

# Optimal matching for heterogeneous treatment effect estimation

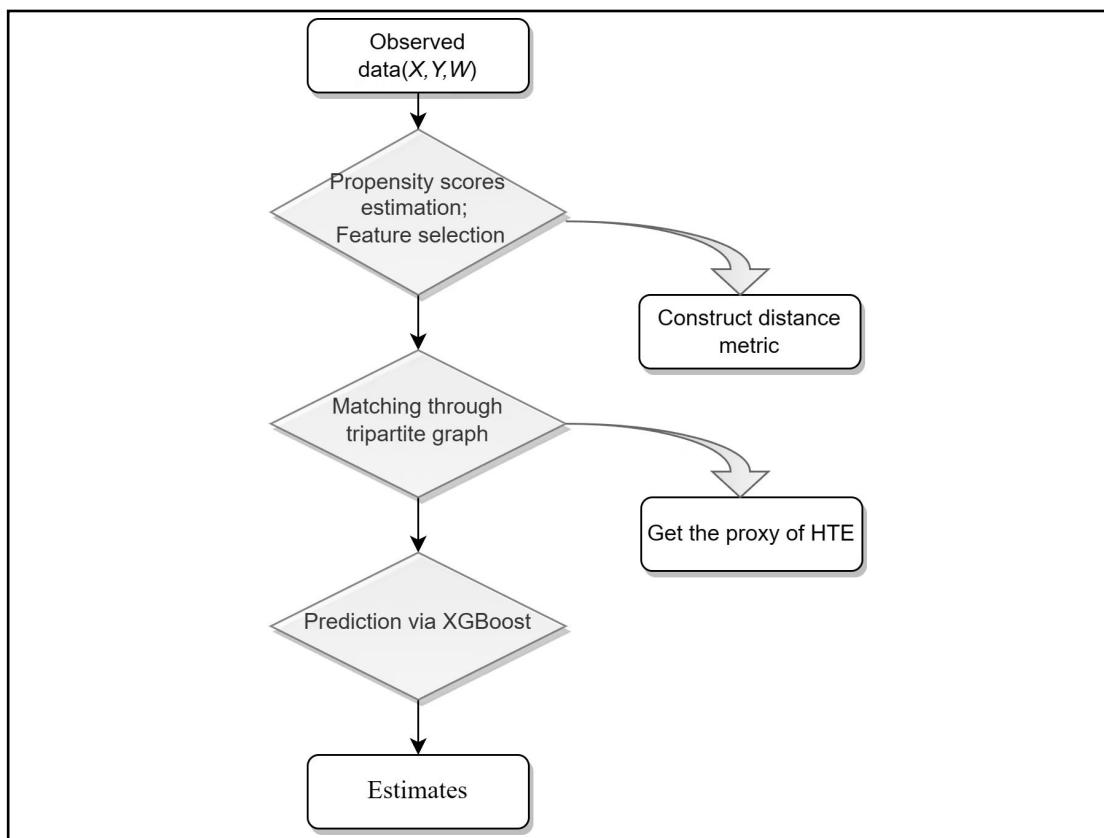
Yun Cai, and Shuguang Zhang 

*Department of Statistics and Finance, School of Management, University of Science and Technology of China, Hefei 230026, China*

Correspondence: Shuguang Zhang, E-mail: [sgzhang@ustc.edu.cn](mailto:sgzhang@ustc.edu.cn)

© 2023 The Author(s). This is an open access article under the CC BY-NC-ND 4.0 license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

## Graphical abstract



*Advantages of the proposed framework for predicting heterogeneous treatment effects by matching.*

## Public summary

- Utilizing the minimum average cost flow algorithm to tackle the optimization problem of multiobjective matching yields heightened flexibility and accuracy compared to conventional matching methods.
- Constructing an XGBoost tree using the acquired pseudo individual treatment effects yields better prediction accuracy compared to alternative regression-based methods.
- Both theoretical and experimental results demonstrate that the proposed method boasts a tolerable upper limit of estimation error while incurring minimal average matching costs.

# Optimal matching for heterogeneous treatment effect estimation

Yun Cai, and Shuguang Zhang 

Department of Statistics and Finance, School of Management, University of Science and Technology of China, Hefei 230026, China

 Correspondence: Shuguang Zhang, E-mail: sgzhang@ustc.edu.cn

© 2023 The Author(s). This is an open access article under the CC BY-NC-ND 4.0 license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).



Cite This: JUSTC, 2023, 53(7): 0707 (13pp)



Read Online

**Abstract:** In observational studies, identifying subgroups and exploring heterogeneity is of practical significance. However, causal inference at the individual level is a challenging problem due to the absence of counterfactual outcomes and the presence of selection bias. To address this issue, we propose a general framework called TRIMATCH for estimating heterogeneous treatment effects. First, we find the optimal matching by solving a minimum average cost flow optimization problem in a tripartite graph network structure. Second, with the pseudo individual treatment effects acquired from the previous step, we establish a nonparametric regression model to predict heterogeneous treatment effects for individuals with diverse characteristics. Our experiments demonstrate the effectiveness of the proposed matching method and the interpretability of the results.

**Keywords:** heterogeneous treatment effects; network flow; optimal matching; nonparametric regression

**CLC number:** O212.1

**Document code:** A

**2020 Mathematics Subject Classification:** 62D20

## 1 Introduction

Traditional causal inference is mainly carried out at the population level, focusing on the average treatment effects. However, this technique may fail to capture the nuanced heterogeneity in individual responses. As treatments may elicit varying effects among individuals, it is significant for researchers to infer heterogeneous treatment effects based on individual characteristics, enabling them to identify the beneficiary population. Exploration of the heterogeneity of causal effect leads decision makers to appropriate judgments grounded in the characteristics of individuals. With the deep-going research on heterogeneous treatment effects, the heterogeneity estimation of causal effects plays an increasingly important role in the fields of precision medicine<sup>[1]</sup>, marketing<sup>[2]</sup>, public policy<sup>[3]</sup>, and others that involve heterogeneity in the inference of causal effect at the individual or subgroup level.

The main methods for inferring causality are randomized controlled trials and observational data studies. Because randomized trials are time-consuming and costly, researchers mine causality from observational data. In causal inference studies, we can generally only observe the potential results of individuals under a single treatment condition; we cannot directly know the counterfactual results of individuals. The core problem of estimating heterogeneous treatment effects is how to obtain these counterfactual results. In observational studies, missing counterfactual outcomes and confounding bias are two major challenges in the evaluation of heterogeneous treatment effects. In the previous literature, many methods have emerged to estimate heterogeneous treatment effects from

observational data, such as tree-based methods<sup>[4, 5]</sup>, matching methods<sup>[6–8]</sup>, meta-learning methods<sup>[9, 10]</sup>, and neural networks<sup>[11]</sup>. Although the real treatment effect is unknown, it is still worth discussing how to estimate it and evaluate the efficacy of model estimation.

Among the methods for estimating heterogeneous treatment effects, matching is a relatively straightforward and interpretable method, but it does not consider extrapolation regions where there are no reasonable matching pairs. Hence, a possible solution to this challenge is to employ a hybrid approach that combines matching with machine learning methods.

In this article, we adopt a combination of matching and machine learning methods to detect heterogeneous treatment effects, introducing a new method we call “TRIMATCH”. Specifically, we obtain the pseudo individual treatment effects via the matching method and then establish a nonparametric regression model based on the pseudo treatment effects. Different from previous literature on matching, TRIMATCH considers both the closeness of the match and the balance of the covariates from the perspective of the average distance. It is necessary to seek a Pareto solution when dealing with multiobjective optimization problems, which entails minimizing the average closeness of matched pairs on a multivariate distance and minimizing the average imbalance. This solution can be derived by utilizing the minimum average cost flow algorithm, which is a tried and tested approach. In contrast to the previous method of bipartite matching, we adopt a new matching network: the tripartite graph proposed by Zhang et al.<sup>[12]</sup> Based on the tripartite matching structure, Zhang et al.<sup>[12]</sup> used the minimum cost flow algorithm to solve the

multiobjective matching problem from the perspective of the minimum matching distance. They proved, theoretically and experimentally, that the tripartite graph matching method is better than the traditional bipartite graph matching method, but this method focuses on fixed proportion matching; that is, the number of matches  $k$  of each individual is unified. Fixed proportion matching has limitations due to its high dependence on the sample distribution.

Drawn upon the tripartite graph, the objective of the TRIMATCH matching algorithm in this paper is to find the matching set with variable matching ratios that minimizes the average matching distance. Theoretical and experimental findings indicate that the algorithm for minimum average matching distance surpasses its counterpart, the algorithm for minimum matching distance, in generating a greater number of matching combinations while simultaneously minimizing average matching distance<sup>[7]</sup>. By incorporating this tripartite structure, we can better consider the matching goals separately and improve the matching quality of the matching algorithm. Compared with the regression-based estimation approach, the new approach prioritizes balancing the covariates between the treatment and control groups in the observed data using a matching method. Doing so minimizes the impact of sample imbalance on subsequent predictions. Additionally, TRIMATCH is suitable for the analysis of heterogeneous treatment effects in the case of high-dimensional covariates.

The specific contributions of this article are as follows:

(I) In terms of estimating heterogeneous treatment effects, we propose a data-driven method for estimating heterogeneous treatment effects from observational data by integrating matching methods in causal inference with machine learning techniques. This approach, called TRIMATCH, combines the strengths of both methods to create a robust and accurate estimator. On the one hand, machine learning methods have good generalization performance and can perform reasonable extrapolation based on matched samples, which is difficult for traditional matching methods. On the other hand, machine learning methods are sensitive to sample distribution, and some “black box models” lack interpretability. The integration of these matching methods not only reduces the negative impact of sample imbalance on machine learning method estimation but also improves the interpretability of the model. This study uses simulation experiments to verify the effectiveness and accuracy of TRIMATCH and proves theoretically that the method has a tolerable upper limit of estimation error. This work’s integration of machine learning and traditional causal methods is conducive to the development of causal theory in the era of big data. The work also provides a reference for exploring the combination of traditional causal inference methods and machine learning methods.

(II) Regarding research on matching methods, this article considers the tradeoff between multiple matching objectives from the perspective of average cost and solves the optimization problem using a method for minimum average cost flow. In previous literature, researchers usually selected matching samples by minimizing the total distance of matching. Gao et al.<sup>[7]</sup> proposed a single-objective matching method based on the minimum average cost flow algorithm. Based on this research, this article extends the research perspective to

multiobjective matching. We present theoretical and experimental results showing that the matching set obtained by the proposed method has superior properties compared to prior methods.

The rest of this paper is organized as follows. Section 2 reviews methods for estimating heterogeneous treatment effects using matching techniques. Section 3 introduces a step-by-step approach to estimating heterogeneous treatment effects using TRIMATCH. Section 4 discusses the impact of six different caliper settings on simulated datasets and shows the proposed method’s performance on simulation datasets. In Section 5, the method is applied to real-world data, and the results are analyzed. Section 6 concludes with a discussion.

## 2 Related work

Matching is an effective method for estimating treatment effects in observational studies, which can be achieved by matching treated and control groups with similar covariate distributions<sup>[6]</sup>. The match-based approach reduces the estimation bias brought by observed confounders. The keys to solving the matching problem are to choose a distance metric to define “closeness” and a matching algorithm to implement matching.

For distance metrics, the most straightforward way to match is to perform exact matching based on some discrete covariates that define the distance as 0 when the covariates are equal and infinity otherwise. Since exact matching does not apply to the matching problem on continuous covariates, Iacus et al.<sup>[13]</sup> proposed a new method known as coarsened exact matching (CEM), which converts continuous variables into ordered multicategorical variables under the exact matching framework. This method considers the extrapolation region where unmatched units exist. However, as the covariate dimension increases, it is challenging to match exactly based on multiple covariates or coarsened covariates due to computational complexity and data problems. For multivariate matching, Euclidean distance<sup>[14]</sup> or Mahalanobis distance<sup>[6]</sup> was typically selected in the previous literature. Numerous studies have shown that matching based on the propensity score can reduce bias to ensure balance on the variables highly correlated with the treatment assignment<sup>[15, 16]</sup>. The prognostic score can be considered an analog to the propensity score, and it can be used to produce a balance on prognostically relevant variables<sup>[17]</sup>. Since both the propensity score and the prognostic score must be estimated, score model misspecification can result in low-quality matches. Compared to propensity score or prognostic score matching, matching on both scores may improve the accuracy of treatment effect estimation and enhance robustness to score model misspecification<sup>[18, 19]</sup>. In the matching process, previous publications recommend including confounding factors associated with treatment assignment and the outcome variable as much as possible. They also recommend excluding some variables to improve the estimate’s accuracy, such as the post-treatment variables and the instrumental variables<sup>[6, 20, 21]</sup>.

Various matching algorithms exist to achieve different matching goals. Based on a selected distance measure, the most direct and simple method is nearest-neighbor matching.

Unfortunately, in cases where the sizes of the treatment and control groups are disproportionate, utilizing a one-to-one matching strategy can result in sample wastage. Rubin<sup>[14]</sup> proposed 1 to  $k$  (or  $k$  to 1) matching where  $k$  is a fixed value, and the choice of  $k$  is generally contingent upon the tradeoff between estimated bias and variance. If some individuals are spatially distant from their  $k$  nearest neighbors within a given dataset, a fixed ratio of matches will result in matches of inferior quality. Considering the uneven distribution of observational data, variable ratio matching has been proposed to enhance bias reduction<sup>[22, 23]</sup>. To acquire satisfactory matching pairs with certain desired properties, researchers usually impose constraints on the matching algorithms, such as fine and near-fine balance, exact and near-exact matching on crucial nominal variables, or maximum number of matches per individual. Soft or hard constraints can be applied to filter out matching pairs that do not meet the desired criteria corresponding to the matching goal, and there are tradeoffs between different matching targets<sup>[24]</sup>. The problem of optimal matching satisfying specific properties can be expressed as an optimization problem subject to constraints, and such problems can be solved by mixed integer programming or network flow methods. Yu and Rosenbaum<sup>[25]</sup> explored the connection between directional penalties and integer programming technique to improve covariate balance in a matched sample. Morucci et al.<sup>[26]</sup> found hyper-box-shaped regions where the treatment effects are roughly constant throughout with mixed integer program method. Another approach to solve the matching problem is the network flow algorithm, which is usually reformulated into a minimum cost flow problem<sup>[27, 28]</sup>.

### 3 Methodology

#### 3.1 Overall design of the framework

We follow the potential outcomes framework proposed by Rubin<sup>[29]</sup>. Suppose we have  $N$  independent observations from  $(Y, W, X)$ , where  $Y$  is the observed outcome,  $W$  is a binary treatment of interest, and  $X$  is a  $P$ -dimensional vector of covariates, each unit's potential outcomes are given by  $(Y_i(0), Y_i(1))$ , whose components represent the outcomes of units assigned to the control group and the treatment group. The potential outcome can be modeled as  $Y_i = \mu(X_i) + W \cdot \tau(X_i) + v_i$ , where  $E(v_i) = 0$ <sup>[30]</sup>.

For each unit  $i$ , the individual treatment effect (ITE) is defined as

$$\tau_i := Y_i(1) - Y_i(0).$$

Since two potential outcomes for the same individual cannot be observed simultaneously, the true individual treatment effect is unknown. However, it is possible to estimate the average treatment effect among individuals with the same vector of covariates  $X$ , which is referred to as the conditional average treatment effect (CATE). Under the following classical assumptions, the CATE can be estimated to capture heterogeneity among causal effects within the population, which is defined as

$$\tau(x) := E(Y(1) - Y(0)|X = x).$$

**Assumption 1** (Stable unit treatment value assumption). The potential outcomes for any unit do not depend on the treatments assigned to other units. There are no different versions of the treatment.

**Assumption 2** (Unconfoundedness). Given the covariates  $X$ , treatment assignment  $W$  is independent of the potential outcome for each unit.  $(Y(1), Y(0)) \perp\!\!\!\perp W|X$ .

#### 3.2 Matching as a minimum average cost flow problem

The matching problem is usually modeled as an optimization problem subject to various constraints, and it can be solved in various ways, such as network optimization techniques. Motivated by Zhang et al.<sup>[12]</sup>, we adopt a tripartite graph technique to construct the match network structure rather than using a bipartite graph. In contrast to classical bipartite matching, tripartite matching can achieve the goals of matching proximity and covariate balance simultaneously without conflict. The matching problem can be transformed into a minimum average cost flow problem based on a tripartite graph, which is solved by general algorithms for that problem.

##### 3.2.1 Tripartite matching structure

Considering the general structure of the network, the basic framework of the network comprises a set of vertices  $\mathcal{N}$  and edges  $\mathcal{E}$ , where  $\mathcal{E} \subseteq \mathcal{N} \times \mathcal{N}$ . Each edge  $e \in \mathcal{E}$  is of the form  $e = (i, j)$ , which represents the edge between  $i$  and  $j$ , where  $i, j \in \mathcal{N}$ . In the left part of the tripartite graph, there are  $T$  treated subjects, denoted by  $\mathcal{T} = \{t_1, \dots, t_T\}$ , and  $C$  control subjects, denoted by  $\mathcal{C} = \{c_1, \dots, c_C\}$ , where  $C > T$  in general. Each subject (treatment or control) has a duplicate in the right part of tripartite graph, which is denoted by  $t'_i$  or  $c'_j$ . A source  $s$  and a sink  $s'$  exist in addition to the abovementioned vertices at the beginning and end of the graph, respectively. The structure of tripartite matching is shown in Fig. 1.

##### 3.2.2 Distance metric

To ensure the efficacy of matching estimators, matching designs generally aim to achieve covariate balance and matching closeness within an appropriate matching ratio. The former ensures that the distribution of covariates observed in the treatment group is similar to that in the matched control group, making it closer to a randomized study. The latter requires close pairing for key observed covariates, which can reduce the estimation bias brought by confounders.

Our proposed matching algorithm considers both of these

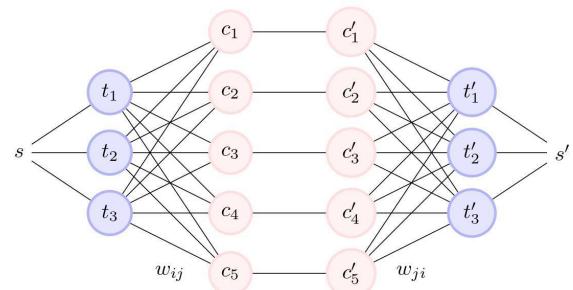


Fig. 1. Tripartite network structure.

matching goals simultaneously. We use  $\delta_{t,c_j}$  to represent the cost of edge  $(t_i, c_j)$  as the metric of pairing proximity, and we use  $\Delta_{t'_i, c'_j}$  to represent the cost of edge  $(c'_j, t'_i)$  as the metric of covariate imbalance. The penalty term denoted by  $\lambda$  achieves a tradeoff between the two distances.

Regarding matching closely, when the covariates are high-dimensional, it is not necessary to match based on all covariates, although Glazerman et al.<sup>[31]</sup> suggest including as many variables as possible that are highly predictive of outcomes and treatment assignment. Meanwhile, Pearl<sup>[32]</sup> advise excluding the instrumental variables affecting only the treatment assignment because they tend to amplify the bias of treatment effect estimators. We implement feature selection via random forests to screen for covariates with the smallest mean squared error in predicting outcome variables. For the multivariate distance, the Mahalanobis distance or robust Mahalanobis distance is typically used as the measure of distance because the Mahalanobis distance is not affected by dimension and takes into account the correlation between variables in multivariate settings, making it superior to the Euclidean distance. When the key covariates are all nominal variables, we can also utilize the Hamming distance as the distance measure or use exact matching.

Regarding balance matching, the propensity score is a balance score that refers to the estimated probability of a subject receiving the treatment arrangement given the covariates, defined as  $e(x) = P(W = 1 | X = x)$ . Matching on the propensity score tends to produce a balance between the treated and untreated groups, consequently removing bias that can arise due to the influence of covariates that are highly predictive of treatment assignment.

### 3.2.3 Minimum average cost flow

The network flow optimization algorithm is an important method for solving the variable ratio matching problem. Unlike the minimum cost flow-based algorithm matching, our objective is to minimize the average distance of the matches, which is equivalent to the minimum average cost flow optimization problem. As suggested by Gao et al.<sup>[7]</sup>, a matching method that aims to minimize the average match distance results in a greater number of matched pairs and lower average distance between the pairs than any method that focuses on minimizing the total match distance.

The problem of minimal average cost flow was put forward by Chen<sup>[33]</sup> to find the feasible flow from a specified node to another node with the minimum average transportation cost, which is essentially a problem of flow allocation and route planning with a given amount of flow. Let  $G = \{N, \mathcal{E}, b, \text{cap}, s, s', B\}$  be a network with  $|N|$  nodes,  $|\mathcal{E}|$  edges, and fixed cost  $B$ . For any edge in  $\mathcal{E}$ , the transportation cost per unit flow is  $b$ , the capacity of each edge is  $\text{cap}$ , and the flow is  $f$ . The total flow from the source node  $s$  to the sink node  $s'$  in network  $G$  is  $F$ , and the flow between two nodes  $u$  and  $v$  in the network is  $f(u, v)$  for  $u, v \in N$ , so the minimum average cost flow problem can be formulated as follows:

$$\begin{aligned} \min_f & \left\{ \left( \sum_{(u,v)} b(u,v) \times f(u,v) \right) + B \right\} / F \\ \text{s.t. } & \sum_v f(u,v) - \sum_v f(v,u) = \begin{cases} F & \text{if } u = s; \\ 0 & \text{if } u \notin s, s'; \\ -F & \text{if } u = s'; \end{cases} \\ & 0 \leq f(u,v) \leq \text{cap}(u,v), \quad (u,v) \in \mathcal{N}. \end{aligned} \quad (1)$$

The minimum cost flow problem is essentially a linear programming problem and can be solved by general algorithms in the previous literature. However, the minimum average cost flow problem is a nonlinear programming problem for which no general solution is currently known. Inspired by Chen<sup>[33]</sup>, we know that when the marginal cost is equal to the average cost, the corresponding flow in the flow network is the optimal solution of the minimum average cost flow problem. We use binary search to improve the operation efficiency, and we find the optimal solution by reformulating the minimum average cost flow problem into a minimum cost flow problem given the optimal feasible flow.

Based on the network flow model, the problem of finding the minimum average matching distance can be transformed into a minimum average cost flow problem. In the network structure depicted by Fig. 1, we use  $\mathcal{G}_1$  to denote the left part of the graph, consisting of the edges  $(s, t_i)$  and  $(t_i, c_j)$ , while using  $\mathcal{G}_2$  to represent the right part of the graph, comprising the edges  $(c'_j, t'_i)$  and  $(t'_i, s')$ . Since the objective is to find the optimal solution that minimizes the matching cost, we only need to consider the cost between the control group and the treatment group. In the tripartite graph, the cost of the edge between the treated units and controls corresponds to the matching distance, and the cost of all other edges is 0. The flow of the edge between the treated units and controls,  $f(e) \in \{0, 1\}, e \in \{(t_i, c_j)\}$ , indicates whether the individuals connecting this edge are matched. The minimum average cost flow problem can be reformulated as follows:

$$\min_f \frac{\sum_{(t_i, c_j) \in \mathcal{G}_1} \delta_{t_i, c_j} f(t_i, c_j) + \lambda \sum_{(c'_j, t'_i) \in \mathcal{G}_2} \Delta_{t'_i, c'_j} f(c'_j, t'_i)}{\sum_{(t_i, c_j)} f(t_i, c_j)} \quad (2a)$$

$$\text{s.t. } 0 \leq f(e) \leq 1, \forall e \in \{(t_i, c_j), (c'_j, t'_i)\}, \quad (2b)$$

$$1 \leq f(e) \leq k, \forall e \in \{(s, t_i), (t'_i, s')\}, \quad (2c)$$

$$0 \leq f(e) \leq k, \forall e \in \{(c_j, c'_j)\}, \quad (2d)$$

$$\sum_i f(s, t_i) = \sum_i f(t'_i, s'), \quad (2e)$$

$$\sum_{a(a,b) \in \mathcal{E}} f(a, b) = \sum_{b(b,c) \in \mathcal{E}} f(b, c), \quad \forall b \in N \setminus \{s, s'\}. \quad (2f)$$

The objective function contains two distance metrics representing the cost of two parts of the network structure. Part  $\mathcal{G}_1$  prioritizes pairing individuals who are close on key covariates, and  $\mathcal{G}_2$  gives preference to matching groups with smaller propensity score gaps. The former determines the

matching set, while the latter forces the matching set to satisfy the property of covariate balance. Hence, the resulting match is represented by the feasible flow of edge  $(t_i, c_j)$  in  $\mathcal{G}_1$ . In real-world data, the number of control group individuals is typically larger than that of treatment group individuals, and researchers are more concerned with the treatment effects on treatment groups in practical applications. Thus, in the matching process, every unit in the treatment group is guaranteed to be matched, while some control group units are allowed to remain unmatched. Constraint (2c) ensures that each treated unit is matched with at most  $k$  control units and at least one control unit. Constraints (2d) ensure that each control unit is matched with at most  $k$  treated units. There is a tradeoff between the matching accuracy and the number of matching pairs that should be considered in the selection of the  $k$  value; that is, a value of  $k$  that is too large will increase the bias of the estimate and reduce the computational efficiency, while an overly small  $k$  value may impractically decrease the number of pairs. What is worse, there is no feasible solution for the target function. Following Brito et al.<sup>[34]</sup>, we choose  $k \approx \log(n)$  as the maximum number of matches per individual. Constraint (2e) ensures that the flow out of source  $s$  is equal to the flow into sink  $s'$ . Constraint (2f) guarantees that the inflow and outflow of each unit are equal, except for the source and sink; that is, the net flow is 0.

**Proposition 3.1.** If there is a feasible flow in the above-mentioned network structure, then:

- (i) The flow through  $\mathcal{G}_1$  equals the flow through  $\mathcal{G}_2$ ; that is, the number of pairs generated by the two parts of the network is equal.
- (ii) For the same control unit in  $\mathcal{G}_1$  and  $\mathcal{G}_2$ , the number of pairs is equivalent, but the treated units paired with the control unit are not necessarily identical.
- (iii) The total distance of the optimal solution generated by tripartite matching is less than or equal to the minimum total distance generated by bipartite matching.

**Proof.** Eq. (2e) constrains the outbound flow from the source  $s$  to equal the inbound flow from the destination  $s'$ , so claim (i) holds. Eq. (2f) implies that for each control group, the flow out of  $c_j$  equals the flow into  $c'_j$ , and  $f(t_i, c_j) = 1$  or 0, so for the same control unit in  $\mathcal{G}_1$  and  $\mathcal{G}_2$ , the number of pairs of their matches are equal, while the treated units paired with the control unit are not necessarily identical, which supports claim (ii). If we force them to match the same set of treatment groups for the same control unit in  $\mathcal{G}_1$  and  $\mathcal{G}_2$ , by adding constraints  $f(t_i, c_j) = f(c'_j, t'_i)$ , this is equivalent to bipartite matching. In other words, under the same conditions, bipartite matching is more stringent than tripartite matching. Therefore, the optimal solution of bipartite matching must be a feasible solution of tripartite matching, but it is not necessarily the optimal solution of tripartite matching. Therefore, the minimum average cost of tripartite matching must be at most the minimum average cost of bipartite matching, which proves claim (iii).

As Korte et al.<sup>[35]</sup> showed, the algorithmic complexity of the minimum cost flow problem is  $O(|\mathcal{N}||\mathcal{E}| + |\mathcal{N}|^2 \log(|\mathcal{N}|))$ , which ranges from  $O(|\mathcal{N}|^2 \log(|\mathcal{N}|))$  to  $|\mathcal{N}|^2$  depending on the complexity of the network structure. The network  $(\mathcal{N}, \mathcal{E})$  is dense if  $|\mathcal{E}| = O(|\mathcal{N}|^2)$ , and sparse if  $|\mathcal{E}| = O(|\mathcal{N}|)$ . The

minimum average cost flow problem only takes  $O(\log(|\mathcal{E}| + \log(d_1)) \times T(|\mathcal{E}|, |\mathcal{N}|))$  arithmetic operations (see Ref. [33, Lemma 9]), where  $d_1$  is the largest capacity of the edge in our network, and  $T(|\mathcal{E}|, |\mathcal{N}|)$  is the time required for solving a minimum cost flow problem in a network with  $|\mathcal{E}|$  edges and  $|\mathcal{N}|$  nodes. In our network structure, we can reduce the complexity of the algorithm by increasing the sparsity of the graph. In other words, incorporating constraints can reduce the number of edges in the network. One viable approach to sparsifying the network is to coerce individuals into exact matches with key nominal variables. Another option is to impose a propensity caliper restriction on potential matches, that is, to discard matching pairs whose propensity score gap exceeds the given caliper threshold.

Depending on the matching objective, different types and widths of calipers can be chosen. In previous literature, this caliper is generally 0.2 pooled standard deviations of the estimated propensity score or prognostic score. Leacy and Stuart<sup>[18]</sup> proved that the model based on the joint use of propensity and prognostic scores performed better than the single score model. In our network structure, different calipers can be used in two parts of the network, with a loose caliper in  $\mathcal{G}_1$  and a tight caliper in  $\mathcal{G}_2$ . To remove the ineligible edges from the network, we can apply hard constraints by setting the cost of the edges that do not meet requirements to infinity or setting the capacity cap of the undesirable edges to 0. In addition, soft constraints are also used to remove those ineligible edges with a penalty on the cost of the edges. When comparing hard and soft constraints in optimization problems, the conditions of soft constraints are less stringent, and the methods employed to satisfy them are more adaptable. The objective may be feasible even when soft constraints are violated, making them a preferable choice in this paper. Nevertheless, when dealing with large volumes of data, hard constraints can contribute to simplifying the model and increasing the efficiency of the algorithm.

### 3.3 CATE estimation

Based on the result of matching, each unit's counterfactual outcome can be estimated by the mean of the outcomes of the individuals in its matching set. Let  $\Pi_i$  be the matching set of unit  $i$ , and  $|\Pi_i|$  denote the number of matches of  $i$ :

$$\Pi_i := \{j : f(i, j) = 1, W_j = 1 - W_i\}.$$

Let  $\tilde{\tau}$  be the proxy of the individual treatment effect, which can be computed by

$$\tilde{\tau}_i = (2W_i - 1) \left( Y_i - \frac{1}{|\Pi_i|} \sum_{j \in \Pi_i} Y_j \right).$$

Using a potential outcomes framework, Rosenbaum and Rubin<sup>[15]</sup> proved that  $Y$  and  $W$  are also conditionally independent, conditional on the propensity score under the unconfoundedness assumption,

$$(Y(1), Y(0)) \perp\!\!\!\perp W \mid e(x).$$

Therefore, we can derive an estimator of the conditional average treatment effect:

$$\begin{aligned}\tau(x) = & E(Y(1) - Y(0) | X_{re}, X_{ir}, e(X)) = \\ & E(Y(1) - Y(0) | X_{re}, e(X)),\end{aligned}\quad (3)$$

where  $X_{re}$  refers to covariates relevant to the outcome, and  $X_{ir}$  refers to covariates irrelevant to the outcome.

To estimate the conditional average treatment effect of each unit, we build the XGBoost tree for regression on the pseudo individual treatment effect as the proxy of the true individual treatment effect. The regressors include the key covariates screened for matching previously, which are highly predictive for the outcome variable. In the matching step, covariates that lack correlation with the outcome variable have been eliminated via implementation of the random forest methodology. The propensity score summarizes information highly predictive of the treatment indicator. When using full covariates for regression analysis, the curse of dimensionality often arises. To tackle this issue, we can remove variables that are irrelevant to the outcome variable. This will not only help to effectively reduce the dimensionality of predictive factors but also avoid the impact of redundant variables. The propensity score serves as a useful tool to balance covariates because it acts as a balancing score. Incorporating the propensity score into the predictive factors can help control confounding factors and improve the accuracy of predictions. Moreover, nonparametric regression that employs propensity scores resembles “one-to-one” matching based on propensity scores, establishing a one-to-one mapping relationship.

As stated in Proposition 3.2, the expected error upper bound of heterogeneous treatment effects can be effectively reduced by narrowing the distance within the matching group and the difference in prognostic score. See Appendix A.1 for a specific proof of Proposition 3.2.

**Proposition 3.2.** Let  $g$  be a regressor based on pseudo treatment effects  $\tilde{\tau}$ , and  $\hat{\tau}$  be the final estimated conditional average treatment effect, defined as  $\hat{\tau} = g(X, e(X))$ . Let  $\epsilon := \max |g(X, e(X)) - \tilde{\tau}|$ ,  $b_{\Pi_i} := \sum_{j \in \Pi_i} \frac{(2W_i - 1)(\mu_i - \mu_j)}{|\Pi_i|}$ , and  $v_{\Pi_i} := \sum_{j \in \Pi_i} \frac{(2W_i - 1)(v_i - v_j)}{|\Pi_i|}$ . Let  $\delta_{ij}^{(k)}$  be the  $k$ th closest distance

to unit  $i$  based on covariates. Assume that  $\tau$  is Lipschitz continuous; that is, there exists some constant  $L$  such that  $|\tau(X_i) - \tau(X_j)| \leq L|X_i - X_j|$ . Using the basic framework of the above model, if the tripartite matching problem is solvable within the caliper denoted by  $\alpha$ , then

$$E(\tilde{\tau}(X_i) - \tau(X_i)|X_i) \leq (2e(X_i) - 1)b_{ij} + (1 - e(X_i))L\delta_{ij}^{(k)}, \quad (4)$$

$$\begin{aligned}E((\hat{\tau}(X_i) - \tau(X_i))^2 | X_i) \leq \\ \{[b_{ij}^2 + \frac{k_i + 1}{k_i}\sigma^2 + (1 - e(X_i))(2L^2\delta_{ij}^{(k)} + b_{ij}^2)]^{\frac{1}{2}} + \epsilon\}^2.\end{aligned}\quad (5)$$

## 4 Simulation

This section analyzes the results from several simulation experiments carried out to compare the performance of TRIMATCH with some of the canonical models proposed for estimating heterogeneous treatment effects, including match-based and tree-based models. The simulated data were generated from the following model:

$$X_i \stackrel{i.i.d.}{\sim} F, \quad i = 1, \dots, n, \quad (6)$$

$$e(X_i) = \text{expit}(\gamma X_i), \quad (7)$$

$$W_i \stackrel{i.i.d.}{\sim} \text{Bern}(e(X_i)), \quad (8)$$

$$Y_i = \mu(X_i) + W_i \cdot \tau(X_i) + v_i, \quad v_i \stackrel{i.i.d.}{\sim} N(0, 1). \quad (9)$$

Note that  $e(\cdot)$  denotes a propensity score,  $\mu(\cdot)$  denotes a prognostic score defined as  $E(Y(0)|X)$ , and  $\tau(\cdot)$  denotes the function we want estimate, namely, the true individual treatment effect.

We compare the following estimators to TRIMATCH: ① propensity score matching (PSM); ② prognostic score matching (PGM); ③  $1:k$  matching on the joint use of propensity and prognostic scores (PP)<sup>[36]</sup>; ④ variable ratio matching based on the minimum average cost flow algorithm of a bipartite graph (COMBO)<sup>[7]</sup>; ⑤ a nonparametric Bayesian regression approach using dimensionally adaptive random basis elements with Bayesian additive regression trees (BART)<sup>[37]</sup>; and ⑥ an estimator based on an R-learner implemented by random forests with causal forests (CF)<sup>[4]</sup>. For all the  $1:k$  estimators mentioned above, Ye et al.<sup>[36]</sup> suggested that the empirical choice  $k \approx \log(n)$  yields good estimation accuracy at tolerable computational costs. In our proposed model, the magnitude of the penalty coefficient determines the tradeoff between the two matching objectives. Notably, a penalty coefficient of 1 yields matching outcomes that perform well in three crucial aspects: the number of matching groups, the gap in propensity score, and the distance of key covariates. We set the penalty coefficient  $\lambda$  at 1 for all subsequent experiments. For detailed experimental findings, kindly refer to Appendix A.2.

We examine the performance of each model in the following experimental settings under two scenarios: with linear confounding factors and with nonlinear confounding factors.

Denote the sample size by  $n \in \{1000, 2000\}$ , and the number of features observed for each unit by  $p \in \{10, 50, 100, 200\}$ .

### (I) Linear confounders:

$$X \stackrel{i.i.d.}{\sim} \mathcal{U}[-1, 1]^{n \times p}, \quad (10)$$

$$W|X \sim \text{Bern}\left(\frac{\exp(X_1 - 2X_2 - 1.5)}{1 + \exp(X_1 - 2X_2 - 1.5)}\right), \quad (11)$$

$$\tau(X) = -2X_1 - 1.5X_2 + X_3 + 0.9X_4 - 0.2X_5 + 1, \quad (12)$$

$$\mu(X) = 2X_1 + 1.5X_2. \quad (13)$$

### (II) Nonlinear confounders:

$$X \stackrel{i.i.d.}{\sim} \mathcal{U}[-1, 1]^{n \times p}, \quad (14)$$

$$W|X \sim \text{Bern}\left(\frac{\exp(X_1 + 2.5X_2 - 0.3X_2^2 + 0.4X_3^3 - 1.5)}{1 + \exp(X_1 + 2.5X_2 - 0.3X_2^2 + 0.4X_3^3 - 1.5)}\right), \quad (15)$$

$$\tau(X) = -2X_1 - 1.5X_2 + X_3 + 0.9X_4 - 0.2X_5 + 0.5X_3^2 - 0.3X_5^2 + 1, \quad (16)$$

$$\mu(X) = 2X_1 + 1.5X_2 + 0.5X_1^2 + 2.4X_2^2 - 0.8X_1X_2. \quad (17)$$

We estimate the propensity scores using logistic regression, which is implemented by the function *bart* in the R package *dbarts*<sup>[37]</sup>. Each method's performance was assessed according to mean absolute error (MAE), defined as  $\frac{1}{n} \sum_i |\hat{\tau}_i - \tau_i|$ , and mean square error (MSE) between the true treatment effects and estimated treatment effects, defined as  $\frac{1}{n} \sum_i (\hat{\tau}_i - \tau_i)^2$ .

#### 4.1 Comparison with different calipers

To determine the best choice of calipers and hyperparameters for the proposed TRIMATCH method, we compare the results of matching based on the average difference of the true propensity score (TP), the average difference of the estimated propensity score (EP) between the treatment group and the control group after matching, the average Mahalanobis distance of the treatment group and the control group based on key covariates (MAHAL), the distance of each key covariate between the control group and the treatment group, and the number of matched pairs (PAIRS). We conduct six simulation settings with different calipers as follows:

- a. Matching  $(t_i, c_j)$  without calipers on the graph's left part, and matching  $(c'_j, t'_i)$  within 0.2 pooled standard deviations of the estimated propensity score caliper on the graph's right part.
- b. Matching  $(t_i, c_j)$  within 0.2 pooled standard deviations of the estimated propensity score caliper on the graph's left part, and matching  $(c'_j, t'_i)$  without a caliper on the graph's right part.
- c. Matching  $(t_i, c_j)$  within 0.2 pooled standard deviations of the estimated propensity score caliper on the graph's left part, and matching  $(c'_j, t'_i)$  within 0.05 pooled standard deviations of the estimated prognostic score caliper on the graph's right part.
- d. Matching  $(t_i, c_j)$  within 0.2 pooled standard deviations of the estimated prognostic score caliper on the graph's left part, and matching  $(c'_j, t'_i)$  within 0.05 pooled standard deviations of the estimated propensity score caliper on the graph's right part.
- e. Matching  $(t_i, c_j)$  within 0.2 pooled standard deviations of

the estimated propensity score caliper on the graph's left part, and matching  $(c'_j, t'_i)$  within 0.05 pooled standard deviations of the estimated propensity score caliper on the graph's right part.

- f. Matching  $(t_i, c_j)$  within 0.2 pooled standard deviations of the estimated prognostic score caliper on the graph's left part, and matching  $(c'_j, t'_i)$  within 0.05 pooled standard deviations of the estimated prognostic score caliper on the graph's right part.

We generated experimental data with a sample size of 1000 and covariate dimension 10 based on the linear confounder model, and the simulation experiment was carried out 100 times with each of the above six different caliper settings. The final experimental results are shown in Table 1. From the simulation results, matching under caliper setting (a) can effectively reduce the gap between the treatment group and the matched control group and significantly improve the matching quality. Among the double caliper settings, matching under caliper setting (d) shows satisfactory performance in covariate balance and pairing proximity. The prognostic score extracts information about the covariates associated with the outcome, and the propensity score summarizes information relevant to the treated variables, so we recommend matching with the prognostic score caliper for closeness, using the propensity score caliper for balance. That is, calipers correspond to the matching objective function and the matching goals. Whether using a single caliper or double calipers, it is advisable to impose a propensity score caliper on the right part of the graph.

#### 4.2 Comparison with different estimators

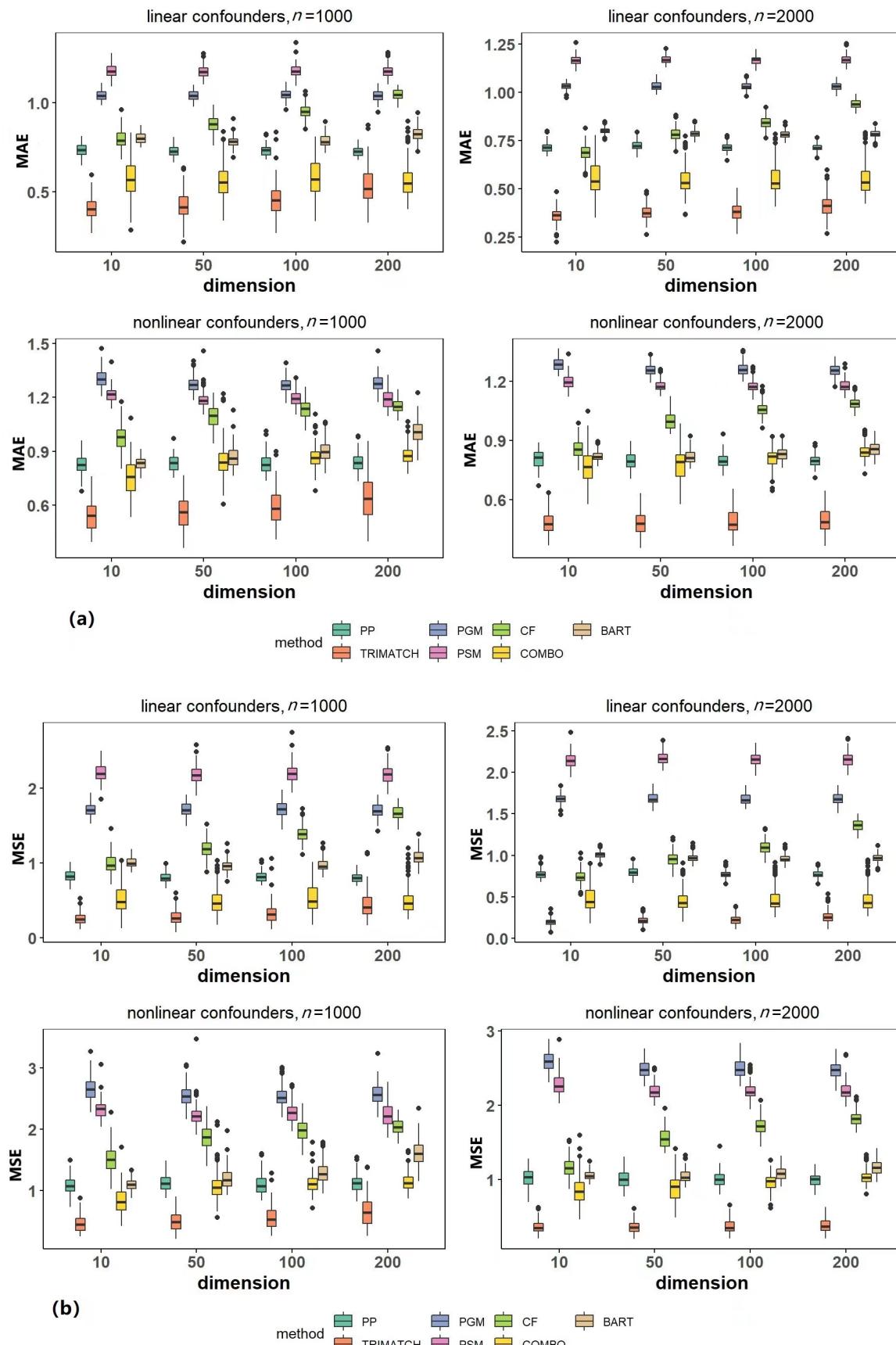
Acknowledging the simulation study results in the previous section, we carry out matching with 0.2 pooled standard deviations of the estimated propensity score on the right part of the graph to balance covariates. We consider the performance of each model in the context of linear and nonlinear confounders. All experiments were repeated 100 times based on datasets with a sample size of 1000 or 2000. For all matching procedures, we use the function *bart* in the R package *dbarts* to estimate the propensity score and the function *glmnet* in the R package *glmnet* to estimate the prognostic score.

The box plots in Fig. 2 depict the prediction performance

**Table 1.** Matching performance with different caliper settings.

	original	a	b	c	d	e	f
PAIRS	–	<b>1026.43</b>	824.83	755.14	991.39	991.39	755.14
TP	1.16	<b>0.08</b>	0.21	0.16	0.11	0.20	0.15
EP	1.59	0.45	<b>0.14</b>	0.21	0.49	0.14	0.54
MAHAL	10.50	<b>2.32</b>	4.39	3.41	3.03	3.57	2.77
std.diff.X1	0.47	<b>0.04</b>	0.09	0.11	0.04	0.09	0.06
std.diff.X2	1.03	<b>0.03</b>	0.17	0.09	0.05	0.15	0.09
std.diff.X3	0.05	<b>0.02</b>	0.05	0.04	0.03	0.04	0.03
std.diff.X4	0.06	<b>0.02</b>	0.05	0.05	0.03	0.04	0.03
std.diff.X5	0.06	<b>0.03</b>	0.06	0.05	0.04	0.04	0.04

The differences in key covariates and propensity scores between matched pairs are standardized by the pooled standard deviation of the respective variable. For each evaluation metric, the best experimental performance under different caliper settings is highlighted in bold.



**Fig. 2.** Comparison between different methods. (a) Mean absolute errors of 100 simulations for various methods. (b) Mean squared errors of 100 simulations for various methods

of different estimators under different sample sizes and data generation models. We use the mean absolute error and mean squared error to measure the accuracy of the model estimates. For both linear confounders and nonlinear confounders, among all the methods, TRIMATCH has the lowest errors in the evaluation metrics of mean absolute error and mean square error, and TRIMATCH's variance is also relatively low, indicating that the proposed method has the optimal performance in the stability and accuracy of estimation. With increasing dimension and complexity of the data model, TRIMATCH's accuracy decreases, but only by a small magnitude, and the standard deviation of the estimation error also maintains a low level, which reflects the robustness of our method. Compared with the fixed ratio matching methods, the error estimated by the variable ratio matching method based on the minimum average cost flow algorithm is smaller, which reflects the superiority of the variable ratio matching method. Compared with the COMBO method, TRIMATCH eliminates the influence of irrelevant factors in the measurement of matching distance and considers the balance between covariates, which reduces the estimation bias compared to the COMBO method and improves the stability of the estimation results.

## 5 Application

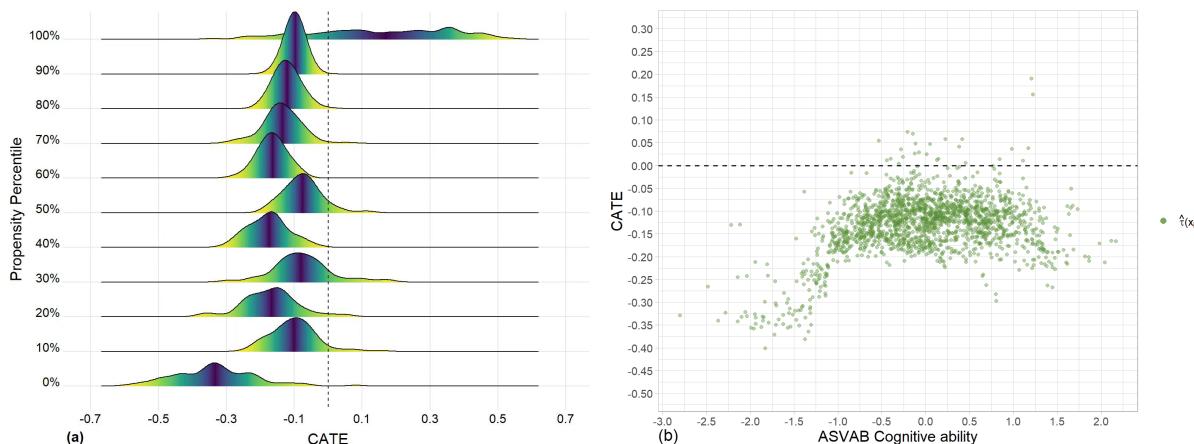
In this section, we apply TRIMATCH to a real-world dataset to study the heterogeneous treatment effects of colleges on reducing low-wage work, which was previously analyzed by Brand et al.<sup>[38]</sup>. The data are derived from the National Longitudinal Survey of Youth 1979 (NLSY79), and we use the same samples and variables as Brand et al.<sup>[38]</sup>. Following their study, the outcome variable of interest  $Y_i$  is the proportion of time in low-wage work over the person's career, while the treatment  $W_i$  denotes whether the individual has completed college by age 25. The covariates  $X$  include ① sociodemographic factors (male, Black, Hispanic, southern residence at age 14, rural residence at age 14), ② family background factors (parents' household income, fathers' highest education, mothers' highest education, father upper-white collar

occupation, two-parent family at age 14, sibship size), ③ cognitive and psychosocial factors (cognitive ability ASVAB, high school college preparatory program, Rotter's Locus of Control Scale, juvenile delinquency activity scale, educational expectations, educational aspirations, friends' educational aspirations), ④ school factors (school disadvantage scale), and ⑤ family formation factors (marital status at age 18, had a child by age 18). The dataset consists of 1764 observations, with 347 in the treated group ( $W = 1$ ) and 1417 in the controlled group ( $W = 0$ ).

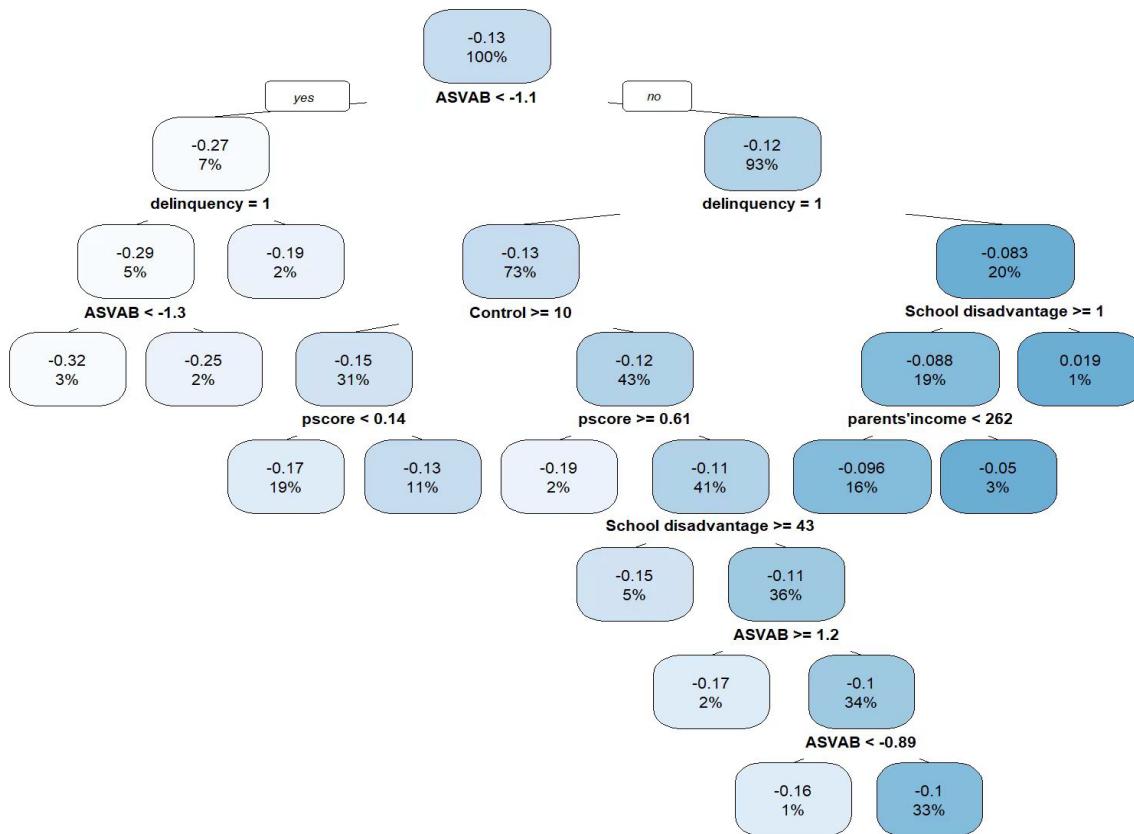
According to the correlation coefficient matrix of covariates, the factor of parental income is associated with several covariates. It is positively related to father's or mothers' highest education, father's upper-white collar occupation, intact family, high school college preparatory program, educational expectations, educational aspirations, friends' educational aspirations, and propensity score. However, parental income is negatively related to sibship size, Rotter's Locus of Control Scale, Black, Hispanic, rural residence at age 14, and marital and childbearing status at age 18.

**Fig. 3a** depicts the estimated heterogeneous treatment effect posterior distributions corresponding to different approximated propensity score percentiles (i.e., to individuals in the sample whose estimated propensity scores are equal or closest to the PS percentiles). Apart from the 100th percentile, the mean of the conditional average treatment effect distributions corresponding to the other percentiles is less than 0, which reflects the positive impact of college completion on reducing the proportion of time in low-wage work. The distribution of the estimated heterogeneous treatment effects near the 0th percentile and 100th quantile of propensity score is more dispersed, indicating greater uncertainty in the estimates of individuals located in the poor overlap region.

When ASVAB cognitive ability is less than  $-0.5$ , **Fig. 3b** shows a roughly positive correlation between conditional average and ASVAB cognitive ability, which means that the lower the cognitive ability, the more significant the effects of college completion on reducing the proportion of low-wage work time. When ASVAB cognitive ability is at least  $-0.5$ , the heterogeneous treatment effects fluctuated more



**Fig. 3.** (a) The graph displays the distribution of estimated heterogeneous treatment effects corresponding to the approximated propensity score percentiles. (b) Scatter plots of estimated treatment effect (averaged over the 100 iterations) against ASVAB cognitive ability.



**Fig. 4.** Prediction model of heterogeneous treatment effects.

consistently and were less than zero for most individuals, implying that college education is beneficial in reducing the proportion of time in low-wage work.

As shown in Fig. 4, from the overall sample perspective, college completion positively impacts reducing the proportion of low-paying jobs in an individual's career. This effect is more pronounced for individuals with lower cognitive abilities, as indicated by an ASVAB score of less than  $-1.1$ . This implies that completing college may be particularly beneficial for individuals with lower cognitive abilities, as it may help them secure higher-paying jobs. Juvenile delinquency activity scale is also an important node for dividing homogeneous subgroups, and those who had delinquent behaviors benefit more from completing college than those without delinquent behaviors.

## 6 Conclusions

In this article, we propose a two-step approach called TRIMATCH for estimating individual treatment effects. First, we construct a variable ratio matching model based on a tripartite graph, utilize the minimum average cost flow algorithm to obtain the optimal matching, and create a “near-randomized” sample by matching to obtain a proxy for the individual treatment effect. Second, based on the matching outcomes, we employ the extreme gradient boosting tree technique to build a prediction model for individual treatment effects.

We validate the effectiveness and accuracy of the

proposed approach through both theoretical analysis and empirical experiments. Our experimental findings demonstrate that TRIMATCH effectively enhances the quality of matching not only by reducing the average distance within the matching set but also by improving the balance of covariates between the treatment and control groups after matching. Furthermore, our approach surpasses published alternative methods in terms of the accuracy of individual treatment effect estimation.

The variable ratio matching problem is commonly reformulated as a minimum average cost flow problem based on bipartite graphs. The minimum average cost flow algorithm offers advantages over the minimum cost flow algorithm because the former can produce more feasible flows for a lower average cost. Tripartite graph matching surpasses bipartite graph matching in its ability to address multiobjective matching problems and yields a lower overall matching cost. Our proposed TRIMATCH offers flexibility and generality in the following ways: (i) It does not impose any restrictions on the type of covariates, making it suitable for both discrete and continuous variables. (ii) The matching ratio is variable and can be adapted to the data distribution by selecting an appropriate matching set size. (iii) In the first step of matching process, TRIMATCH can be combined with covariate balancing techniques such as exact matching and refined balance matching to enhance the balance of the matched covariates. In the second step, other machine learning methods can replace the

extreme gradient boosting tree method to construct the prediction model. Compared to other regression-based causal inference methods, matching enables the adjustment of the covariate equilibrium, and the joint use of matching methods and regression methods helps to reduce the model's sensitivity to unobserved confounders.

TRIMATCH exhibits impressive matching quality, but it requires further optimization in terms of time and space complexity. To optimize the matching algorithm, we can start from the network structure and the solution algorithm, respectively, and combine the network sparsity method to simplify the network structure, explore the optimization time of other nonlinear programming solution algorithms and seek the possibility of further optimization. In the second step of the proposed method, the extreme gradient boosting tree can be replaced by other machine learning methods, combined with meta-learners to explore whether there are more accurate individual treatment effects estimation methods. Moreover, our method is currently applicable only to the case of binary treatment variables. However, TRIMATCH can be extended to the case of multiple treatments by mapping different treatment groups to different layers of the network structure, where the number of layers corresponds to that of categories of treatment variables. Furthermore, different weights could be assigned to network layers to prioritize matching tasks accordingly.

## Acknowledgements

This work was supported by the Science and Technology Planning Project of Anhui Province (202106f01050008).

## Conflict of interest

The authors declare that they have no conflict of interest.

## Biographies

**Yun Cai** is currently a postgraduate at the School of Management, University of Science and Technology of China. Her research mainly focuses on causal inference.

**Shuguang Zhang** is currently a Professor at the School of Management, University of Science and Technology of China (USTC). He received his Ph.D. degree in Statistics from USTC in 1992. His research mainly focuses on stochastic partial differential equations, backward stochastic differential equations, mathematical finance, and financial engineering.

## References

- [1] Chantrill L A, Nagrial A M, Watson C, et al. Precision medicine for advanced pancreas cancer: The individualized molecular pancreatic cancer therapy (IMPaCT) trial. *Clinical Cancer Research*, 2015, 21 (9): 2029–2037.
- [2] Sun W, Wang P, Yin D, et al. Causal inference via sparse additive models with application to online advertising. *Proceedings of the AAAI Conference on Artificial Intelligence*, 2015, 29 (1): 297–303.
- [3] Athey S, Imbens G W. The state of applied econometrics: Causality and policy evaluation. *Journal of Economic Perspectives*, 2017, 31 (2): 3–32.
- [4] Wager S, Athey S. Estimation and inference of heterogeneous treatment effects using random forests. *Journal of the American Statistical Association*, 2018, 113 (523): 1228–1242.
- [5] Richard Hahn P, Murray J S, Carvalho C M. Bayesian regression tree models for causal inference: Regularization, confounding, and heterogeneous effects (with Discussion). *Bayesian Analysis*, 2020, 15 (3): 965–1056.
- [6] Stuart E A. Matching methods for causal inference: A review and a look forward. *Statistical Science*, 2010, 25 (1): 1–21.
- [7] Gao Z, Hastie T, Tibshirani R. Assessment of heterogeneous treatment effect estimation accuracy via matching. *Statistics in Medicine*, 2021, 40 (17): 3990–4013.
- [8] Long M, Sun L, Li Q.  $k$ -Resolution sequential randomization procedure to improve covariates balance in a randomized experiment. *Statistics in Medicine*, 2021, 40 (25): 5534–5546.
- [9] Künzle S R, Sekhon J S, Bickel P J, et al. Metalearners for estimating heterogeneous treatment effects using machine learning. *Proceedings of the National Academy of Sciences*, 2019, 116 (10): 4156–4165.
- [10] Curth A, van der Schaar M. Nonparametric estimation of heterogeneous treatment effects: From theory to learning algorithms. In: Proceedings of the 24th International Conference on Artificial Intelligence and Statistics. San Diego, CA: PMLR, 2021: 1810–1818.
- [11] Nie X, Wager S. Quasi-oracle estimation of heterogeneous treatment effects. *Biometrika*, 2021, 108 (2): 299–319.
- [12] Zhang B, Small D S, Lasater K B, et al. Matching one sample according to two criteria in observational studies. *Journal of the American Statistical Association*, 2023, 118: 1140–1151.
- [13] Iacus S M, King G, Porro G. Causal inference without balance checking: Coarsened exact matching. *Political Analysis*, 2012, 20: 1–24.
- [14] Rubin D B. Matching to remove bias in observational studies. *Biometrics*, 1973, 29 (1): 159–183.
- [15] Rosenbaum P R, Rubin D B. The central role of the propensity score in observational studies for causal effects. *Biometrika*, 1983, 70 (1): 41–55.
- [16] Rubin D B. Using propensity scores to help design observational studies: Application to the tobacco litigation. *Health Services and Outcomes Research Methodology*, 2001, 2 (3): 169–188.
- [17] Hansen B B. The prognostic analogue of the propensity score. *Biometrika*, 2008, 95 (2): 481–488.
- [18] Leacy F P, Stuart E A. On the joint use of propensity and prognostic scores in estimation of the average treatment effect on the treated: A simulation study. *Statistics in Medicine*, 2014, 33 (20): 3488–3508.
- [19] Antonelli J, Cefalu M, Palmer N, et al. Doubly robust matching estimators for high dimensional confounding adjustment. *Biometrics*, 2018, 74 (4): 1171–1179.
- [20] Rosenbaum P R, Rubin D B. Reducing bias in observational studies using subclassification on the propensity score. *Journal of the American Statistical Association*, 1984, 79 (387): 516–524.
- [21] Wooldridge J M. Should instrumental variables be used as matching variables? *Research in Economics*, 2016, 70 (2): 232–237.
- [22] Rosenbaum P R. Optimal matching for observational studies.

- Journal of the American Statistical Association*, **1989**, *84* (408): 1024–1032.
- [23] Zubizarreta J, Keele L. Optimal multilevel matching in clustered observational studies: A case study of the effectiveness of private schools under a large-scale voucher system. *Journal of the American Statistical Association*, **2017**, *112* (518): 547–560.
- [24] Pimentel S D, Kelz R R. Optimal tradeoffs in matched designs comparing US-trained and internationally trained surgeons. *Journal of the American Statistical Association*, **2022**, *115* (532): 1675–1688.
- [25] Yu R, Rosenbaum P R. Directional penalties for optimal matching in observational studies. *Biometrics*, **2019**, *75* (4): 1380–1390.
- [26] Morucci M, Orlandi V, Roy S, et al. Adaptive hyperbox matching for interpretable individualized treatment effect estimation. In: Proceedings of the 36th Conference on Uncertainty in Artificial Intelligence (UAI). Toronto, Canada: PMLR, **2020**: 1089–1098.
- [27] Hansen B B, Klopfer S O. Optimal full matching and related designs via network flows. *Journal of Computational and Graphical Statistics*, **2006**, *15* (3): 609–627.
- [28] Pimentel S D, Kelz R R, Silber J H, et al. Large, sparse optimal matching with refined covariate balance in an observational study of the health outcomes produced by new surgeons. *Journal of the American Statistical Association*, **2015**, *110* (510): 515–527.
- [29] Rubin D B. Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of Educational Psychology*, **1974**, *66* (5): 688–701.
- [30] Robinson P M. Root-N-consistent semiparametric regression. *Econometrica*, **1988**, *56*: 931–954.
- [31] Glazerman S, Levy D M, Myers D. Nonexperimental versus experimental estimates of earnings impacts. *The Annals of the American Academy of Political and Social Science*, **2003**, *589* (1): 63–93.
- [32] Pearl J. On a class of bias-amplifying variables that endanger effect estimates. In: Proceedings of the Twenty-Sixth Conference on Uncertainty in Artificial Intelligence. Arlington, VA: AUAI Press, **2010**: 417–424.
- [33] Chen Y L. The minimal average cost flow problem. *European Journal of Operational Research*, **1995**, *81* (3): 561–570.
- [34] Brito M R, Chávez E L, Quiroz A J, et al. Connectivity of the mutual  $k$ -nearest-neighbor graph in clustering and outlier detection. *Statistics & Probability Letters*, **1997**, *35* (1): 33–42.
- [35] Korte B, Vygen J. Combinatorial Optimization: Theory and Algorithms. Berlin: Springer, **2011**.
- [36] Ye S S, Chen Y, Padilla O H M. Non-parametric interpretable score based estimation of heterogeneous treatment effects. arXiv:2110.02401, **2021**.
- [37] Chipman H A, George E I, McCulloch R E. BART: Bayesian additive regression trees. *The Annals of Applied Statistics*, **2010**, *4*: 266–298.
- [38] Brand J E, Xu J, Koch B, et al. Uncovering sociological effect heterogeneity using machine learning. arXiv: 1909.09138, **2019**.

## Appendix

### A.1 Proof of Proposition 3.2

Let  $\bar{Y}_{II_i} = \sum_{j \in II_i} \frac{Y_j}{|II_i|}$ ,  $\bar{\tau}(X_j) = \sum_{j \in II_i} \frac{\tau(X_j)}{|II_i|}$  for simplicity. If  $k_i$  is the number of neighbors of individual  $i$ , then  $v_{ij} \sim N\left(0, \frac{k_i+1}{k_i} \sigma^2\right)$ .

The proof of inequality (4) is as follows.

$$\begin{aligned} E(\bar{\tau}(X_i) - \tau(X_i)|X_i) &= \\ E[W_i(Y_i - \bar{Y}_j) + (1 - W_i)(\bar{Y}_j - Y_i) - \tau(X_i)|X_i] &= E[(2W_i - 1)(b_{ij} + v_{ij}) + (1 - W_i)\left(\sum_{j \in II_i} \frac{\tau(X_j) - \tau(X_i)}{|II_i|}\right)|X_i] = \\ E[(2W_i - 1)b_{ij}|X_i] + E[(1 - W_i)(\bar{\tau}(X_j) - \tau(X_i))|X_i] &\leq (2e(X_i) - 1)b_{ij} + (1 - e(X_i))L\delta_{ij}^{(k)}. \end{aligned} \quad (\text{A1})$$

The proof of inequality (5) is as follows.

$$\begin{aligned} E((\hat{\tau}(X_i) - \tau(X_i))^2|X_i) &= \\ E[(\hat{\tau}(X_i) - \bar{\tau}(X_i) + \bar{\tau}(X_i) - \tau(X_i))^2|X_i] &= E[(\hat{\tau}(X_i) - \bar{\tau}(X_i))^2|X_i] + 2E[(\hat{\tau}(X_i) - \bar{\tau}(X_i))(\bar{\tau}(X_i) - \tau(X_i))|X_i] + E[(\bar{\tau}(X_i) - \tau(X_i))^2|X_i]. \end{aligned} \quad (\text{A2})$$

By the Hölder inequality, we can obtain

$$E((\hat{\tau}(X_i) - \tau(X_i))^2|X_i) \leq \epsilon^2 + 2\epsilon \sqrt{E[(\hat{\tau}(X_i) - \tau(X_i))^2|X_i]} + E[(\hat{\tau}(X_i) - \tau(X_i))^2|X_i] = (\sqrt{E[(\hat{\tau}(X_i) - \tau(X_i))^2|X_i]} + \epsilon)^2. \quad (\text{A3})$$

$$\begin{aligned} E[(\hat{\tau}(X_i) - \tau(X_i))^2|X_i] &= \\ E[(W_i(Y_i - \bar{Y}_j) + (1 - W_i)(\bar{Y}_j - Y_i) - \tau(X_i))^2|X_i] &= E[((2W_i - 1)(b_{ij} + v_{ij}) + (1 - W_i)\left(\sum_{j \in II_i} \frac{\tau(X_j) - \tau(X_i)}{|II_i|}\right))^2|X_i] = \\ E[(b_{ij} + v_{ij})^2|X_i] + E[(1 - W_i)(\bar{\tau}(X_j) - \tau(X_i))^2|X_i] - 2E[(1 - W_i)(b_{ij} + v_{ij})(\bar{\tau}(X_j) - \tau(X_i))|X_i] &= \\ b_{ij}^2 + \frac{k_i+1}{k_i} \sigma^2 + (1 - e(X_i))(\bar{\tau}(X_j) - \tau(X_i))^2 - 2b_{ij}(1 - e(X_i))(\bar{\tau}(X_j) - \tau(X_i)) &= \\ b_{ij}^2 + \frac{k_i+1}{k_i} \sigma^2 + (1 - e(X_i))[(\bar{\tau}(X_j) - \tau(X_i))^2 - 2b_{ij}(\bar{\tau}(X_j) - \tau(X_i))] &= \\ b_{ij}^2 + \frac{k_i+1}{k_i} \sigma^2 + (1 - e(X_i))[(\bar{\tau}(X_j) - \tau(X_i)) - b_{ij}]^2 - b_{ij}^2 &. \end{aligned} \quad (\text{A4})$$

By the Cauchy-Schwarz inequality, we continue to obtain

$$E[(\bar{\tau}(X_i) - \tau(X_i))^2 | X_i] \leq b_{ij}^2 + \frac{k_i + 1}{k_i} \sigma^2 + (1 - e(X_i)) [2(\bar{\tau}(X_j) - \tau(X_i))^2 + 2b_{ij}^2 - b_{ij}^2] \leq b_{ij}^2 + \frac{k_i + 1}{k_i} \sigma^2 + (1 - e(X_i))(2L^2 \delta_{ij}^{(k)} + b_{ij}^2). \quad (\text{A5})$$

Summarizing (A2) to (A5), we have

$$E((\hat{\tau}(X_i) - \tau(X_i))^2 | X_i) \leq (\sqrt{E[(\bar{\tau}(X_i) - \tau(X_i))^2 | X_i]} + \epsilon)^2 \leq \{[b_{ij}^2 + \frac{k_i + 1}{k_i} \sigma^2 + (1 - e(X_i))(2L^2 \delta_{ij}^{(k)} + b_{ij}^2)]^{\frac{1}{2}} + \epsilon\}^2. \quad (\text{A6})$$

## A.2 100 Monte Carlo simulations with different penalty coefficients $\lambda$

We compared the matching performance under the linear and nonlinear confounder data models using different  $\lambda$ ; the results are shown in the Tables A1 and A2. The coefficient  $\lambda$  represents the tradeoff between matching proximity and covariance balance. The tables illustrate that when  $\lambda \leq 10$ , increasing the  $\lambda$  value reduces the difference in estimated and true propensity scores between the matched treatment and control groups. Meanwhile, the gap of each key covariate between the matched groups increases, and the average Mahalanobis distance of the key covariates also increases. However, when  $\lambda > 10$ , the matching results become insensitive to the penalty coefficient's magnitude.

Upon comparing the matching sets selected by  $\mathcal{G}_1$  and  $\mathcal{G}_2$  for  $\lambda > 10$ , we observe a high degree of overlap. In other words, the matching sets with closer covariate matching distances tend to have closely aligned propensity scores, resulting in a relatively balanced distribution of covariates within the matched treatment and control groups. This result may be attributed to the underlying data generation model, where the covariates linked to the treatment variable represent a subset of the treatment variable's impact factors. The distance measure used in  $\mathcal{G}_1$  incorporates pertinent information related to the propensity score in  $\mathcal{G}_2$ , leading to a matching outcome that is not overly sensitive to the value of  $\lambda$ . Notably, when  $\lambda$  equals 1, the matching method already demonstrates high performance, leaving limited room for improvement. Additionally, in our matching design, we impose a fixed cost on edges that violate the caliper setting in  $\mathcal{G}_2$ . Consequently, when  $\lambda$  is too large, the initial cost of each edge is also large, and the fixed penalty cost has little impact on the cost of each edge in  $\mathcal{G}_2$ .

Nevertheless, the analysis highlights a substantial decrease in the distance observed between the covariance and propensity score of the control group after matching, which also reflects that the tripartite graph model proficiently accounts for both matching proximity and covariance balance.

**Table A1.** Matching performance with different  $\lambda$  under linear confounders.

$\lambda$	original	0.01	0.1	1	10	100	1000
PAIRS	–	773.21	994.48	1025.88	<b>1026.43</b>	1026.43	1026.43
TP	1.16	0.12	0.09	0.07	<b>0.07</b>	0.08	0.08
EP	1.59	0.56	0.51	0.47	0.45	0.45	<b>0.45</b>
MAHAL	10.5	<b>2.04</b>	2.14	2.20	2.30	2.36	2.37
std.diff.X1	0.47	0.04	0.03	<b>0.03</b>	0.03	0.04	0.04
std.diff.X2	1.03	0.08	0.05	0.03	<b>0.02</b>	0.03	0.03
std.diff.X3	0.05	0.01	0.01	<b>0.01</b>	0.02	0.02	0.02
std.diff.X4	0.06	0.01	<b>0.01</b>	0.01	0.02	0.02	0.02
std.diff.X5	0.06	0.05	0.04	0.03	<b>0.03</b>	0.03	0.03

The note is the same as that of Table 1.

**Table A2.** Matching performance with different  $\lambda$  under nonlinear confounders.

$\lambda$	original	0.01	0.1	1	10	100	1000
PAIRS	–	813.55	857.65	961.98	969.53	<b>969.7</b>	969.7
TP	0.65	<b>0.03</b>	0.03	0.09	0.13	0.16	0.17
EP	1.3	0.55	0.5	0.36	0.29	0.28	<b>0.28</b>
MAHAL	10.12	<b>1.52</b>	1.55	1.69	1.85	2	2.04
std.diff.X1	0.56	<b>0.02</b>	0.03	0.07	0.11	0.13	0.14
std.diff.X2	0.28	<b>0.02</b>	0.02	0.04	0.06	0.07	0.08
std.diff.X3	0.12	<b>0.01</b>	0.01	0.02	0.03	0.04	0.04
std.diff.X4	0.06	<b>0.01</b>	0.01	0.02	0.02	0.03	0.03
std.diff.X5	0.06	0.03	<b>0.03</b>	0.03	0.04	0.04	0.04

The note is the same as that of Table 1.