

BOOTSTRAP INFERENCE FOR QUANTILE TREATMENT EFFECTS IN
RANDOMIZED EXPERIMENTS WITH MATCHED PAIRS

By

Liang Jiang, Xiaobin Liu, Peter C.B. Phillips, and Yichong Zhang

August 2020

COWLES FOUNDATION DISCUSSION PAPER NO. 2249



COWLES FOUNDATION FOR RESEARCH IN ECONOMICS
YALE UNIVERSITY
Box 208281
New Haven, Connecticut 06520-8281

<http://cowles.yale.edu/>

Bootstrap Inference for Quantile Treatment Effects in Randomized Experiments with Matched Pairs*

Liang Jiang[†] Xiaobin Liu[‡] Peter C.B. Phillips[§] Yichong Zhang[¶]

*We thank Kengo Kato and Yu-Chin Hsu for their valuable comments and suggestions. We are grateful to Esther Duflo and Cynthia Kinnan for providing the data used in the empirical application. Yichong Zhang acknowledges the financial support from Singapore Ministry of Education Tier 2 grant under grant MOE2018-T2-2-169 and the Lee Kong Chian fellowship. Any and all errors are our own.

[†]Fudan University. E-mail address: jiangliang@fudan.edu.cn.

[‡]The corresponding author. School of Economics, Academy of Financial Research, and Institute for Fiscal Big-Data & Policy of Zhejiang University. E-mail address: liuxiaobin@zju.edu.cn.

[§]Yale University, University of Auckland, University of Southampton, and Singapore Management University. E-mail address: peter.phillips@yale.edu

[¶]Singapore Management University. E-mail address: yczhang@smu.edu.sg.

Abstract

This paper examines methods of inference concerning quantile treatment effects (QTEs) in randomized experiments with matched-pairs designs (MPDs). We derive the limit distribution of the QTE estimator under MPDs, highlighting the difficulties that arise in analytical inference due to parameter tuning. We show that the naive weighted bootstrap fails to approximate the limit distribution of the QTE estimator under MPDs because it ignores the dependence structure within the matched pairs. To address this difficulty we propose two bootstrap methods that can consistently approximate the limit distribution: the gradient bootstrap and the weighted bootstrap of the inverse propensity score weighted (IPW) estimator. The gradient bootstrap is free of tuning parameters but requires knowledge of the pair identities. The weighted bootstrap of the IPW estimator does not require such knowledge but involves one tuning parameter. Both methods are straightforward to implement and able to provide pointwise confidence intervals and uniform confidence bands that achieve exact limiting coverage rates. We demonstrate their finite sample performance using simulations and provide an empirical application to a well-known dataset in microfinance.

Keywords: Bootstrap inference, matched pairs, quantile treatment effect, randomized control trials

JEL codes: C14, C21

1 Introduction

Matched-pairs designs (MPDs) have recently seen widespread and increasing use in various randomized experiments conducted by economists. By MPD we mean a randomization scheme that first pairs units based on the closeness of their baseline covariates and then randomly assigns one unit in the pair to be treated. In development economics, researchers routinely pair villages, neighborhoods, micro-enterprises, or townships in their experiments (Banerjee, Duflo, Glennerster, and Kinnan, 2015; Crepon, Devoto, Duflo, and Pariente, 2015; Glewwe, Park, and Zhao, 2016; Groh and McKenzie, 2016). In labor economics, especially in the field of education, researchers pair schools or students to evaluate the effects of various education interventions (Angrist and Lavy, 2009; Beuermann, Cristia, Cueto, Malamud, and Cruzaguayo, 2015; Fryer, 2017; Fryer, Devi, and Holden, 2017; Bold, Kimenyi, Mwabu, Nganga, and Sandefur, 2018; Fryer, 2018). Bruhn and McKenzie (2009) surveyed leading experts in development field experiments and reported that 56% of them explicitly match pairs of observations on baseline characteristics.

Researchers often use randomized experiments to estimate quantile treatment effects (QTEs) as well as average treatment effects (ATEs). Quantile effects can capture heterogeneity in both the sign and magnitude of treatment effects, which may vary according to position within the distribution of outcomes. A common practice in conducting inference on QTEs is to use bootstrap rather than analytical methods because the latter usually require tuning parameters in implementation. However, the treatment assignment in MPDs introduces negative *dependence* because exactly half of the units are treated. Standard bootstrap inference procedures that rely on cross-sectional *independence* are therefore conservative and lack power. This difficulty raises the question of how to conduct bootstrap inference for QTEs in MPDs in a manner that mitigates these shortcomings.

The present paper addresses this question by showing that both the gradient bootstrap and the weighted bootstrap of the inverse propensity score weighted (IPW) estimator can

consistently approximate the limit distribution of the original QTE estimator under MPDs, thereby eliminating asymptotic size distortion in inference. In particular, for testing null hypotheses that the QTEs equal some pre-specified values involving single or multiple quantile indexes (or some pre-specified function over a compact set of quantile indexes), the usual pointwise confidence interval or uniform confidence band constructed by using the corresponding bootstrap standard errors achieves a limiting rejection probability under the null equal to the nominal level.

Our starting point is to derive the limit distribution of the two-sample-difference type QTE estimator in MPDs uniformly over a compact set of quantile indexes. Analytic computation of the variance of the QTE estimator using this limit theory requires estimation of two infinite dimensional nuisance parameters. By implication two tuning parameters are needed for every quantile index of interest. This procedure is inevitably cumbersome and provides the motivation to develop bootstrap methods of inference that reduce the need for tuning parameters.

As noted above, observations under MPDs are generally dependent within the pairs, whereas the usual bootstrap counterparts are asymptotically independent conditional on the data. In accord with this contrasting property of the bootstrap we show that the naive weighted bootstrap fails to approximate the limit distribution of the QTE estimator. Consequently, usual bootstrap tests of the null hypothesis that the QTE equals a pre-specified value are conservative and lack power.

To tackle this shortcoming we propose a gradient bootstrap method and show that it can consistently approximate the limit distribution of the QTE estimator under MPDs uniformly over a compact set of quantile indexes. Hagemann (2017) proposed using the gradient bootstrap for the cluster-robust inference in linear quantile regression models. Like Hagemann (2017), we rely on the gradient bootstrap to avoid estimating the Hessian matrix that involves the infinite-dimensional nuisance parameters. The gradient bootstrap procedure is therefore free of tuning parameters. On the other hand and differing from Hagemann

(2017), we construct a specific perturbation of the score based on pair and adjacent pairs of observations, which can capture the dependence structure in the original data.

To implement our gradient bootstrap method, researchers need to know the identities of pairs. Such information may not be available when they are using an experiment that was run by someone else in the past and the randomization procedure may not have been fully described. To address this issue, we propose a weighted bootstrap of the IPW QTE estimator, which can be implemented without such knowledge. We show that such a bootstrap can consistently estimate the asymptotic distribution of the QTE estimator under MPDs. There is a cost to not using information about pair identities as the method requires one tuning parameter for the nonparametric estimation of the propensity score. In spite of this additional cost, this weighted bootstrap method still has an advantage over direct analytic inference because practical implementation of the latter requires more than one tuning parameter.

The contributions in the present paper relate to other recent research. Bai, Shaikh, and Romano (2019) first pointed out that in MPDs the two-sample t -test for the null hypothesis that the ATE equals a pre-specified value is conservative. They then proposed adjusting the standard error of the estimator and studied the validity of the permutation test. This paper complements those results by considering the QTEs and by developing new methods of bootstrap inference. Unlike the permutation test, our methods of bootstrap inference do not require studentization, which is cumbersome in the QTE context. In addition, our weighted bootstrap method complements their results by providing a way to perform inference relating to both ATEs and QTEs when pair identities are unknown. In other work, Bai (2019) investigated the optimality of MPDs in randomized experiments. Zhang and Zheng (2020) considered bootstrap inference under covariate-adaptive randomization. A key difference in our contribution is that in MPDs the number of strata is proportional to the sample size, whereas in covariate-adaptive randomization that number is fixed. In consequence, the present work uses fundamentally different asymptotic arguments and bootstrap methods from those employed by Zhang and Zheng (2020). The present paper also fits within a

growing literature that studies inference in randomized experiments (e.g., Hahn, Hirano, and Karlan (2011), Athey and Imbens (2017), Abadie, Chingos, and West (2018), Bugni, Canay, and Shaikh (2018), Tabord-Meehan (2018), and Bugni, Canay, and Shaikh (2019), among others).

The remainder of the paper is organized as follows. Section 2 describes the model setup and notation. Section 3 develops the asymptotic properties of our QTE estimator. In Section 4 we study the naive weighted bootstrap, the gradient bootstrap, and the weighted bootstrap of the IPW estimator. Section 5 provides computational details and recommendations for practitioners. Section 6 reports simulation results. Section 7 gives an empirical application of our methods of bootstrap inference to the data in Banerjee et al. (2015), examining both the ATEs and QTEs of microfinance on the take-up rates of microcredit. Section 8 concludes. Proofs of all results are in the supplement appendix.

2 Setup and Notation

Denote the potential outcomes for treated and control groups as $Y(1)$ and $Y(0)$, respectively. The treatment status is written as A , where $A = 1$ means treated and $A = 0$ means untreated. The researcher can only observe $\{Y_i, X_i, A_i\}_{i=1}^{2n}$ where $Y_i = Y_i(1)A_i + Y_i(0)(1 - A_i)$, and $X_i \in \mathfrak{R}^{d_x}$ is a collection of baseline covariates, where d_x is the dimension of X . The parameter of interest is the τ th QTE, denoted as

$$q(\tau) = q_1(\tau) - q_0(\tau),$$

where $q_1(\tau)$ and $q_0(\tau)$ are the τ th quantiles of $Y(1)$ and $Y(0)$, respectively. The testing problems of interest involve single, multiple, or even a continuum of quantile indexes, as in

the following null hypotheses

$$\mathcal{H}_0 : q(\tau) = \underline{q} \quad \text{versus} \quad q(\tau) \neq \underline{q},$$

$$\mathcal{H}_0 : q(\tau_1) - q(\tau_2) = \underline{q} \quad \text{versus} \quad q(\tau_1) - q(\tau_2) \neq \underline{q}, \text{ and}$$

$$\mathcal{H}_0 : q(\tau) = \underline{q}(\tau) \quad \forall \tau \in \Upsilon \quad \text{versus} \quad q(\tau) \neq \underline{q}(\tau) \text{ for some } \tau \in \Upsilon,$$

for some pre-specified value \underline{q} or function $\underline{q}(\tau)$, where Υ is some compact subset of $(0, 1)$.

The units are grouped into pairs based on the closeness of their baseline covariates, which will be made clear next. We denote the pairs of units as

$$(\pi(2j - 1), \pi(2j)) \text{ for } j \in [n],$$

where $[n] = \{1, \dots, n\}$ and π is a permutation of $2n$ units based on $\{X_i\}_{i=1}^{2n}$ as specified in Assumption 1(iv) below. Within the pair, one unit is randomly assigned to treatment and the other to control. Specifically, we make the following assumption on the data generating process (DGP) and the treatment assignment rule.

Assumption 1. (i) $\{Y_i(1), Y_i(0), X_i\}_{i=1}^{2n}$ is *i.i.d.*

$$(ii) \{Y_i(1), Y_i(0)\}_{i=1}^{2n} \perp\!\!\!\perp \{A_i\}_{i=1}^{2n} | \{X_i\}_{i=1}^{2n}.$$

(iii) Conditionally on $\{X_i\}_{i=1}^{2n}$, $(\pi(2j - 1), \pi(2j))$ for $j \in [n]$, are *i.i.d.* and each uniformly distributed over the values in $\{(1, 0), (0, 1)\}$.

$$(iv) \frac{1}{n} \sum_{j=1}^n \|X_{\pi(2j)} - X_{\pi(2j-1)}\|_2^r \xrightarrow{P} 0 \text{ for } r = 1, 2.$$

Assumption 1 is used in Bai et al. (2019) to which we refer readers for more discussion. In Assumption 1(iv), $\|\cdot\|_2$ denotes Euclidean distance. However, all our results hold if $\|\cdot\|_2$ is replaced by any distance that is equivalent to it, such as L_∞ distance, L_1 distance, and the Mahalanobis distance when all the eigenvalues of the covariance matrix are bounded and

bounded away from zero. Later in Section 4 and following Assumption 4 we provide two cases for which Assumption 1(iv) holds.

3 Estimation

Let $\hat{q}_1(\tau)$ and $\hat{q}_0(\tau)$ be the τ th percentiles of outcomes in the treated and control groups, respectively. Then, the τ th QTE estimator we consider is just

$$\hat{q}(\tau) = \hat{q}_1(\tau) - \hat{q}_0(\tau).$$

To facilitate further analysis and motivate our bootstrap procedure, we note that $\hat{q}(\tau)$ can be equivalently computed by direct quantile regression. Let

$$(\hat{\beta}_0(\tau), \hat{\beta}_1(\tau)) = \arg \min_b \sum_{i=1}^{2n} \rho_\tau(Y_i - \dot{A}'b),$$

where $\dot{A}_i = (1, A_i)^T$ and $\rho_\tau(u) = u(\tau - 1\{u \leq 0\})$. Then, $\hat{q}(\tau) = \hat{\beta}_1(\tau)$ and $\hat{q}_0(\tau) = \hat{\beta}_0(\tau)$.

Assumption 2. For $a = 0, 1$, denote $F_a(\cdot)$, $F_a(\cdot|x)$, $f_a(\cdot)$, and $f_a(\cdot|x)$ as the CDF of $Y_i(a)$, the conditional CDF of $Y_i(a)$ given $X_i = x$, the PDF of $Y_i(a)$, and the conditional PDF of $Y_i(a)$ given $X_i = x$, respectively.

(i) $f_a(q_a(\tau))$ is bounded and bounded away from zero uniformly over $\tau \in \Upsilon$, and $f_a(q_a(\tau)|x)$ is uniformly bounded for $(x, \tau) \in \text{Supp}(X) \times \Upsilon$.

(ii) There exists a function $C(x)$ such that

$$\sup_{\tau \in \Upsilon} |f_a(q_a(\tau) + v|x) - f_a(q_a(\tau)|x)| \leq C(x)|v| \quad \text{and} \quad \mathbb{E}C(X_i) < \infty.$$

(iii) Let \mathcal{N}_0 be a neighborhood of 0. Then, there exists a constant C such that for any

$x, x' \in \text{Supp}(X)$

$$\sup_{\tau \in \Upsilon, v \in \mathcal{N}_0} |f_a(q_a(\tau) + v|x') - f_a(q_a(\tau) + v|x)| \leq C \|x' - x\|_2$$

and

$$\sup_{\tau \in \Upsilon, v \in \mathcal{N}_0} |F_a(q_a(\tau) + v|x) - F_a(q_a(\tau) + v|x')| \leq C \|x - x'\|_2.$$

Assumption 2(i) is the standard regularity condition widely assumed in quantile estimation. The Lipschitz conditions in Assumptions 2(ii) and 2(iii) are similar in spirit to those assumed in Bai et al. (2019, Assumption 2.1) and ensure that units that are “close” in terms of their baseline covariates are suitably comparable. For $a = 0, 1$, let $m_{a,\tau}(x, q) = \mathbb{E}(\tau - 1\{Y(a) \leq q\} | X = x)$ and $m_{a,\tau}(x) = m_{a,\tau}(x, q_a(\tau))$.

Theorem 3.1. *Suppose Assumptions 1 and 2 hold. Then, uniformly over $\tau \in \Upsilon$,*

$$\sqrt{n}(\hat{q}(\tau) - q(\tau)) \rightsquigarrow \mathcal{B}(\tau),$$

where $\mathcal{B}(\tau)$ is a Gaussian process with covariance kernel $\Sigma(\cdot, \cdot)$ such that

$$\begin{aligned} \Sigma(\tau, \tau') &= \frac{\min(\tau, \tau') - \tau\tau' - \mathbb{E}m_{1,\tau}(X)m_{1,\tau'}(X)}{f_1(q_1(\tau))f_1(q_1(\tau'))} + \frac{\min(\tau, \tau') - \tau\tau' - \mathbb{E}m_{0,\tau}(X)m_{0,\tau'}(X)}{f_0(q_0(\tau))f_0(q_0(\tau'))} \\ &\quad + \frac{1}{2} \mathbb{E} \left(\frac{m_{1,\tau}(X)}{f_1(q_1(\tau))} - \frac{m_{0,\tau}(X)}{f_0(q_0(\tau))} \right) \left(\frac{m_{1,\tau'}(X)}{f_1(q_1(\tau'))} - \frac{m_{0,\tau'}(X)}{f_0(q_0(\tau'))} \right). \end{aligned}$$

Several remarks are in order. First, the asymptotic variance of $\hat{q}(\tau)$ under MPDs is

$$\Sigma(\tau, \tau) = \frac{\tau - \tau^2 - \mathbb{E}m_{1,\tau}^2(X)}{f_1^2(q_1(\tau))} + \frac{\tau - \tau^2 - \mathbb{E}m_{0,\tau}^2(X)}{f_0^2(q_0(\tau))} + \frac{1}{2} \mathbb{E} \left(\frac{m_{1,\tau}(X)}{f_1(q_1(\tau))} - \frac{m_{0,\tau}(X)}{f_0(q_0(\tau))} \right)^2.$$

Further note that the asymptotic variance of $\hat{q}(\tau)$ under simple random sampling is

$$\Sigma^\dagger(\tau, \tau) = \frac{\tau - \tau^2}{f_1^2(q_1(\tau))} + \frac{\tau - \tau^2}{f_0^2(q_0(\tau))}.$$

It is clear that

$$\Sigma^\dagger(\tau, \tau) - \Sigma(\tau, \tau) = \frac{1}{2} \mathbb{E} \left(\frac{m_{1,\tau}(X)}{f_1(q_1(\tau))} + \frac{m_{0,\tau}(X)}{f_0(q_0(\tau))} \right)^2 \geq 0. \quad (3.1)$$

Equality in the last expression holds when both $m_{1,\tau}(X)$ and $m_{0,\tau}(X)$ are zero, which implies that X is irrelevant to the τ th quantiles of $Y(0)$ and $Y(1)$.

Second, the asymptotic variance $\Sigma(\tau, \tau)$ coincides with the semiparametric efficiency bound of the QTE estimator established in Firpo (2007) and Donald and Hsu (2014) for observational data under unconfoundedness.¹ Hahn (1998) pointed out that, even in the case of simple random sampling, to achieve the semiparametric efficiency bound one needs to use the IPW estimator with a nonparametrically estimated propensity score. We view the MPD as an alternative to achieving such efficiency without nonparametric estimation.²

Third, to provide an analytic estimate of the asymptotic variance $\Sigma(\tau, \tau)$ it is necessary at least to estimate the infinite dimensional nuisance parameters $f_1(q_1(\tau))$ and $f_0(q_0(\tau))$, which requires two tuning parameters. Hence, if a researcher is interested in testing a null hypothesis that involves G quantile indexes, $2G$ tuning parameters are needed to estimate $2G$ densities, a cumbersome task in practical work; and to construct a uniform confidence band for the QTE analytically, two tuning parameters are needed at each grid point of the quantile indexes. Moreover, if pair identities are unknown, analytic methods of inference potentially require nonparametric estimation of the quantities $m_{a,\tau}(\cdot)$ for $a = 0, 1$ as well. There are other practical difficulties. Nonparametric estimation is sometimes sensitive to the choice of tuning parameters and rule-of-thumb tuning parameter selection may not be appropriate for every data generating process (DGP) or every quantile. Use of cross-validation in selecting the tuning parameters is possible in principle but in practice time-consuming.

¹The propensity score is just a constant of $1/2$.

²Whether the efficiency bound remains the same under MPDs is still an open question and is an interesting topic for future research.

These practical difficulties of analytic methods of inference provide a strong motivation to investigate bootstrap inference procedures are much less reliant on tuning parameters.

4 Bootstrap Inference

This section examines three bootstrap inference procedures for the QTEs in MPDs. We first show that a naive weighted bootstrap method fails to approximate the limit distribution of the QTE estimator derived in Section 3. We then propose two bootstrap methods that can consistently estimate the asymptotic distribution of the QTE estimator.

4.1 Naive Weighted Bootstrap Inference

We first consider the naive weighted bootstrap estimators of $\hat{\beta}_0(\tau)$ and $\hat{\beta}_1(\tau)$. Let

$$(\hat{\beta}_0^w(\tau), \hat{\beta}_1^w(\tau)) = \arg \min_b \sum_{i=1}^{2n} \xi_i \rho_\tau(Y_i - A'b),$$

where ξ_i is the bootstrap weight defined in the next assumption.

Assumption 3. *Suppose $\{\xi_i\}_{i=1}^{2n}$ is a sequence of nonnegative i.i.d. random variables with unit expectation and variance and a sub-exponential upper tail.*

Denote $\hat{q}^w(\tau) = \hat{\beta}_1^w(\tau)$ and recall that $\hat{q}(\tau) = \hat{\beta}_1(\tau)$.

Theorem 4.1. *If Assumptions 1–3 hold, then conditional on the data and uniformly over $\tau \in \Upsilon$,*

$$\sqrt{n}(\hat{q}^w(\tau) - \hat{q}(\tau)) \rightsquigarrow \mathcal{B}^w(\tau),$$

where $\mathcal{B}^w(\tau)$ is a Gaussian process with covariance kernel $\Sigma^\dagger(\cdot, \cdot)$ such that

$$\Sigma^\dagger(\tau, \tau') = \frac{\min(\tau, \tau') - \tau\tau'}{f_1(q_1(\tau))f_1(q_1(\tau'))} + \frac{\min(\tau, \tau') - \tau\tau'}{f_0(q_0(\tau))f_0(q_0(\tau'))}.$$

Three remarks are in order. First, $\Sigma^\dagger(\tau, \tau')$ is just the covariance kernel of the QTE estimator when simple random sampling (instead of the MPD) is used as the treatment assignment rule. It follows that the naive weighted bootstrap fails to approximate the limit distribution of $\hat{q}(\tau)$ ($\hat{\beta}_1(\tau)$). The intuition is straightforward. Given the data, the bootstrap weights are i.i.d. and thus unable to mimic the cross-sectional dependence in the original sample.

Second, it is possible to consider the conventional nonparametric bootstrap in which the bootstrap sample is generated from the empirical distribution of the data. If the observations are i.i.d., van der Vaart and Wellner (1996, Section 3.6) showed that the conventional bootstrap is first-order equivalent to a weighted bootstrap with Poisson(1) weights. However, in the current setting, $\{A_i\}_{i \in [2n]}$ are dependent. It is technically challenging to show rigorously that the above equivalence still holds and this is left for future research.

Third, an alternative procedure is to bootstrap the pairs of observations, i.e., to use the same bootstrap weights for observations indexed by $\pi(2j - 1)$ and $\pi(2j)$. But such a bootstrap alone is unable to mimic the dependence structure in the original sample. In fact, the gradient bootstrap procedure proposed below follows this idea and uses the same weight for the observations in the same pair to construct the score $S_{n,1}^*$ defined in (4.5). But in order to construct a final score that can mimic the dependence in the data we need an extra score $S_{n,2}^*$, which is defined in (4.6).

4.2 Gradient Bootstrap Inference

We now approximate the asymptotic distribution of the QTE estimator via the gradient bootstrap. Let $u = \sqrt{n}(b - \beta(\tau))$ be a localizing estimation error parameter. From the derivations in Theorem 3.1, we see that

$$\sqrt{n}(\hat{\beta}(\tau) - \beta(\tau)) = \arg \min_u \sum_{i=1}^{2n} \rho_\tau \left(Y_i - \dot{A}^T \beta(\tau) - \frac{\dot{A}^T u}{\sqrt{n}} \right),$$

where

$$\sum_{i=1}^{2n} \left[\rho_{\tau}(Y_i - \dot{A}^T \beta(\tau) - \frac{\dot{A}^T u}{\sqrt{n}}) - \rho_{\tau}(Y_i - \dot{A}^T \beta(\tau)) \right] \approx -u' \begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix} S_n(\tau) + \frac{u^T Q(\tau) u}{2}, \quad (4.1)$$

$$S_n(\tau) = \begin{pmatrix} \sum_{i=1}^{2n} \frac{A_i}{\sqrt{n}} (\tau - 1\{Y_i(1) \leq q_1(\tau)\}) \\ \sum_{i=1}^{2n} \frac{(1-A_i)}{\sqrt{n}} (\tau - 1\{Y_i(0) \leq q_0(\tau)\}) \end{pmatrix},$$

and

$$Q(\tau) = \begin{pmatrix} f_1(q_1(\tau)) + f_0(q_0(\tau)) & f_1(q_1(\tau)) \\ f_1(q_1(\tau)) & f_1(q_1(\tau)) \end{pmatrix}.$$

Minimizing the right side of (4.1) gives

$$\sqrt{n}(\hat{\beta}(\tau) - \beta(\tau)) \approx Q^{-1}(\tau) \begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix} S_n(\tau). \quad (4.2)$$

The gradient bootstrap proposes to perturb the objective function by some random error $S_n^*(\tau)$, which will be specified later. This error in turn perturbs the score function $S_n(\tau)$. The corresponding bootstrap estimator $\hat{\beta}^*(\tau)$ solves the following optimization problem

$$\hat{\beta}^*(\tau) = \arg \min_b \sum_{i=1}^{2n} \rho_{\tau}(Y_i - \dot{A}'b) - \sqrt{n}b^T \begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix} S_n^*(\tau). \quad (4.3)$$

By a change of variable and (4.1) we have

$$\sqrt{n}(\hat{\beta}^*(\tau) - \beta(\tau)) \approx \arg \min_u -u' \begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix} [S_n(\tau) + S_n^*(\tau)] + \frac{u^T Q(\tau) u}{2},$$

which implies

$$\sqrt{n}(\hat{\beta}^*(\tau) - \beta(\tau)) \approx Q^{-1}(\tau) \begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix} [S_n(\tau) + S_n^*(\tau)]. \quad (4.4)$$

Taking the difference between (4.2) and (4.4), we have

$$\sqrt{n}(\hat{\beta}^*(\tau) - \hat{\beta}(\tau)) \approx Q^{-1}(\tau) \begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix} S_n^*(\tau).$$

The second element of $\hat{\beta}^*(\tau)$ in (4.3) is the bootstrap version of the QTE estimator, which is denoted $\hat{q}^*(\tau)$. By solving (4.3) we avoid estimating the Hessian $Q(\tau)$, which involves the infinite-dimensional nuisance parameters. Then, for the gradient bootstrap to consistently approximate the limit distribution of the original estimator $\hat{\beta}(\tau)$, we need only construct $S_n^*(\tau)$ in such a way that its weak limit given the data coincides with that of the original score $S_n(\tau)$.

Accordingly, we now show how to specify $S_n^*(\tau)$. Let $\{\eta_j\}_{j=1}^n$ and $\{\hat{\eta}_k\}_{k=1}^{\lfloor n/2 \rfloor}$ be two mutually independent i.i.d. sequences of standard normal random variables. Use the indexes $(j, 1), (j, 0)$ to denote the indexes in $(\pi(2j - 1), \pi(2j))$ with $A = 1$ and $A = 0$, respectively. For example, if $A_{\pi(2j)} = 1$ and $A_{\pi(2j-1)} = 0$, then $(j, 1) = \pi(2j)$ and $(j, 0) = \pi(2j - 1)$. Similarly, use indexes $(k, 1), \dots, (k, 4)$ to denote the first index in $(\pi(4k - 3), \dots, \pi(4k))$ with $A = 1$, the first index with $A = 0$, the second index with $A = 1$, and the second index with $A = 0$, respectively. Now let

$$S_n^*(\tau) = \frac{S_{n,1}^*(\tau) + S_{n,2}^*(\tau)}{\sqrt{2}},$$

where

$$S_{n,1}^*(\tau) = \frac{1}{\sqrt{n}} \left(\frac{\sum_{j=1}^n \eta_j (\tau - 1\{Y_{(j,1)} \leq \hat{q}_1(\tau)\})}{\sum_{j=1}^n \eta_j (\tau - 1\{Y_{(j,0)} \leq \hat{q}_0(\tau)\})} \right) \quad (4.5)$$

and

$$S_{n,2}^*(\tau) = \frac{1}{\sqrt{n}} \left(\frac{\sum_{k=1}^{\lfloor n/2 \rfloor} \hat{\eta}_k [(\tau - 1\{Y_{(k,1)} \leq \hat{q}_1(\tau)\}) - (\tau - 1\{Y_{(k,3)} \leq \hat{q}_1(\tau)\})]}{\sum_{k=1}^{\lfloor n/2 \rfloor} \hat{\eta}_k [(\tau - 1\{Y_{(k,2)} \leq \hat{q}_0(\tau)\}) - (\tau - 1\{Y_{(k,4)} \leq \hat{q}_0(\tau)\})]} \right). \quad (4.6)$$

In Section 5 we show how to compute the bootstrap estimator $\hat{\beta}^*(\tau)$ directly from the sub-gradient condition of (4.3). This method avoids the optimization inherent in (4.3) and computation is fast. The following assumption imposes the condition that baseline covariates in adjacent pairs are also ‘close’.

Assumption 4. *Suppose that $\frac{1}{n} \sum_{k=1}^{\lfloor n/2 \rfloor} \|X_{(k,l)} - X_{(k,l')}\|_2^r \xrightarrow{P} 0$ for $r = 1, 2$ and $l, l' \in [4]$.*

Assumption 4 and Assumption 1(iv) are jointly equivalent to Bai et al. (2019, Assumption 2.4). We refer readers to Bai et al. (2019) for further discussion of this assumption. In particular, Bai et al. (2019, Theorems 4.1 and 4.2) established two cases under which both Assumption 4 and Assumption 1(iv) hold. We repeat their results below for completeness.

Case (1). Suppose X is a scalar and $\mathbb{E}X^2 < \infty$. Let π be any permutation of $2n$ elements such that $X_{\pi(1)} \leq \dots \leq X_{\pi(2n)}$. Then, both Assumption 4 and Assumption 1(iv) hold.

Case (2). Suppose $\text{Supp}(X) \subset [0, 1]^{d_x}$. Let $\check{\pi}$ be any permutation of $2n$ elements minimizing $\frac{1}{n} \sum_{j=1}^n \|X_{\check{\pi}(2j-1)} - X_{\check{\pi}(2j)}\|_2$, let $\bar{X}_j = \frac{1}{2} (X_{\check{\pi}(2j-1)} + X_{\check{\pi}(2j)})$, and let $\bar{\pi}$ be any permutation of n elements minimizing $\frac{1}{n} \sum_{j=1}^n \|\bar{X}_{\bar{\pi}(j)} - \bar{X}_{\bar{\pi}(j-1)}\|_2$. Then, the permutation π with $\pi(2j) = \check{\pi}(2\bar{\pi}(j))$ and $\pi(2j-1) = \check{\pi}(2\bar{\pi}(j) - 1)$ satisfies Assumption 4 and Assumption 1(iv).

Denote $\hat{q}^*(\tau) = \hat{\beta}_1^*(\tau)$ and recall that $\hat{q}(\tau) = \hat{\beta}_1(\tau)$. We now have the following result.

Theorem 4.2. *Suppose Assumptions 1, 2, and 4 hold. Then, conditional on the data and uniformly over $\tau \in \Upsilon$, $\sqrt{n}(\hat{q}^*(\tau) - \hat{q}(\tau)) \rightsquigarrow \mathcal{B}(\tau)$, where $\mathcal{B}(\tau)$ is the same Gaussian process defined in Theorem 3.1.*

Three remarks on Theorem 4.2 are in order. First, the bootstrap estimator $\hat{q}^*(\tau)$ has the following objectives: (i) to avoid estimating densities; and (ii) to mimick the distribution of the original estimator $\hat{\beta}(\tau)$ under MPDs. Objective (i) relates to the Hessian (Q) and (ii) to the score (S_n) of the quantile regression. The gradient bootstrap provide a flexible approach to achieve both goals.

Second, Bai et al. (2019) showed that adjacent pairs can be used to construct a valid standard error for the ATE estimator under MPDs. Our approach follows their lead and bootstraps pairs and adjacent pairs of units. Theorem 4.2 shows that the limit distribution of the resulting bootstrapped perturbation $S_n^*(\tau)$ given that the data can consistently approximate that of the original score $S_n(\tau)$ uniformly over $\tau \in \Upsilon$. For inference concerning the ATE, it is not necessary to use the gradient bootstrap as the Hessian does not contain any infinite-dimensional nuisance parameters. In fact, the way we compute the perturbation $S_n^*(\tau)$ leads directly to a variance estimator $\hat{\nu}_n^2$ for the ATE estimator $\hat{\Delta} = \frac{1}{n} \sum_{j=1}^n (Y_{(j,1)} - Y_{(j,0)})$, where

$$\hat{\nu}_n^2 = \frac{1}{2n} \sum_{j=1}^n (Y_{(j,1)} - Y_{(j,0)} - \hat{\Delta})^2 + \frac{1}{2n} \sum_{k=1}^{\lfloor n/2 \rfloor} [(Y_{(k,1)} - Y_{(k,3)}) - (Y_{(k,2)} - Y_{(k,4)})]^2.$$

By some manipulation, one can show that $\hat{\nu}_n^2$ is numerically the same as the estimate used in the adjusted t -test of Bai et al. (2019, Section 3.3).

Third, to implement the gradient bootstrap, researchers need to know pair identities. That information may not be available when the base experiment was run by others and the randomization procedure not fully detailed. In such cases, we propose bootstrapping the IPW estimator of the QTE, whose validity is established in the next section.

4.3 Weighted Bootstrap of Inverse Propensity Score Weighted Estimator

As indicated in Section 3, the QTE estimator under MPDs achieves the semiparametric efficiency bound established for independent observational data. If we use independent bootstrap weights and seek to maintain efficiency, we need to bootstrap an estimator that can achieve the semiparametric efficiency bound under independent data. As pointed out by Hahn (1998) and Firpo (2007), the IPW estimator with a nonparametrically estimated propensity score satisfies this requirement. Accordingly, we now propose a weighted bootstrap version of the IPW estimator to approximate the limit distribution of the QTE estimator in MPDs.

The sieve method is used to estimate the propensity score. Let $b(X)$ be the K -dimensional sieve basis on X and \hat{A}_i the estimated propensity score for the i th individual. Then,

$$\hat{A}_i = b(X_i)' \hat{\theta}, \quad (4.7)$$

where ξ_i is the bootstrap weight defined in Assumption 3 and $\hat{\theta} = \arg \min_{\theta} \sum_{i=1}^{2n} \xi_i (A_i - b(X_i)' \theta)^2$.

Because the true propensity score is $1/2$, by setting the first component of $b(X)$ to unity, we have $1/2 = b'(X)\theta_0$ where $\theta_0 = (0.5, 0, \dots, 0)^T$. The linear probability model for the propensity score is correctly specified. It is possible to use sieve logistic regression to compute the propensity score, as done by Hirano, Imbens, and Ridder (2003), Firpo (2007), and Donald and Hsu (2014). The main benefit of using logistic regression is to guarantee that the estimated propensity score lies between zero and one. For simplicity, we use a linear sieve regression here.

The weighted bootstrap IPW estimator can be computed as

$$\hat{q}_{ipw}^w(\tau) = \hat{q}_{ipw,1}^w(\tau) - \hat{q}_{ipw,0}^w(\tau),$$

where

$$\hat{q}_{ipw,1}^w(\tau) = \arg \min_q \sum_{i=1}^{2n} \frac{\xi_i A_i}{\hat{A}_i} \rho_\tau(Y_i - q) \quad \text{and} \quad \hat{q}_{ipw,0}^w(\tau) = \arg \min_q \sum_{i=1}^{2n} \frac{\xi_i (1 - A_i)}{1 - \hat{A}_i} \rho_\tau(Y_i - q). \quad (4.8)$$

Assumption 5. (i) *The support of X is compact. The first component of $b(X)$ is 1.*

(ii) $\max_{k \in [K]} \mathbb{E} b_k^2(X_i) \leq \bar{C} < \infty$ for some constant $\bar{C} > 0$. $\sup_{x \in \text{Supp}(X)} \|b(x)\|_2 = \zeta(K)$.

(iii) $K^2 \zeta(K)^2 \log(n) = o(n)$.

(iv) *With probability approaching one, there exist constants \underline{C} and \bar{C} such that*

$$0 < \underline{C} \leq \lambda_{\min} \left(\frac{1}{n} \sum_{i=1}^{2n} \xi_i b(X_i) b(X_i)' \right) \leq \lambda_{\max} \left(\frac{1}{n} \sum_{i=1}^{2n} \xi_i b(X_i) b(X_i)' \right) \leq \bar{C} < \infty,$$

where $\lambda_{\min}(\mathcal{M})$ and $\lambda_{\max}(\mathcal{M})$ denote the minimum and maximum eigenvalues of matrix \mathcal{M} .

(v) *There exist $\gamma_1(\tau) \in \mathbb{R}^K$ and $\gamma_0(\tau) \in \mathbb{R}^K$ such that*

$$B_{a,\tau}(x) = m_{a,\tau}(x) - b'(x) \gamma_a(\tau), \quad a = 0, 1,$$

and $\sup_{a=0,1, \tau \in \Upsilon, x \in \text{Supp}(X)} |B_{a,\tau}(x)| = o(1/\sqrt{n})$.

Two remarks are in order. First, requiring X to have a compact support is common in nonparametric sieve estimation. Second, the quantity $\zeta(K)$ depends on the choice of basis functions. For example, $\zeta(K) = O(K^{1/2})$ for B-splines and $\zeta(K) = O(K)$ for power series³. Taking B-splines as an example, Assumption 5(iii) requires $K = o(n^{1/3})$. Assumption 5(iv) is standard because $K \ll n$. Assumption 5(v) requires that the approximation error of $m_{a,\tau}(x)$ via a linear sieve function is sufficiently small. For instance, suppose $m_{a,\tau}(x)$ is s-times

³See Chen (2007) for a full discussion of the sieve method.

continuously differentiable in x with all derivatives uniformly bounded by some constant \bar{C} , then $\sup_{a=0,1,\tau \in \Upsilon, x \in \text{Supp}(X)} |B_{a,\tau}(x)| = O(K^{-s/d_x})$. Assumptions 5(iii) and 5(v) imply that $K = n^h$ for some $h \in (d_x/(2s), 1/3)$, which implicitly requires $s > 3d_x/2$. The choice of K reflects the usual bias-variance trade-off and is the only tuning parameter that researchers need to specify when implementing this bootstrap method.

Theorem 4.3. *Suppose Assumptions 1–3 and 5 hold, then conditionally on the data and uniformly over $\gamma \in \Upsilon$, $\sqrt{n}(\hat{q}_{ipw}^w(\tau) - \hat{q}(\tau)) \rightsquigarrow \mathcal{B}(\tau)$, where $\mathcal{B}(\tau)$ is the same Gaussian process as defined in Theorem 3.1.*

The benefit of the weighted bootstrap of the IPW estimator is that it does not require knowledge of the pair identities. The cost is that we have to nonparametrically estimate the propensity score, which requires one tuning parameter and is subject to the usual curse of dimensionality. Nonetheless, we still prefer this bootstrap method of inference to the analytic approach. Analytic estimation of the standard error of the QTE estimator without the knowledge of pair identities requires nonparametric estimation of $\{m_{a,\tau}(X), f_a(q_a(\tau))\}_{a=0,1}$, which involves four tuning parameters. The number of tuning parameters further increases with the number of quantile indexes involved in the null hypothesis and uniform confidence bands for QTE over τ requires $4G$ tuning parameters for grid size G . By contrast, implementation of the weighted bootstrap for the IPW estimator requires estimation of the propensity score only once, requiring use of a single tuning parameter.

Inference concerning the *ATE* in MPDs can also be accomplished via the weighted bootstrap of the IPW *ATE* estimator. A similar argument shows that such a bootstrap can consistently approximate the asymptotic distribution of the *ATE* estimator under MPDs. This result complements that established by Bai et al. (2019) because it provides a way to make inferences about the *ATE* in MPDs when information on pair identities is unavailable. That pair identity information is required by Bai et al. (2019) in computing standard errors for their adjusted *t*-test.

5 Computation and Guidance for Practitioners

5.1 Computation of the Gradient Bootstrap

In practice, the order of pairs in the dataset is usually arbitrary and does not satisfy Assumption 4. To apply the gradient bootstrap, researchers first need to re-order the pairs. For the j th pair with units indexed by $(j, 1)$ and $(j, 0)$ in the treatment and control groups, let $\bar{X}_j = \frac{1}{2}\{X_{(j,1)} + X_{(j,0)}\}$. Then, let $\bar{\pi}$ be any permutation of n elements that minimizes

$$\frac{1}{n} \sum_{j=1}^n \|\bar{X}_{\bar{\pi}(j)} - \bar{X}_{\bar{\pi}(j-1)}\|_2.$$

The pairs are re-ordered by indexes $\bar{\pi}(1), \dots, \bar{\pi}(n)$. With an abuse of notation, we still index the pairs after re-ordering by $1, \dots, n$. Note that the original QTE estimator $\hat{q}(\tau) = \hat{q}_1(\tau) - \hat{q}_0(\tau)$ is invariant to the re-ordering.

For the bootstrap sample, we directly compute $\hat{\beta}^*(\tau)$ from the sub-gradient condition of (4.3). Specifically, we compute $\hat{\beta}_0^*(\tau)$ as $Y_{(h_0)}^0$ and $\hat{q}^*(\tau) \equiv \hat{\beta}_1^*(\tau)$ as $Y_{(h_1)}^1 - Y_{(h_0)}^0$, where $Y_{(h_0)}^0$ and $Y_{(h_1)}^1$ are the h_0 th and h_1 th order statistics of outcomes in the treatment and control groups, respectively,⁴ and h_0 and h_1 are two integers satisfying

$$n\tau + T_{n,a}^*(\tau) + 1 \geq h_a \geq n\tau + T_{n,a}^*(\tau), \quad a = 0, 1, \quad (5.1)$$

with

$$\begin{aligned} \begin{pmatrix} T_{n,1}^*(\tau) \\ T_{n,0}^*(\tau) \end{pmatrix} &= \sqrt{n} S_n^*(\tau) = \frac{1}{\sqrt{2}} \left[\begin{pmatrix} \sum_{j=1}^n \eta_j (\tau - 1\{Y_{(j,1)} \leq \hat{q}_1(\tau)\}) \\ \sum_{j=1}^n \eta_j (\tau - 1\{Y_{(j,0)} \leq \hat{q}_0(\tau)\}) \end{pmatrix} \right. \\ &\quad \left. + \begin{pmatrix} \sum_{k=1}^{\lfloor n/2 \rfloor} \hat{\eta}_k [(\tau - 1\{Y_{(k,1)} \leq \hat{q}_1(\tau)\}) - (\tau - 1\{Y_{(k,3)} \leq \hat{q}_1(\tau)\})] \\ \sum_{k=1}^{\lfloor n/2 \rfloor} \hat{\eta}_k [(\tau - 1\{Y_{(k,2)} \leq \hat{q}_0(\tau)\}) - (\tau - 1\{Y_{(k,4)} \leq \hat{q}_0(\tau)\})] \end{pmatrix} \right]. \end{aligned}$$

⁴We assume $Y_{(1)}^a \leq \dots \leq Y_{(n)}^a$ for $a = 0, 1$.

As the probability of $n\tau + T_{n,a}^*(\tau)$ being an integer is zero, h_a is uniquely defined with probability one.

We summarize the steps in the bootstrap procedure as follows.

1. Re-order the pairs.
2. Compute the original estimator $\hat{q}(\tau) = \hat{q}_1(\tau) - \hat{q}_0(\tau)$.
3. Let B be the number of bootstrap replications. Let \mathcal{G} be a grid of quantile indexes. For $b \in [B]$, generate $\{\eta_j\}_{j \in [n]}$ and $\{\hat{\eta}_k\}_{k \in [n/2]}$. Compute $\hat{q}^{*b}(\tau) = Y_{(h_1)}^1 - Y_{(h_0)}^0$ for $\tau \in \mathcal{G}$, where h_0 and h_1 are computed in (5.1). Obtain $\{\hat{q}^{*b}(\tau)\}_{\tau \in \mathcal{G}}$.
4. Repeat the above step for $b \in [B]$ and obtain B bootstrap estimators of the QTE, denoted as $\{\hat{q}^{*b}(\tau)\}_{b \in [B], \tau \in \mathcal{G}}$.

5.2 Computation of the Weighted Bootstrap of the IPW estimator

We first provide more details on the sieve basis. Let $b(x) \equiv (b_1(x), \dots, b_K(x))'$, where $\{b_k(\cdot)\}_{k=1}^K$ are K basis functions of a linear sieve space \mathcal{B} . Given that all d_x elements of X are continuously distributed, the sieve space \mathcal{B} can be constructed as follows.

1. For each element $X^{(l)}$ of X , $l = 1, \dots, d_x$, let \mathcal{B}_l be the univariate sieve space of dimension J_n . One example of \mathcal{B}_l is the linear span of the J_n dimensional polynomials given by

$$\mathcal{B}_l = \left\{ \sum_{k=0}^{J_n} \alpha_k x^k, x \in \text{Supp}(X^{(l)}), \alpha_k \in \mathfrak{R} \right\};$$

Another is the linear span of r -order splines with J_n nodes given by

$$\mathcal{B}_l = \left\{ \sum_{k=0}^{r-1} \alpha_k x^k + \sum_{j=1}^{J_n} b_j [\max(x - t_j, 0)]^{r-1}, x \in \text{Supp}(X^{(l)}), \alpha_k, b_j \in \mathfrak{R} \right\},$$

where the grid $-\infty = t_0 \leq t_1 \leq \dots \leq t_{J_n} \leq t_{J_n+1} = \infty$ partitions $\text{Supp}(X^{(l)})$ into

$J_n + 1$ subsets $I_j = [t_j, t_{j+1}) \cap \text{Supp}(X^{(l)})$, $j = 1, \dots, J_n - 1$, $I_0 = (t_0, t_1) \cap \text{Supp}(X^{(l)})$, and $I_{J_n} = (t_{J_n}, t_{J_n+1}) \cap \text{Supp}(X^{(l)})$.

- Let \mathcal{B} be the tensor product of $\{\mathcal{B}_l\}_{l=1}^{d_x}$, which is defined as a linear space spanned by the functions $\prod_{l=1}^{d_x} g_l$, where $g_l \in \mathcal{B}_l$. The dimension of \mathcal{B} is then $K \equiv d_x J_n$.

Given the sieve basis, we can estimate the propensity score following (4.7). We then obtain $\hat{q}_{ipw,1}^w(\tau)$ and $\hat{q}_{ipw,0}^w(\tau)$ by solving the sub-gradient conditions for the two optimizations in (4.8). Specifically, we have $\hat{q}_{ipw,1}^w(\tau) = Y_{h'_1}$ and $\hat{q}_{ipw,0}^w(\tau) = Y_{h'_0}$, where the indexes h'_0 and h'_1 satisfy $A_{h'_a} = a$, $a = 0, 1$,

$$\tau \left(\sum_{i=1}^{2n} \frac{\xi_i A_i}{\hat{A}_i} \right) - \frac{\xi_{h'_1}}{\hat{A}_{h'_1}} \leq \sum_{i=1}^{2n} \frac{\xi_i A_i}{\hat{A}_i} 1\{Y_i < Y_{h'_1}\} \leq \tau \left(\sum_{i=1}^{2n} \frac{\xi_i A_i}{\hat{A}_i} \right), \quad (5.2)$$

and

$$\tau \left(\sum_{i=1}^{2n} \frac{\xi_i (1 - A_i)}{1 - \hat{A}_i} \right) - \frac{\xi_{h'_0}}{1 - \hat{A}_{h'_0}} \leq \sum_{i=1}^{2n} \frac{\xi_i (1 - A_i)}{1 - \hat{A}_i} 1\{Y_i < Y_{h'_0}\} \leq \tau \left(\sum_{i=1}^{2n} \frac{\xi_i (1 - A_i)}{1 - \hat{A}_i} \right). \quad (5.3)$$

In the implementation, we set $\{\xi_i\}_{i \in [2n]}$ as i.i.d. standard exponential random variables. In this case, all the equalities in (5.2) and (5.3) hold with probability zero. Thus, h'_1 and h'_0 are uniquely defined with probability one.

We summarize the bootstrap procedure as follows.

- Compute the original estimator $\hat{q}(\tau) = \hat{q}_1(\tau) - \hat{q}_0(\tau)$.
- Let B be the number of bootstrap replications. Let \mathcal{G} be a grid of quantile indexes. For $b \in [B]$, generate $\{\xi_i\}_{i \in [2n]}$ as a sequence of i.i.d. exponential random variables. Estimate the propensity score following (4.7). Compute $\hat{q}_{ipw}^{w,b}(\tau) = Y_{h'_1} - Y_{h'_0}$ for $\tau \in \mathcal{G}$, where h'_0 and h'_1 are computed in (5.2) and (5.3), respectively. Obtain $\{\hat{q}_{ipw}^{w,b}(\tau)\}_{\tau \in \mathcal{G}}$.
- Repeat the above step for $b \in [B]$ and obtain B bootstrap estimators of the QTE, denoted as $\{\hat{q}_{ipw}^{w,b}(\tau)\}_{b \in [B], \tau \in \mathcal{G}}$.

For comparison, we also consider the naive weighted bootstrap in our simulations. Its computation follows a procedure similar to the above with only one difference: the nonparametric estimate \hat{A}_i of the propensity score is replaced by the truth, that is, 1/2.

5.3 Bootstrap Confidence Intervals

Given the bootstrap estimates, we discuss how to conduct bootstrap inference for the null hypotheses with single, multiple, and a continuum of quantile indexes. We take the gradient bootstrap as an example. If the IPW bootstrap is used, one can just replace $\{\hat{q}^{*b}(\tau)\}_{b \in [B], \tau \in \mathcal{G}}$ by $\{\hat{q}_{ipw}^{w,b}(\tau)\}_{b \in [B], \tau \in \mathcal{G}}$ in the following cases.

Case (1). We aim to test the single null hypothesis that $\mathcal{H}_0 : q(\tau) = \underline{q}$ vs. $q(\tau) \neq \underline{q}$. Let $\mathcal{G} = \{\tau\}$ in the procedures described above. Further denote $\mathcal{Q}(\nu)$ as the ν th empirical quantile of the sequence $\{\hat{q}^{*b}(\tau)\}_{b \in [B]}$. Let $\alpha \in (0, 1)$ be the significance level. We suggest using the bootstrap estimator to construct the standard error of $\hat{q}(\tau)$ as $\hat{\sigma} = \frac{\mathcal{Q}(0.975) - \mathcal{Q}(0.025)}{C_{0.975} - C_{0.025}}$, where C_μ is the μ th standard normal critical value. Then the valid confidence interval and Wald test using this standard error are

$$CI_1(\alpha) = (\hat{q}(\tau) - C_{1-\alpha/2}\hat{\sigma}, \hat{q}(\tau) + C_{\alpha/2}\hat{\sigma}),$$

and $1\left\{\left|\frac{\hat{q}(\tau) - \underline{q}}{\hat{\sigma}}\right| \geq C_{1-\alpha/2}\right\}$, respectively.

Further denote the standard and percentile bootstrap confidence intervals as CI_2 and CI_3 , respectively, where

$$CI_2(\alpha) = (2\hat{q}(\tau) - \mathcal{Q}(1 - \alpha/2), 2\hat{q}(\tau) - \mathcal{Q}(\alpha/2))$$

and

$$CI_3(\alpha) = (\mathcal{Q}(\alpha/2), \mathcal{Q}(1 - \alpha/2)).$$

Theoretically, CI_1 , CI_2 , and CI_3 are all valid. When $\alpha = 0.05$, CI_1 , CI_2 , and CI_3 are centered at $\hat{q}(\tau)$, $2\hat{q}(\tau) - \frac{1}{2}\{\mathcal{Q}(0.975) + \mathcal{Q}(0.025)\}$, and $\frac{1}{2}\{\mathcal{Q}(0.975) + \mathcal{Q}(0.025)\}$, respectively, but share the same length $\mathcal{Q}(0.975) - \mathcal{Q}(0.025)$. In (unreported) simulations, we found that in small samples, CI_1 usually has the best size control while CI_2 over-rejects and CI_3 under-rejects.

Case (2). We aim to test the null hypothesis that $\mathcal{H}_0 : q(\tau_1) - q(\tau_2) = \underline{q}$ vs. $q(\tau_1) - q(\tau_2) \neq \underline{q}$. In this case, let $\mathcal{G} = \{\tau_1, \tau_2\}$. Further, let $\mathcal{Q}(\nu)$ denote the ν th empirical quantile of the sequence $\{\hat{q}^{*b}(\tau_1) - \hat{q}^{*b}(\tau_2)\}_{b \in [B]}$, and let $\alpha \in (0, 1)$ be the significance level. For the same reason discussed in case (1), we suggest using the bootstrap standard error to construct the valid confidence interval and Wald test as

$$CI_1(\alpha) = (\hat{q}(\tau_1) - \hat{q}(\tau_2) - C_{1-\alpha/2}\hat{\sigma}, \hat{q}(\tau_1) - \hat{q}(\tau_2) + C_{\alpha/2}\hat{\sigma}),$$

and $1\left\{\left|\frac{\hat{q}(\tau_1) - \hat{q}(\tau_2) - \underline{q}}{\hat{\sigma}}\right| \geq C_{1-\alpha/2}\right\}$, respectively, where $\hat{\sigma} = \frac{\mathcal{Q}(0.975) - \mathcal{Q}(0.025)}{C_{0.975} - C_{0.025}}$.

Case (3). We aim to test the null hypothesis that

$$\mathcal{H}_0 : q(\tau) = \underline{q}(\tau) \forall \tau \in \Upsilon \text{ vs. } q(\tau) \neq \underline{q}(\tau) \exists \tau \in \Upsilon.$$

In theory, we should let $\mathcal{G} = \Upsilon$. In practice, we let $\mathcal{G} = \{\tau_1, \dots, \tau_G\}$ be a fine grid of Υ where G should be as large as computationally possible. Further, let $\mathcal{Q}_\tau(\nu)$ denote the ν th empirical quantile of the sequence $\{\hat{q}^{*b}(\tau)\}_{b \in [B]}$ for $\tau \in \mathcal{G}$. Compute the standard error of $\hat{q}(\tau)$ as

$$\hat{\sigma}_\tau = \frac{\mathcal{Q}_\tau(0.975) - \mathcal{Q}_\tau(0.025)}{C_{0.975} - C_{0.025}}.$$

The uniform confidence band with an α significance level is constructed as

$$CB(\alpha) = \{\hat{q}(\tau) - \mathcal{C}_\alpha \hat{\sigma}_\tau, \hat{q}(\tau) + \mathcal{C}_\alpha \hat{\sigma}_\tau : \tau \in \mathcal{G}\},$$

where the critical value \mathcal{C}_α is computed as

$$\mathcal{C}_\alpha = \inf \left\{ z : \frac{1}{B} \sum_{b=1}^B 1 \left\{ \sup_{\tau \in \mathcal{G}} \left| \frac{\hat{q}^{*b}(\tau) - \tilde{q}(\tau)}{\hat{\sigma}_\tau} \right| \leq z \right\} \geq 1 - \alpha \right\}$$

and $\tilde{q}(\tau)$ is first-order equivalent to $\hat{q}(\tau)$ in the sense that $\sup_{\tau \in \Upsilon} |\tilde{q}(\tau) - \hat{q}(\tau)| = o_p(1/\sqrt{n})$. We suggest choosing $\tilde{q}(\tau) = \frac{1}{2}\{\mathcal{Q}_\tau(0.975) + \mathcal{Q}_\tau(0.025)\}$ over other choices such as $\tilde{q}(\tau) = \mathcal{Q}_\tau(0.5)$ and $\tilde{q}(\tau) = \hat{q}(\tau)$ due to its better finite-sample performance. We reject \mathcal{H}_0 at an α significance level if $\underline{q}(\cdot) \notin CB(\alpha)$.

5.4 Practical Recommendations

Our practical recommendations are straightforward. If pair identities are known, we suggest using the gradient bootstrap for inference. If pair identities are unknown, we suggest using the weighted bootstrap of the IPW estimator with a nonparametrically estimated propensity score for inference.

6 Simulation

In this section, we assess the finite-sample performance of the methods discussed in Section 4 with a Monte Carlo simulation study. In all cases, potential outcomes for $a \in \{0, 1\}$ and $1 \leq i \leq 2n$ are generated as

$$Y_i(a) = \mu_a + m_a(X_i) + \sigma_a(X_i) \varepsilon_{a,i}, \quad a = 0, 1, \tag{6.1}$$

where $\mu_a, m_a(X_i), \sigma_a(X_i)$, and $\varepsilon_{a,i}$ are specified as follows. In each of the specifications below, $n \in \{50, 100\}$ and $(X_i, \varepsilon_{0,i}, \varepsilon_{1,i})$ are i.i.d. The number of replications is 10,000. For bootstrap replications we set $B = 5,000$.

Model 1 $X_i \sim \text{Unif}[0, 1]$; $m_0(X_i) = 0$; $m_1(X_i) = 10(X_i^2 - \frac{1}{3})$; $\varepsilon_{a,i} \sim N(0, 1)$ for $a = 0, 1$; $\sigma_0(X_i) = \sigma_0 = 1$ and $\sigma_1(X_i) = \sigma_1$.

Model 2 As in Model 1, but $\sigma_0(X_i) = (1 + X_i^2)$ and $\sigma_1(X_i) = (1 + X_i^2)\sigma_1$.

Model 3 $X_i = (\Phi(V_{i1}), \Phi(V_{i2}))'$, where $\Phi(\cdot)$ is the standard normal cumulative distribution function and

$$V_i \sim N \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \right),$$

$m_0(X_i) = \gamma'X_i - 1$; $m_1(X_i) = m_0(X_i) + 10(\Phi^{-1}(X_{i1})\Phi^{-1}(X_{i2}) - \rho)$; $\varepsilon_{a,i} \sim N(0, 1)$ for $a = 0, 1$; $\sigma_0(X_i) = \sigma_0 = 1$ and $\sigma_1(X_i) = \sigma_1$. We set $\gamma = (1, 1)'$, $\sigma_1 = 1$, $\rho = 0.2$.

Model 4 As in Model 3, but with $\gamma = (1, 4)'$, $\sigma_1 = 2$, $\rho = 0.7$.

Pairs are determined similarly to those in Bai et al. (2019). Specifically, if X_i is a scalar, then pairs are determined by sorting $\{X_i\}_{i \in [2n]}$ as described in Case (1) in Section 4.2. If X_i is multi-dimensional, then the pairs are determined by the permutation π described in Case (2) in Section 4.2, which can be obtained by using the *R* package *nbpMatching*. After forming the pairs, we assign treatment status within each pair through a random draw from the uniform distribution over $\{(0, 1), (1, 0)\}$.

We examine the performance of various tests for ATEs and QTEs at the nominal level $\alpha = 5\%$. For the ATE, we consider the hypothesis that

$$\mathbb{E}(Y(1) - Y(0)) = \text{truth} + \Delta \quad \text{vs.} \quad \mathbb{E}(Y(1) - Y(0)) \neq \text{truth} + \Delta.$$

For the QTE, we consider the hypotheses that

$$q(\tau) = \text{truth} + \Delta \quad \text{vs.} \quad q(\tau) \neq \text{truth} + \Delta,$$

for $\tau = 0.25, 0.5,$ and $0.75,$

$$q(0.25) - q(0.75) = \text{truth} + \Delta \quad \text{vs.} \quad q(0.25) - q(0.75) \neq \text{truth} + \Delta, \quad (6.2)$$

and

$$q(\tau) = \text{truth} + \Delta \quad \forall \tau \in [0.25, 0.75] \quad \text{vs.} \quad q(\tau) \neq \text{truth} + \Delta \quad \exists \tau \in [0.25, 0.75]. \quad (6.3)$$

To illustrate size and power of the tests, we set $\mathcal{H}_0 : \Delta = 0$ and $\mathcal{H}_1 : \Delta = 1/2$. The true value for the ATE is 0, whereas the true values for the QTEs are simulated with a 10,000 sample size and replications. The computational procedures described in Section 5 are followed to perform the bootstrap and calculate the test statistics. To test the single null hypothesis involving one or two quantile indexes, we use the Wald tests specified in Section 5.3. To test the null hypothesis involving a continuum of quantile indexes, we use the uniform confidence band $CB(\alpha)$ defined in Case (3) in the same section.

The results for the ATEs appear in Table 1. Each row presents a different model and each column reports the rejection probabilities for the various methods. The column ‘Naive’ refers to the two-sample t -test and ‘Adj’ refers to the adjusted t -test in Bai et al. (2019); the column ‘IPW’ corresponds to the t -test with standard errors generated by the weighted bootstrap of the IPW ATE estimator. In all cases, we find that (i) the two-sample t -test has rejection probability under \mathcal{H}_0 far below the nominal level and is the least powerful test among the three, and (ii) the adjusted t -test has rejection probability under \mathcal{H}_0 close to the nominal level and is not conservative. These results are consistent with those in Bai et al. (2019). The IPW t -test proposed in this paper has performance similar to the adjusted t -

test.⁵ Under \mathcal{H}_0 , the test has rejection probability close to 5%; under \mathcal{H}_1 , it is more powerful than the naive method and has power similar to the adjusted t -test. These findings indicate that the IPW t -test provides an alternative to the adjusted t -test when pair identities are unknown.

Table 1: The Empirical Size and Power of Tests for ATEs

Model	$\mathcal{H}_0: \Delta = 0$						$\mathcal{H}_1: \Delta = 1/2$					
	$n = 50$			$n = 100$			$n = 50$			$n = 100$		
	Naive	Adj	IPW	Naive	Adj	IPW	Naive	Adj	IPW	Naive	Adj	IPW
1	1.32	5.47	5.44	1.22	5.75	6.00	11.80	29.10	29.44	27.67	49.79	50.46
2	1.85	5.35	5.59	1.64	5.63	5.89	10.43	23.26	24.24	23.72	40.42	41.68
3	1.20	4.76	4.92	0.77	4.68	5.16	1.31	5.66	5.91	1.92	8.13	8.74
4	2.32	6.47	6.01	1.25	5.33	4.74	1.08	5.16	4.35	0.93	5.65	4.89

Notes: The table presents the rejection probabilities for tests of ATEs. The columns ‘Naive’ and ‘Adj’ correspond to the two-sample t -test and the adjusted t -test in Bai et al. (2019), respectively; the column ‘IPW’ corresponds to the t -test using the standard errors estimated by the weighted bootstrap of the IPW ATE estimator.

The results for QTEs are summarized in Tables 2 and 3. Each table has four panels (Models 1-4). Each row in the panel displays the rejection probabilities for the tests using the standard errors estimated by various bootstrap methods. Specifically, the rows ‘Naive weight’, ‘Gradient’, and ‘IPW’ respectively correspond to the results of the naive weighted bootstrap, the gradient bootstrap, and the weighted bootstrap of the IPW QTE estimator.

Table 2 reports empirical size and power of the tests with a single null hypothesis involving one or two quantile indexes. Columns ‘0.25’, ‘0.50’, and ‘0.75’ correspond to tests with quantiles at 25%, 50%, and 75%. Column ‘Dif’ corresponds to the test with null hypothesis (6.2). As expected given Theorem 4.1, the test with standard errors estimated by the naive

⁵Throughout this section, we use B-splines to nonparametrically estimate the propensity score in the weighted bootstrap of the IPW estimator. If $\dim(X_i)=1$, we choose the bases $\{1, X, [\max(X - qx_0, X - qx_{0.5})]^2\}$ where qx_0 and $qx_{0.5}$ are quantiles of X at 0 and 50%, respectively; if $\dim(X_i)=2$, we choose the bases $\{1, \max(X_1 - qx_{1,0}, X_1 - x_{1,0.5}), \max(X_2 - qx_{2,0}, X_2 - qx_{2,0.5}), X_1 X_2\}$. The choices of the sieve basis functions and K are adhoc. It is possible to use data-driven methods to select them but a rigorous analysis of the validity of various data-driven methods is beyond the scope of this paper.

method performs poorly in all cases. It is conservative under \mathcal{H}_0 and lacks power under \mathcal{H}_1 . In contrast, the test using the standard errors estimated by either the gradient bootstrap or the IPW method has a rejection probability under \mathcal{H}_0 that is close to the nominal level in almost all specifications. When the number of pairs is 50, the tests in the ‘Dif’ column constructed based on either the gradient or the IPW method are slightly conservative. Sizes approach the nominal level when n increases to 100.

Table 3 reports empirical size and power of the uniform confidence bands for the hypothesis specified in (6.3) with a grid $\mathcal{G} = \{0.25, 0.27, \dots, 0.47, 0.49, 0.5, 0.51, 0.53, \dots, 0.73, 0.75\}$. The test using standard errors estimated by the naive method has rejection probabilities under \mathcal{H}_0 far below the nominal level in all specifications. In Models 1-2, the test using standard errors estimated by either the gradient bootstrap or the IPW bootstrap yields a rejection probability under \mathcal{H}_0 that is very close to the nominal level even when the number of pairs is as small as 50. Nonetheless, in Models 3-4, the tests constructed based on both methods are conservative when the number of pairs equals 50. When the number of pairs increases to 100, both tests perform much better and have a rejection probability under \mathcal{H}_0 that is close to the nominal level. Under \mathcal{H}_1 , the tests based on both the gradient and IPW methods are more powerful than those based on the naive method.

In summary, the simulation results in Tables 2 and 3 are consistent with the results in Theorems 4.2 and 4.3: both the gradient bootstrap and the IPW bootstrap provide valid pointwise and uniform inference for QTEs under MPDs. The findings also show that when the information on pair identities is unavailable the IPW method continues to provide a sound basis for inference.

Table 2: The Empirical Size and Power of Tests for QTEs

	$\mathcal{H}_0: \Delta = 0$								$\mathcal{H}_1: \Delta = 1/2$							
	$n = 50$				$n = 100$				$n = 50$				$n = 100$			
	0.25	0.50	0.75	Dif	0.25	0.50	0.75	Dif	0.25	0.50	0.75	Dif	0.25	0.50	0.75	Dif
<i>Model 1</i>																
Naive weight	3.00	2.00	2.22	1.98	3.12	2.06	1.93	1.73	16.67	6.05	5.56	3.96	34.93	11.56	8.11	7.35
Gradient	5.13	4.82	4.92	3.66	5.07	5.62	5.30	4.04	23.76	13.03	11.27	8.18	42.92	22.91	17.30	14.57
IPW	5.47	5.31	6.17	4.24	5.26	5.83	5.65	3.95	24.81	13.48	12.12	8.40	43.93	23.33	17.21	13.91
<i>Model 2</i>																
Naive weight	3.08	2.32	2.55	1.96	3.64	2.53	2.08	1.87	14.82	6.54	4.71	3.68	30.29	11.50	7.46	6.88
Gradient	4.57	4.63	4.39	3.44	5.00	5.42	5.28	3.68	19.51	12.25	8.76	6.57	35.38	20.86	14.79	12.25
IPW	4.93	5.12	5.78	4.45	5.17	5.73	5.88	4.00	20.29	12.90	10.40	7.35	36.38	21.53	15.14	12.53
<i>Model 3</i>																
Naive weight	2.11	1.03	2.10	0.92	1.56	1.37	1.58	0.86	4.98	2.85	1.92	0.98	6.57	7.14	1.73	1.43
Gradient	5.24	3.06	3.14	1.76	4.83	4.20	4.27	3.01	9.71	7.43	3.22	2.39	13.80	16.72	5.67	4.40
IPW	4.76	3.19	5.61	2.60	4.77	3.71	4.95	3.02	8.75	7.81	5.35	3.09	13.04	15.42	6.06	4.21
<i>Model 4</i>																
Naive weight	2.59	1.71	1.98	1.65	2.65	1.66	1.55	1.23	6.09	1.94	1.76	1.28	9.85	2.98	1.19	1.18
Gradient	4.75	4.00	3.33	2.82	4.70	4.74	5.06	3.88	9.37	5.76	3.35	2.87	14.67	8.88	5.27	4.25
IPW	3.97	3.97	4.91	3.68	4.23	4.51	5.01	3.48	8.08	5.37	4.79	3.26	13.50	8.33	5.17	3.51

Note: The table presents the rejection probabilities for tests of QTEs. The columns ‘0.25’, ‘0.50’, and ‘0.75’ correspond to tests with quantiles at 25%, 50%, and 75%, respectively; the column ‘Dif’ corresponds to the test with the null hypothesis specified in (6.2). The rows ‘Naive weight’, ‘Gradient’, and ‘IPW’ correspond to the results of the naive weighted bootstrap, the gradient bootstrap, and the weighted bootstrap of the IPW QTE estimator, respectively.

Table 3: The Empirical Size and Power of Uniform Inferences for QTEs

	$\mathcal{H}_0: \Delta = 0$		$\mathcal{H}_1: \Delta = 1/2$	
	$n = 50$	$n = 100$	$n = 50$	$n = 100$
<i>Model 1</i>				
Naive weight	1.07	1.52	7.50	18.12
Gradient	4.08	4.64	17.88	33.30
IPW	4.49	4.94	16.30	32.40
<i>Model 2</i>				
Naive weight	1.37	1.85	6.73	16.50
Gradient	3.66	4.57	14.30	27.64
IPW	4.25	4.91	14.27	27.47
<i>Model 3</i>				
Naive weight	0.63	0.63	1.43	3.50
Gradient	1.90	3.07	5.19	13.33
IPW	2.19	2.99	4.25	11.34
<i>Model 4</i>				
Naive weight	0.99	1.00	1.40	3.05
Gradient	2.87	3.72	4.47	8.57
IPW	2.78	3.36	3.18	6.98

Notes: The table presents the rejection probabilities of the uniform confidence bands for the hypothesis specified in (6.3). The rows ‘Naive weight’, ‘Gradient’ and ‘IPW’ correspond respectively to the results of the naive weighted bootstrap, the gradient bootstrap, and the weighted bootstrap of the IPW QTE estimator.

7 Empirical Application

Questions surrounding the effectiveness of microfinance as a development tool has sparked a great deal of interest from both policymakers and economists. To answer such questions a growing number of studies have implemented randomized experiments in different settings (see Banerjee, Karlan, and Zinman, 2015, and the references therein). In particular, Banerjee et al. (2015) adopted MPD in their randomization. In this section, we apply the bootstrap methods of inference developed in this paper to their data to examine both the ATEs and

QTEs on the take-up rates of microcredit to assess the effectiveness of microfinance.⁶

The sample consists of 104 areas in the city of Hyderabad in India. Based on average per capita consumption and per-household outstanding debt, the areas were grouped into pairs of similar neighborhoods. This segmentation gives 52 pairs in the sample; one area in each pair was randomly assigned to the treatment group and the other to the control group. In the treatment areas, a group-lending microcredit program was implemented. Banerjee et al. (2015) then examined the impacts of expanding access to microfinance on various outcome variables at two endlines.

Table 4: Summary Statistics

	Total	Treatment group	Control group
<i>Loan take-up rate</i>			
Spandana	0.128(0.140)	0.193(0.131)	0.062(0.117)
Any MFI	0.224(0.152)	0.265(0.151)	0.182(0.143)
<i>Matching variable</i>			
Consumption	1026.4(184.4)	1047.8(195.7)	1005.0(171.5)
Debt	36184.7(36036.5)	32694.1(17755.5)	39675.3(47776.8)
Observations	104	52	52

Notes: Unit of observation: area. The table presents the means and standard deviations (in parentheses) of two outcome variables: the take-up rate of loans from Spandana and the take-up rate of loans from any MFI, and two pair-matching variables: average per capita consumption and per-household debt.

Here we focus on the impacts of microfinance on two area-level outcome variables at the first endline. One is the area’s take-up rate of loans from Spandana, a microfinance organization that implemented the group-lending microcredit program. The other is the area’s take-up rate of loans from any microfinance institutions (MFIs). Table 4 gives descriptive statistics (means and standard deviations) for these two outcome variables as well as the matching variables used by Banerjee et al. (2015) to form the pairs in their experiments.

⁶The public-use data provided by the authors does not contain information on pair assignment. We thank Esther Duflo and Cynthia Kinnan for providing us with this information.

Table 5: ATEs of Microfinance on Take-up Rates of Microcredit

	Naive	Adj	IPW
Spandana	0.131(0.024)	0.131(0.022)	0.131(0.022)
Any MFI	0.083(0.029)	0.083(0.024)	0.083(0.027)

Notes: The table presents the ATE estimates of the effect of microfinance on the take-up rates of microcredit. Standard errors are in parentheses. The columns “Naive” and “Adj” correspond to the two-sample t -test and the adjusted t -test in Bai et al. (2019), respectively. The column “IPW” corresponds to the t -test using the standard errors estimated by the weighted bootstrap of the IPW ATE estimator.

Table 5 reports the results on the ATE estimates of the effect of microfinance on the take-up rates of microcredit with the standard errors (in parentheses) calculated by three methods. Specifically, the columns ‘Naive’ and ‘Adj’ correspond to the two-sample t -test and the adjusted t -test in Bai et al. (2019), respectively; the column ‘IPW’ corresponds to the t -test using standard errors estimated by the weighted bootstrap of the IPW ATE estimator.⁷ The results lead to the following observations. First, consistent with the findings in Banerjee et al. (2015), the ATE estimates show that expanding access to microfinance has highly significant average effects on the take-up rates of microcredit from both Spandana and any MFIs. Second, the standard errors in the adjusted t -test are lower than those in the naive t -test. This result is consistent with the finding in Bai et al. (2019). More importantly, the standard errors estimated by the IPW weighted bootstrap are also lower than those in the naive t -test and similar to those for the adjusted t -test. For example, in the test of the ATE on the take-up rate of microcredit from Spandana, the IPW method reduces the standard error by 8% compared with the naive one. The magnitude of the reduction is the same as that in the adjusted t -test. These results corroborate our earlier finding that the IPW method is an alternative to the approach adopted in Bai et al. (2019), especially when the information on pair identities is unavailable.

⁷Throughout this section, to nonparametrically estimate the propensity score in the IPW weighted bootstrap, we first standardize the data to have mean zero and variance one and then fit the standardized data via the sieve estimation based on the B-splines with the same basis as used in Section 6.

Table 6: QTEs of Microfinance on Take-up Rates of Microcredit

	Naive weight	Gradient	IPW
<i>Panel A. Spandana</i>			
25%	0.082(0.021)	0.082(0.026)	0.082(0.020)
50%	0.182(0.024)	0.182(0.021)	0.182(0.023)
75%	0.229(0.047)	0.229(0.046)	0.229(0.047)
<i>Panel B. Any MFI</i>			
25%	0.056(0.045)	0.056(0.043)	0.056(0.042)
50%	0.082(0.040)	0.082(0.034)	0.082(0.040)
75%	0.141(0.054)	0.141(0.054)	0.141(0.049)

Notes: The table presents the QTE estimates of the effect of microfinance on the take-up rates of microcredit at quantiles 25%, 50%, and 75%. Standard errors are in parentheses. The columns “Naive weight,” “Gradient,” and “IPW” correspond to the results of the naive weighted bootstrap, the gradient bootstrap, and the weighted bootstrap of the IPW QTE estimator, respectively.

Next, we estimate the QTEs of microfinance on the take-up rates of microcredit and estimate their standard errors by the three methods discussed in Section 4. Table 6 presents the results on the QTE estimates at quantile indexes 0.25, 0.5, and 0.75 with the standard errors (in parentheses) estimated by three different methods. Specifically, the columns ‘Naive weight’, ‘Gradient’, and ‘IPW’ correspond to the results of the naive weighted bootstrap, the gradient bootstrap,⁸ and the weighted bootstrap of the IPW QTE estimator, respectively. These results lead to the following two observations.

First, consistent with the theory in Section 4, the standard errors estimated by the gradient bootstrap or the IPW weighted bootstrap are mostly lower than those estimated by the naive weighted bootstrap. For example in Panel A, at the median, compared with the naive weighted bootstrap, the gradient bootstrap reduces the standard errors by 12.5% and the IPW weighted bootstrap reduces the standard errors by over 4%. In Panel B, all the standard errors computed using methods Gradient and IPW are smaller than those

⁸Using the original pair identities and matching variables in Banerjee et al. (2015), we can re-order the pairs according to the procedure described in Section 5.1. We follow Banerjee et al. (2015) in using Euclidean distance to measure the distance between the covariates in distinctive pairs.

computed using the naive method.

Second, there seem to be considerable heterogeneity in the effects of microfinance. Specifically, the treatment effects of microfinance on the take-up rates of microcredit increase as the quantile indexes increase and the increases are economically substantial. For example, in Panel A, the treatment effect increases by about 122% from the 25th percentile to the median and by about 26% from the median to the 75th percentile. In Panel B, the treatment effect at the 25th percentile is positive but not statistically significantly different from zero. The treatment effect increases by over 46% from the 25th percentile to the median and by about 72% from the median to the 75th percentile. These findings may imply that expanding access to microfinance has small, if not negligible, effects on the take-up rates of microcredit for areas in the lower tail of the distribution but that these effects become stronger for upper-ranked areas, thereby exhibiting the so-called Matthew effect.

The second observation in Table 6 indicates that the heterogeneous effects of microfinance on the take-up rates of microcredit are economically substantial. Are they statistically significant too? In Table 7, we provide statistical tests for the heterogeneity of the QTEs. Specifically, we test the null hypotheses that $q(0.50) - q(0.25) = 0$ and $q(0.75) - q(0.50) = 0$. We find that only the difference between the 25th and median QTEs in Panel A is statistically significant. This finding implies that the statistical evidence of heterogeneous treatment effects of microfinance is strong only for the areas in the lower tail of the distribution and when the loans are from Spandana.

8 Conclusion

This paper has studied estimation and inference of QTEs under MPDs and developed new bootstrap methods to improve statistical performance. Derivation of the limit distribution of QTE estimators under MPDs reveals that analytic methods of inference based on asymptotic theory requires estimation of two infinite-dimensional nuisance parameters for every quantile

Table 7: Tests for the Difference between Two QTEs of Micofinance

	Naive weight	Gradient	IPW
<i>Panel A. Spandana</i>			
$q(0.50) - q(0.25)$	0.099(0.023)	0.099(0.024)	0.099(0.022)
$q(0.75) - q(0.50)$	0.047(0.046)	0.047(0.046)	0.047(0.045)
<i>Panel B. Any MFI</i>			
$q(0.50) - q(0.25)$	0.026(0.043)	0.026(0.044)	0.026(0.044)
$q(0.75) - q(0.50)$	0.059(0.049)	0.059(0.046)	0.059(0.046)

Notes: The table presents tests for the difference between two QTEs of microfinance on the take-up rates of microcredit. Standard errors are in parentheses. The columns ‘Naive weight’, ‘Gradient’, and ‘IPW’ correspond to the results of the naive weighted bootstrap, the gradient bootstrap, and the weighted bootstrap of the IPW QTE estimator, respectively.

index of interest. A further limitation is that the naive weighted bootstrap fails to approximate the limit distribution of the QTE estimator as it does not preserve the dependence structure in the original sample. Instead, we propose a gradient bootstrap approach that can consistently approximate the limit distribution of the original estimator and is free of tuning parameters. Implementation of the gradient bootstrap requires knowledge of pair identities. So when such information is unavailable we propose a weighted bootstrap procedure based on the IPW estimator of the QTE and show that it can consistently approximate the limit distribution of the original QTE estimator. Simulations provide finite-sample evidence of these procedures that support the asymptotic findings. In our empirical application of these bootstrap methods to the real dataset in Banerjee et al. (2015) we find considerable evidence of heterogeneity in the effects of microfinance on the take-up rates of microcredit. In both the simulations and the empirical application, the two recommended bootstrap methods of inference perform well in the sense that they usually provide smaller standard errors and greater inferential accuracy than those obtained by naive bootstrap methods.

References

- Abadie, A., M. M. Chingos, and M. R. West (2018). Endogenous stratification in randomized experiments. *The Review of Economics and Statistics* 100(4), 567–580.
- Angrist, J. D. and V. Lavy (2009). The effects of high stakes high school achievement awards: Evidence from a randomized trial. *The American Economic Review* 99(4), 1384–1414.
- Athey, S. and G. Imbens (2017). Chapter 3 - the econometrics of randomized experimentsa. In A. V. Banerjee and E. Duflo (Eds.), *Handbook of Field Experiments*, Volume 1 of *Handbook of Economic Field Experiments*, pp. 73 – 140. North-Holland.
- Bai, Y. (2019). Optimality of matched-pair designs in randomized controlled trials. *Available at SSRN 3483834*.
- Bai, Y., A. Shaikh, and J. P. Romano (2019). Inference in experiments with matched pairs. *University of Chicago, Becker Friedman Institute for Economics Working Paper* (2019-63).
- Banerjee, A., E. Duflo, R. Glennerster, and C. Kinnan (2015). The miracle of microfinance? evidence from a randomized evaluation. *American Economic Journal: Applied Economics* 7(1), 22–53.
- Banerjee, A., D. Karlan, and J. Zinman (2015). Six randomized evaluations of microcredit: Introduction and further steps. *American Economic Journal: Applied Economics* 7(1), 1–21.
- Beuermann, D. W., J. Cristia, S. Cueto, O. Malamud, and Y. Cruzaguayo (2015). One laptop per child at home: Short-term impacts from a randomized experiment in peru. *American Economic Journal: Applied Economics* 7(2), 53–80.
- Bold, T., M. S. Kimenyi, G. Mwabu, A. Nganga, and J. Sandefur (2018). Experimental evidence on scaling up education reforms in kenya. *Journal of Public Economics* 168(12), 1–20.

- Bruhn, M. and D. McKenzie (2009). In pursuit of balance: Randomization in practice in development field experiments. *American Economic Journal: Applied Economics* 1(4), 200–232.
- Bugni, F. A., I. A. Canay, and A. M. Shaikh (2018). Inference under covariate-adaptive randomization. *Journal of the American Statistical Association* 113(524), 1741–1768.
- Bugni, F. A., I. A. Canay, and A. M. Shaikh (2019). Inference under covariate-adaptive randomization with multiple treatments. *Quantitative Economics* 10(4), 1747–1785.
- Chen, X. (2007). Large sample sieve estimation of semi-nonparametric models. *Handbook of econometrics* 6, 5549–5632.
- Crepon, B., F. Devoto, E. Duflo, and W. Pariente (2015). Estimating the impact of microcredit on those who take it up: Evidence from a randomized experiment in morocco. *American Economic Journal: Applied Economics* 7(1), 123–150.
- Donald, S. G. and Y.-C. Hsu (2014). Estimation and inference for distribution functions and quantile functions in treatment effect models. *Journal of Econometrics* 178, 383–397.
- Firpo, S. (2007). Efficient semiparametric estimation of quantile treatment effects. *Econometrica* 75(1), 259–276.
- Fryer, R. G. (2017). Management and student achievement: Evidence from a randomized field experiment. *National Bureau of Economic Research*.
- Fryer, R. G. (2018). The 'pupil' factory: Specialization and the production of human capital in schools. *The American Economic Review* 108(3), 616–656.
- Fryer, R. G., T. Devi, and R. Holden (2017). Vertical versus horizontal incentives in education: Evidence from randomized trials. *National Bureau of Economic Research*.

- Glewwe, P., A. F. Park, and M. Zhao (2016). A better vision for development: Eyeglasses and academic performance in rural primary schools in china. *Journal of Development Economics* 122(9), 170–182.
- Groh, M. and D. J. McKenzie (2016). Macroinsurance for microenterprises: A randomized experiment in post-revolution egypt. *Journal of Development Economics* 118(1), 1–38.
- Hagemann, A. (2017). Cluster-robust bootstrap inference in quantile regression models. *Journal of the American Statistical Association* 112(517), 446–456.
- Hahn, J. (1998). On the role of the propensity score in efficient semiparametric estimation of average treatment effects. *Econometrica* 66(2), 315–331.
- Hahn, J., K. Hirano, and D. Karlan (2011). Adaptive experimental design using the propensity score. *Journal of Business & Economic Statistics* 29(1), 96–108.
- Hirano, K., G. W. Imbens, and G. Ridder (2003). Efficient estimation of average treatment effects using the estimated propensity score. *Econometrica* 71(4), 1161–1189.
- Tabord-Meehan, M. (2018). Stratification trees for adaptive randomization in randomized controlled trials. *arXiv preprint arXiv:1806.05127*.
- van der Vaart, A. and J. A. Wellner (1996). *Weak Convergence and Empirical Processes*. Springer, New York.
- Zhang, Y. and X. Zheng (2020). Quantile treatment effects and bootstrap inference under covariate-adaptive randomization. *Quantitative Economics, forthcoming* 11(3), 957–982.