unexpected results in online controlled experiments. *ACM SIGKDD Explorations Newsletter* 12, 2 (2011), 31–35.
- [22] Kun Kuang, Peng Cui, Bo Li, Meng Jiang, and Shiqiang Yang. 2017. Estimating treatment effect in the wild via differentiated confounder balancing. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 265–274.
- [23] Kun Kuang, Peng Cui, Bo Li, Meng Jiang, Shiqiang Yang, and Fei Wang. 2017. Treatment effect estimation with data-driven variable decomposition. In *Proceedings of the 31st AAAI Conference on Artificial Intelligence*. AAAI.
- [24] Kun Kuang, Meng Jiang, Peng Cui, Hengliang Luo, and Shiqiang Yang. 2018. Effective promotional strategies selection in social media: A data-driven approach. *IEEE Transactions on Big Data* 4, 4 (2018), 487–501.
- [25] Kun Kuang, Meng Jiang, Peng Cui, and Shiqiang Yang. 2016. Steering social media promotions with effective strategies. In *Proceedings of the 2016 IEEE 16th International Conference on Data Mining (ICDM'16)*. IEEE, 985–990.
- [26] Robert J. LaLonde. 1986. Evaluating the econometric evaluations of training programs with experimental data. *The American Economic Review* 76, 4 (1986), 604–620.
- [27] Randall Lewis and David Reiley. 2009. Retail advertising works! Measuring the effects of advertising on sales via a controlled experiment on Yahoo! (2009).
- [28] Winston Lin. 2013. Agnostic notes on regression adjustments to experimental data: Reexamining Freedman’s critique. *The Annals of Applied Statistics* 7, 1 (2013), 295–318.
- [29] Christos Louizos, Uri Shalit, Joris M. Mooij, David Sontag, Richard Zemel, and Max Welling. 2017. Causal effect inference with deep latent-variable models. In *Proceedings of the Advances in Neural Information Processing Systems*. 6446–6456.
- [30] Daniel F. McCaffrey, Greg Ridgeway, and Andrew R. Morral. 2004. Propensity score estimation with boosted regression for evaluating causal effects in observational studies. *Psychological Methods* 9, 4 (2004), 403.
- [31] Stephen L. Morgan and Christopher Winship. 2015. *Counterfactuals and Causal Inference*. Cambridge University Press.
- [32] Michal Ozery-Flato, Pierre Thodoroff, and Tal El-Hay. 2018. Adversarial balancing for causal inference. *arXiv:1810.07406*.
- [33] N. Parikh and S. Boyd. 2014. Proximal algorithms. *Foundations and Trends in Optimization* 1, 3 (2014), 127–239.
- [34] Judea Pearl. 2009. *Causality*. Cambridge University Press.
- [35] Hanchuan Peng, Fuhui Long, and Chris Ding. 2005. Feature selection based on mutual information criteria of max-dependency, max-relevance, and min-redundancy. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 27, 8 (2005), 1226–1238.
- [36] Vineeth Rakesh, Ruocheng Guo, Raha Moraffah, Nitin Agarwal, and Huan Liu. 2018. Linked causal variational auto-encoder for inferring paired spillover effects. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*. ACM, 1679–1682.
- [37] Craig A. Rolling and Yuhong Yang. 2014. Model selection for estimating treatment effects. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 76, 4 (2014), 749–769.
- [38] Paul R. Rosenbaum and Donald B. Rubin. 1983. The central role of the propensity score in observational studies for causal effects. *Biometrika* 70, 1 (1983), 41–55.
- [39] Brian C. Sauer, M. Alan Brookhart, Jason Roy, and Tyler VanderWeele. 2013. A review of covariate selection for non-experimental comparative effectiveness research. *Pharmacoepidemiology and Drug Safety* 22, 11 (2013), 1139–1145.
- [40] Uri Shalit, Fredrik D. Johansson, and David Sontag. 2017. Estimating individual treatment effect: Generalization bounds and algorithms. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. 3076–3085.
- [41] Richard M. Shiffrin. 2016. Drawing causal inference from big data. *Proceedings of the National Academy of Sciences* 113, 27 (2016), 7308–7309.
- [42] Kari Torkkola. 2003. Feature extraction by non-parametric mutual information maximization. *Journal of Machine Learning Research* 3, Mar (2003), 1415–1438.
- [43] Tyler J. VanderWeele and Ilya Shpitser. 2011. A new criterion for confounder selection. *Biometrics* 67, 4 (2011), 1406–1413.
- [44] Daniel Westreich, Justin Lessler, and Michele Jonsson Funk. 2010. Propensity score estimation: Neural networks, support vector machines, decision trees (CART), and meta-classifiers as alternatives to logistic regression. *Journal of Clinical Epidemiology* 63, 8 (2010), 826–833.
- [45] Q. Zhao et al. 2019. Covariate balancing propensity score by tailored loss functions. *The Annals of Statistics* 47, 2 (2019), 965–993.
- [46] José R. Zubizarreta. 2015. Stable weights that balance covariates for estimation with incomplete outcome data. *Journal of the American Statistical Association* 110, 511 (2015), 910–922.

Received October 2018; revised July 2019; accepted September 2019