

# Optimality of Matched-Pair Designs in Randomized Controlled Trials\*

Yuehao Bai

Department of Economics

University of Michigan

[yuehaob@umich.edu](mailto:yuehaob@umich.edu)

June 15, 2020

## Abstract

This paper studies the optimality of matched-pair designs in randomized controlled trials (RCTs). Matched-pair designs are examples of stratified randomization, in which the researcher partitions a set of units into strata based on their observed covariates and assign a fraction of units in each stratum to treatment. A matched-pair design is such a procedure with two units per stratum. Despite the prevalence of stratified randomization in RCTs, implementations differ vastly. We provide an econometric framework in which, among all stratified randomization procedures, the optimal one in terms of the mean-squared error of the difference-in-means estimator is a matched-pair design that orders units according to a scalar function of their covariates and matches adjacent units. Our framework captures a leading motivation for stratifying in the sense that it shows that the proposed matched-pair design additionally minimizes the magnitude of the ex-post bias, i.e., the bias of the estimator conditional on realized treatment status. We then consider empirical counterparts to the optimal stratification using data from pilot experiments and provide two different procedures depending on whether the sample size of the pilot is large or small. For each procedure, we develop methods for testing the null hypothesis that the average treatment effect equals a prespecified value. Each test we provide is asymptotically exact in the sense that the limiting rejection probability under the null equals the nominal level. We run an experiment on the Amazon Mechanical Turk using one of the proposed procedures, replicating one of the treatment arms in [DellaVigna and Pope \(2018\)](#), and find the standard error decreases by 29%, so that only half of the sample size is required to attain the same standard error.

**KEYWORDS:** Matched-pair design, stratified randomization, randomized controlled trial, ex-post bias, treatment effect, stratification, pilot experiment, matched pairs

**JEL CLASSIFICATION CODES:** C12, C13, C14, C90

---

\*I am deeply grateful for the encouragement and guidance from my advisors Azeem Shaikh, Stephane Bonhomme, Alex Torgovitsky, and Leonardo Bursztyn. I thank Marinho Bertanha, Wooyong Lee, Joshua Shea, and Max Tabord-Meehan for extensive feedback on earlier drafts of the paper. I would also like to thank seminar participants at many institutions for helpful comments on the paper. I gratefully acknowledge the financial support from the William Rainey Harper/Provost Dissertation Year Fellowship.

# 1 Introduction

This paper studies the optimality of matched-pair designs in randomized controlled trials (RCTs). Matched-pair designs are examples of stratified randomization, in which the researcher partitions a set of units into strata based on their observed covariates and assigns a fraction of units in each stratum to treatment. A matched-pair design is a stratified randomization procedure with two units in each stratum. Stratified randomization is prevalent in economics and more broadly the sciences. A simple search with the keyword “stratified” in the AEA RCT Registry reveals about 500 RCTs. The procedures in these papers, however, differ vastly in terms of variables being stratified on, how strata are formed, and numbers of strata. Among these procedures, matched-pair designs have recently gained popularity. 56% of researchers interviewed in [Bruhn and McKenzie \(2009\)](#) have used matched-pair designs at some point in their research. Moreover, more than 40 ongoing experiments in the AEA RCT Registry use matched-pair designs. See [Section 1.1](#) for a list of papers. Despite the popularity of matched-pair designs, there is little theory justifying their use in RCTs. We provide an econometric framework in which a certain form of matched-pair design emerges as optimal among all stratified randomization procedures. As will be explained below, an attractive feature of our framework is that it captures a leading motivation for stratifying in the sense that it shows that the proposed matched-pair design minimizes the second moment of the ex-post bias, i.e., the bias of the estimator conditional on realized treatment status. We then provide empirical counterparts to the optimal procedure and illustrate one of the proposed procedures by conducting an actual experiment on the Amazon Mechanical Turk (MTurk). In particular, we replicate one of the treatment arms from the experiment in [DellaVigna and Pope \(2018\)](#) and show that the standard error decreases by 29% compared to original results, which means that only half of the sample size is required to attain the same level of precision as in the original paper.

We begin by studying settings where treated fractions are identical across strata. In such settings, it is natural to estimate the average treatment effect (ATE) by the difference in means of the treated and control groups. The properties of the difference-in-means estimator, however, vary substantially with stratifications. In the main text, we further restrict treated fractions to be  $\frac{1}{2}$  within each stratum, but in the appendix, we provide extensions to settings where treated fractions are identical across strata but not equal to  $\frac{1}{2}$  and where they are in addition allowed to vary across a fixed number of subpopulations. Our first result shows the mean-squared error (MSE) of the difference-in-means estimator conditional on the covariates is remarkably minimized by a matched-pair design, where units are ordered by their values of a scalar function of the covariates and paired adjacently. The scalar function is defined by the sum of the expectations of potential outcomes if treated and not treated conditional on the covariates. To the best of our knowledge, our result is the first to characterize the optimal one among all stratified randomization procedures, and additionally, it holds under almost no assumption on the distributions of potential outcomes. In a closely related paper, [Barrios \(2013\)](#) considers minimizing the variance of the difference-in-means estimator. Despite having “optimal stratification” in the title of his paper, he only shows that a certain matched-pair design is optimal among all matched-pair designs, instead of all stratified randomization procedures. Although intuitively attractive, it is not always without loss of generality to restrict attention to matched-pair designs in the first place. [Example E.7](#) shows that under a minimax criterion the optimal stratification might not be a matched-pair design. We show, however, that we could restrict attention to matched-pair designs if the criterion is MSE. In fact, we show that the optimality of

matched-pair designs holds under any expected utility criterion, and even any criterion convex in the distribution of treatment status. See Remark 3.4 for more details. Moreover, Barrios (2013) assumes a homogeneous treatment effect and uses only information about untreated potential outcomes in his analysis, while our optimality result instead holds under heterogeneous treatment effects. Finally, as explained below, we provide novel results relating the MSE to the ex-post bias, as well as novel results on the large sample properties of empirical counterparts to the optimal procedure as well as formal results on inference.

We then study the properties of empirical counterparts to this optimal stratification, in which we replace the unknown scalar function with estimates based on pilot data. Pilot experiments are frequently available in practice. Around 350 out of 3000 experiments in the AEA RCT Registry have pilot experiments. For more examples, see Karlan and Zinman (2008), Karlan and Appel (2016), Karlan and Wood (2017), DellaVigna and Pope (2018), and papers cited in Section 1.1. We first consider a plug-in procedure that estimates the scalar function using data from a pilot experiment and matches the units in the main experiment into pairs based on their values of the estimated function. Under a weak consistency requirement on the plug-in estimator, or more precisely, that it is  $L^2$ -consistent for the scalar function, we show that as the sample sizes of both the pilot and the main experiments increase, the limiting variance of a suitable normalization of the difference-in-means estimator under the plug-in procedure is the same as that under the infeasible optimal procedure. Equivalently, under such a normalization, the limiting MSE of the estimator is the same as that under the optimal stratification. The consistency requirement is satisfied by a large class of nonparametric estimation methods including machine learning methods in high-dimensional settings, i.e., when the dimension of covariates is large. In this sense, when the sample size of the pilot is large, the plug-in procedure is optimal. Of course, this property no longer holds when the sample size of the pilot is small. But even then, researchers may well be content with the plug-in procedure because it results in smaller limiting variance of the difference-in-means estimator than many alternatives. That said, we may be concerned that a poor estimate of the scalar function leads to a matched-pair design under which the MSE of the estimator is large. Therefore, we additionally consider a penalized procedure under which, according to simulation studies with small pilots, the MSE of the estimator is often smaller than those under plug-in and other commonly-used procedures. The procedure is named so because it can be viewed as penalizing the plug-in procedure by the standard error of the plug-in estimate. Another attractive feature of the penalized procedure is that it is optimal in integrated risk in a Bayesian framework with Gaussian priors and linear conditional expectations of potential outcomes.

For each procedure, we develop methods for testing the null hypothesis that the ATE equals a prespecified value. Inference for matched-pair designs is challenging because of the difficulty of consistently estimating the limiting variance of the ATE. Indeed, this is the main reason why Athey and Imbens (2017) suggest not to use matched-pair designs. We get around this problem by a novel standard error adjustment and Lipschitz conditions that guarantee the smoothness of conditional expectations of potential outcomes given the covariates. This condition, together with the observation that paired observations become close in terms of the pairing covariate in the limit, enables us to estimate the limiting variance consistently. Therefore, each test we provide is asymptotically exact in the sense that the limiting rejection probability under the null equals the nominal level. Our results extend those in Bai et al. (2019) to settings where units are matched according to (random) functions of their covariates instead of the covariates themselves. A special feature of inference under the

plug-in procedure is that the same test is valid regardless of the sample size of the pilot. Inference methods under both the plug-in and the penalized procedures are computationally easy.

Our results on optimal stratification formalizes the motivation for using stratified randomization by showing that minimizing the conditional (on covariates) MSE is equivalent to minimizing the conditional second moment of the ex-post bias, i.e., the bias of the estimator conditional on both the covariates and realized treatment status. Furthermore, the two problems are both equivalent to minimizing the conditional variance of the ex-post bias. To illustrate the intuition behind this minimization problem, it is instructive to consider the special case where there is a single binary covariate. Consider an RCT with 100 units, composed of 50 women and 50 men. The intuitive motivation for stratifying by gender is as follows: if all the units are in one stratum, then it could happen that 40 women are treated while only 10 men are so, so that a large part of the difference in treated and control units could be from the difference in gender instead of the treatment itself; on the other hand, if we stratify by gender, then we always end up treating 25 women and 25 men. The intuitive motivation is formalized by the comparison of the ex-post bias. Since the ex-post bias only depends on how many men and women treated instead of their identities, it varies across realized treatment status if all the units are in one stratum, but is identical if we stratify by gender. As a result, the conditional variance of the ex-post bias is positive if all the units are in one stratum but zero if we stratify by gender. When there are more covariates or when some of them are continuous, it is hard to see only by inspection which stratification minimizes the second moment or the variance of the ex-post bias, but the solution is given by the optimal stratification. Our results could also be viewed as formalizing the discussion about which covariates should be stratified on, e.g., the recommendation in [Bruhn and McKenzie \(2009\)](#) and [Glennerster and Takavarasha \(2013\)](#) for using covariates most correlated with the outcome.

One might be tempted to think that stratification only matters in finite sample, or that the limiting distribution of the difference-in-means estimator should be the same no matter how units are stratified. We show that this is not the case. Specifically, we show that under any stratification with a fixed number of strata, the limiting variance of the estimator is weakly greater than and typically strictly greater than that under the optimal stratification. See Remark 5.8 below for more details. In addition, different matched-pair designs lead to different limiting variances of the estimator. Although [Bai et al. \(2019\)](#) show the same limiting variance when units are matched not according to our optimal stratification but to minimize the sum of Euclidean distances of the covariates themselves, their result holds with a fixed number of covariates, while both our optimality result and our plug-in procedure allows the number of covariates to grow with the sample size. See Remark 5.1 below for more details. Sometimes researchers do not stratify but run covariate adjustments afterwards, i.e., regress the outcome on the treatment and some observed covariates. Although the most flexible covariate adjustment could lead to the same limiting variance, they usually require higher order smoothness conditions and a sufficiently fast convergence rate of the nonparametric adjustment component, while we only require the  $L^2$ -consistency of the estimator. See Remark 5.2 below and [Rothe \(2020\)](#) for more details.

While pilot experiments are common in RCTs, there are scenarios in which they are either not available or are performed on a different population from units in the main experiment. For those scenarios, we study a minimax problem that does not rely on pilot data, where we assume the data generating process is chosen by nature adversarially among a large class of distributions that could be characterized by bounded polyhedrons.

In particular, we minimize the variance of the ex-post bias of the difference-in-means estimator conditional on the covariates under the worst possible distribution in this class by choosing across matched-pair designs. The framework accommodates many common shape restrictions on the conditional expectations of potential outcomes given the covariates, including Lipschitz continuity, monotonicity, and convexity. We then rewrite the minimax problem into a mixed integer linear program (MILP) which is computationally easy. Simulation evidence further suggests although the minimax matched-pair design is in general not minimax-optimal among all stratifications except when there is a single covariate, it is often close to being so.

The remainder of the paper is organized as follows. In Section 2, we introduce the setup and notation. We study the optimal stratification in Section 3. In Section 4, we consider empirical counterparts to the optimal stratification, using data from pilot experiments. We consider the plug-in procedure with large pilots and the penalized procedure with small pilots. Section 5 includes asymptotic results and methods for inference for ATE. In Section 6, we illustrate the properties of different procedures in a small simulation study. Section 7 discusses results from the MTurk experiment using the penalized procedure. The experiment shows a 29% reduction in standard error compared to results in the original paper, which means that we need only half of the sample size to attain the same standard error. Section 8 briefly discusses the minimax procedure, the details of which are included in Appendix E. We conclude with recommendations for empirical practice in Section 9.

## 1.1 Related literature

This paper is most closely related to [Barrios \(2013\)](#) and [Tabord-Meehan \(2020\)](#). [Barrios \(2013\)](#) considers minimizing the variance of the difference-in-means estimator but assumes a homogeneous treatment effect and uses only information about untreated potential outcomes in his analysis. Despite having “optimal stratification” in the title, his paper only shows that a certain matched-pair design is optimal among all matched-pair designs, instead of all stratifications. We instead show that a certain matched-pair design is optimal among all stratifications, without assuming a homogeneous treatment effect. Moreover, we provide novel results relating the MSE to the ex-post bias. We also provide formal results on the large sample properties of empirical counterparts to the optimal procedure as well as formal results on inference. [Tabord-Meehan \(2020\)](#) considers optimality within a specific class of stratifications, which is a certain class of stratification trees. Since the number of strata is fixed in his asymptotic framework, his paper precludes matched-pair designs. We instead provide analytical characterization of the optimal one among the set of all stratifications. [Remark 5.9](#) elaborates the details of the comparison between the two papers, and in particular, notes that it is straightforward to combine the procedures in both papers. Under the combined procedure, the asymptotic variance of the fully saturated estimator is no greater than and typically strictly smaller than that when using the procedure in [Tabord-Meehan \(2020\)](#) alone.

Recent examples of stratified randomization in development economics include [Aker et al. \(2012, page 97\)](#), [Alatas et al. \(2012, page 1211\)](#), [Ashraf et al. \(2010, page 2393\)](#), [Dupas and Robinson \(2013, page 168\)](#), [Callen et al. \(2014, page 133\)](#), [Banerjee et al. \(2015, page 31\)](#), [Duflo et al. \(2015, page 96\)](#), [Duflo et al. \(2015, footnote 6\)](#), [Chong et al. \(2016, page 228\)](#), [Berry et al. \(2018, page 75\)](#), [Bursztyn et al. \(2018, page 1570\)](#), [Callen et al. \(2018, page 10\)](#), [Dupas et al. \(2018, page 264\)](#), [Bursztyn et al. \(2019, footnote 15\)](#), [Casaburi](#)

and Macchiavello (2019, page 548), Chen and Yang (2019, page 2308), Dizon-Ross (2019, page 2738), Khan et al. (2019, page 254), and Muralidharan et al. (2019, page 1434). See Bruhn and McKenzie (2009) for more examples in economics and Rosenberger and Lachin (2015) and Lin et al. (2015) for examples in clinical trials. For examples of matched-pair designs, see Riach and Rich (2002), Ashraf et al. (2006), Panagopoulos and Green (2008), Angrist and Lavy (2009), Imai et al. (2009), Sondheimer and Green (2010), List and Rasul (2011), White (2013), Bhargava and Manoli (2015), Banerjee et al. (2015), Crépon et al. (2015), Bruhn et al. (2016), Glewwe et al. (2016), Groh and McKenzie (2016), Bertrand and Duflo (2017), Fryer (2017), Fryer et al. (2017), Heard et al. (2017), Fryer (2018), Bai et al. (2019), and the references therein. See Appendix F for a list of ongoing experiments using matched-pair designs in the AEA RCT Registry. Matched-pair designs are also implemented in leading experimental design packages, including `sampsi_mcc` in Stata. Imbens (2011) and Athey and Imbens (2017) discuss the benefits of stratified randomization in a finite sample framework and a simple example with one binary covariate. These two papers, together with Chapter 10 in Imbens and Rubin (2015), recognize the merit of matched-pair designs in terms of estimation but suggest they come with the cost that the asymptotic variance of the estimator is hard to estimate. Our inference procedure solves this problem and therefore eliminates this cost. Besides Bai et al. (2019), inference under matched-pair designs has also been studied in Abadie and Imbens (2008), who consider another adjustment of standard error, in Fogarty (2018a) and Fogarty (2018b), who provides conservative estimators for the asymptotic variance, and de Chaisemartin and Ramirez-Cuellar (2019), under a sampling scheme different from that in Bai et al. (2019) and a cluster setting.

For general references on RCTs, see Duflo et al. (2007), Bruhn and McKenzie (2009), Glennerster and Takavarasha (2013), Rosenberger and Lachin (2015), Peters et al. (2016), and the Handbook of Field Experiments, Duflo and Banerjee (2017). For earlier work on the optimal design of experiments under parametric models with block structures, see Cox and Reid (2000), Bailey (2004), and Pukelsheim (2006). A series of papers also examine optimal design in RCTs. Hahn et al. (2011) assume independent random sampling across units, whereas stratified randomization induces dependence within each stratum. Chambaz et al. (2015) adaptively assign treatment status for each new observation based on those of the previous units. Kallus (2018) studies optimal treatment assignment from a minimax perspective and optimizes over treatment assignments rather than stratifications. Freedman (2008) and Lin (2013) compare regression-adjusted estimators and the difference-in-means estimator, assuming all the units are in one stratum. Re-randomization, another commonly-used method to balance covariates, is studied in parametric models in Morgan et al. (2012), Morgan and Rubin (2015), Li et al. (2018), Schultzberg and Johansson (2019), and Johansson et al. (2019). Kasy (2016) considers a Bayesian problem in a parametric model, where both the prior and the distributions of potential outcomes are Gaussian with known parameters, and concludes that researchers should never randomize. On the contrary, Wu (1981), Li (1983), and Hooper (1989), and Bai (2019) show the optimality of certain randomization schemes in minimax frameworks. Carneiro et al. (2019) examine the trade-off between collecting more units and more covariates for each unit when designing an RCT under fixed budget. A growing literature, including Manski (2004), Kitagawa and Tetenov (2018), and Mbakop and Tabord-Meehan (2018), considers empirical welfare maximization by assigning treatment status. Banerjee et al. (2019) study optimal experiments under a combination of Bayesian and minimax criteria in terms of welfare.

## 2 Setup and notation

Let  $Y_i$  denote the observed outcome of interest for the  $i$ th unit,  $D_i$  denote the treatment status for the  $i$ th unit and  $X_i = (X_{i,1}, \dots, X_{i,p})' \in \mathbf{R}^p$  denote the observed, baseline covariates for the  $i$ th unit. Further denote by  $Y_i(1)$  the potential outcome of the  $i$ th unit if treated and by  $Y_i(0)$  if not treated. As usual, the observed outcome is related to the potential outcomes and treatment status by the relationship

$$Y_i = Y_i(1)D_i + Y_i(0)(1 - D_i).$$

In addition, we define  $W_i = (Y_i, X_i', D_i)'$ . For ease of exposition, we assume the sample size is even and denote it by  $2n$ . We assume that  $((Y_i(1), Y_i(0), X_i) : 1 \leq i \leq 2n)$  is an i.i.d. sequence of random vectors with distribution  $Q$ . For any random vector indexed by  $i$ ,  $A_i$ , define  $A^{(n)} = (A_1, \dots, A_{2n})'$ . Our parameter of interest is the average treatment effect (ATE) under  $Q$ :

$$\theta(Q) = E_Q[Y_i(1) - Y_i(0)]. \quad (1)$$

For ease of exposition, we will at times suppress the dependence of various quantities on  $Q$ , e.g., use  $\theta$  to refer to  $\theta(Q)$ . In stratified randomization, the first step is to partition the set of units into strata. Formally, we define a stratification  $\lambda = \{\lambda_s : 1 \leq s \leq S\}$  as a partition of  $\{1, \dots, 2n\}$ , i.e.,

- (a)  $\lambda_s \cap \lambda_{s'} = \emptyset$  for all  $s$  and  $s'$  such that  $1 \leq s \neq s' \leq S$ .
- (b)  $\bigcup_{1 \leq s \leq S} \lambda_s = \{1, \dots, 2n\}$ .

Let  $\Lambda_n$  denote the set of all stratifications of  $2n$  units. Many results in the paper will feature matched-pair designs. Recall that a permutation of  $\{1, \dots, 2n\}$  is a function that maps  $\{1, \dots, 2n\}$  onto itself. Let  $\Pi_n$  denote the group of all permutations of  $\{1, \dots, 2n\}$ . A matched-pair design is a stratification with

$$\lambda = \{\{\pi(2s-1), \pi(2s)\} : 1 \leq s \leq n\},$$

where  $\pi \in \Pi_n$ . Further define  $\Lambda_n^{\text{pair}} \subseteq \Lambda_n$  as the set of all matched-pair designs for  $2n$  units.

Define  $n_s = |\lambda_s|$  and  $\tau_s$  as the treated fraction in stratum  $\lambda_s$ . Under stratified randomization, given  $X^{(n)}$ ,  $\lambda$ , and  $(\tau_s : 1 \leq s \leq S)$ , the treatment assignment scheme is as follows: independently for  $1 \leq s \leq S$ , uniformly at random choose  $n_s \tau_s$  units in  $\lambda_s$  and assign  $D_i = 1$  for them, and assign  $D_i = 0$  for the other units. The treatment assignment scheme implies that

$$(Y^{(n)}(0), Y^{(n)}(1)) \perp\!\!\!\perp D^{(n)} | X^{(n)}. \quad (2)$$

It also implies that  $n_s \tau_s$  is an integer for  $1 \leq s \leq S$ . Note that the distribution of  $D^{(n)}$  depends on  $\lambda$ . In the remainder of the paper, we assume the following about the treatment assignment scheme unless indicated otherwise:

**Assumption 2.1.** The treatment assignment scheme satisfies  $\tau_s \equiv \frac{1}{2}$ .

Assumption 2.1 implies that the size of each stratum has to be an even number. Most results below could be extended to settings where  $\tau_s \equiv \tau \in (0, 1)$  or where they are in addition allowed to vary across subpopulations. See Appendix B for more details.

We estimate the ATE by the difference in means between the treated and control groups. Formally, for  $d \in \{0, 1\}$ , define

$$\hat{\mu}_n(d) = \frac{\sum_{1 \leq i \leq 2n} Y_i I\{D_i = d\}}{\sum_{1 \leq i \leq 2n} I\{D_i = d\}} = \frac{1}{n} \sum_{1 \leq i \leq 2n: D_i = d} Y_i.$$

The difference-in-means estimator is defined as

$$\hat{\theta}_n = \hat{\mu}_n(1) - \hat{\mu}_n(0). \quad (3)$$

The difference-in-means estimator is widely used because it is simple and transparent. Under Assumption 2.1, it coincides with the estimator from regressing the outcome on treatment status and strata fixed effects, and the estimator from the fully saturated regression, both of which are also widely used in the analysis of RCTs. See, for example, Duflo et al. (2007), Glennerster and Takavarasha (2013), and Crépon et al. (2015).

### 3 Optimal stratification

For any stratification  $\lambda \in \Lambda_n$ , our objective function is the mean-squared error (MSE) of  $\hat{\theta}_n$  for  $\theta$  conditional on  $X^{(n)}$  under  $\lambda$ :

$$\text{MSE}(\lambda|X^{(n)}) = E_\lambda[(\hat{\theta}_n - \theta)^2|X^{(n)}]. \quad (4)$$

Here, the subscript  $\lambda$  of  $E$  indicates that the expectation depends on  $\lambda$ , since the distribution of treatment status  $D^{(n)}$  depends on  $\lambda$ . We consider minimizing the conditional MSE defined in (4) over the set of all stratifications:

$$\min_{\lambda \in \Lambda_n} \text{MSE}(\lambda|X^{(n)}). \quad (5)$$

The solution will depend on features of the distribution which are generally unknown, and we will consider empirical counterparts to the solution, in which unknown quantities are replaced by estimates using data from pilot experiments, in Section 4. By Assumption 2.1, other aspects of the stratified randomization procedure, especially the treated fractions, are fixed. Therefore, the stratification that solves (5) corresponds to an optimal stratified randomization procedure among all those satisfying Assumption 2.1.

In order to describe an important result that leads to the solution to (5), we define the ex-ante bias of  $\hat{\theta}_n$  for  $\theta$  conditional on  $X^{(n)}$  as

$$\text{Bias}_{n,\lambda}^{\text{ante}}(\hat{\theta}_n|X^{(n)}) = E_\lambda[\hat{\theta}_n|X^{(n)}] - \theta, \quad (6)$$

and the ex-post bias of  $\hat{\theta}_n$  for  $\theta$  conditional on  $X^{(n)}$  and  $D^{(n)}$  as

$$\text{Bias}_{n,\lambda}^{\text{post}}(\hat{\theta}_n|X^{(n)}, D^{(n)}) = E_\lambda[\hat{\theta}_n|X^{(n)}, D^{(n)}] - \theta. \quad (7)$$



Here, ex-ante bias refers to the bias conditional only on covariates, before treatment status is assigned; ex-post bias refers to the bias conditional on both the covariates and treatment status, i.e, after treatment status is assigned. By definition,

$$E_\lambda[\text{Bias}_{n,\lambda}^{\text{post}}(\hat{\theta}_n|X^{(n)}, D^{(n)})|X^{(n)}] = \text{Bias}_{n,\lambda}^{\text{ante}}(\hat{\theta}_n|X^{(n)}), \quad (8)$$

i.e., the expectation of the ex-post bias over the distribution of treatment status equals the ex-ante bias. Note that by (3),

$$\hat{\theta}_n = \frac{1}{n} \sum_{1 \leq i \leq 2n} (Y_i(1)D_i - Y_i(0)(1 - D_i)).$$

Under Assumption 2.1,

$$\text{Bias}_{n,\lambda}^{\text{ante}}(\hat{\theta}_n|X^{(n)}) = \frac{1}{2n} \sum_{1 \leq i \leq 2n} (E[Y_i(1)|X_i] - E[Y_i(0)|X_i]) - \theta, \quad (9)$$

so that ex-ante bias is identical across  $\lambda \in \Lambda_n$ .

To solve (5), we decompose the conditional MSE as follows. First, note that

$$\text{MSE}(\lambda|X^{(n)}) = \text{Bias}_{n,\lambda}^{\text{ante}}(\hat{\theta}_n|X^{(n)})^2 + \text{Var}_\lambda[\hat{\theta}_n|X^{(n)}]. \quad (10)$$

Here,  $\text{Var}_\lambda$  indicates that the distribution of treatment status depends on  $\lambda$ . By (9), the first term on the right-hand side is identical across all  $\lambda \in \Lambda_n$ . Hence, (5) is equivalent to minimizing the second term on the right-hand side of (10), which could be further decomposed into

$$\text{Var}_\lambda[\hat{\theta}_n|X^{(n)}] = E_\lambda[\text{Var}[\hat{\theta}_n|X^{(n)}, D^{(n)}]|X^{(n)}] + \text{Var}_\lambda[E[\hat{\theta}_n|X^{(n)}, D^{(n)}]|X^{(n)}]. \quad (11)$$

By (2), conditional on  $X^{(n)}$  and  $D^{(n)}$ ,  $(Y_i(0), Y_i(1))$ 's are independent across  $i$ , so that for any  $\lambda \in \Lambda_n$ , the first term on the right-hand side of (11) equals

$$E_\lambda \left[ \frac{1}{n^2} \sum_{1 \leq i \leq 2n} (\text{Var}[Y_i(1)|X_i]D_i + \text{Var}[Y_i(0)|X_i](1 - D_i)) \middle| X^{(n)} \right] = \frac{1}{2n^2} \sum_{1 \leq i \leq 2n} (\text{Var}[Y_i(1)|X_i] + \text{Var}[Y_i(0)|X_i]), \quad (12)$$

which is also identical across all  $\lambda \in \Lambda_n$ . Here, we use (2), the facts that  $D_i(1 - D_i) = 0$  for  $1 \leq i \leq 2n$ , and that  $E[D_i|X^{(n)}] = \frac{1}{2}$ . Hence, (5) is further equivalent to minimizing the second term on the right-hand side of (11), which equals

$$\text{Var}_\lambda[\text{Bias}_{n,\lambda}^{\text{post}}(\hat{\theta}_n|X^{(n)}, D^{(n)})|X^{(n)}]. \quad (13)$$

Furthermore, we have

$$\begin{aligned} & \text{Var}_\lambda[E[\hat{\theta}_n|X^{(n)}, D^{(n)}]|X^{(n)}] \\ &= E_\lambda[(E[\hat{\theta}_n|X^{(n)}, D^{(n)}] - E[\hat{\theta}_n|X^{(n)}])^2|X^{(n)}] \\ &= E_\lambda[(\text{Bias}_{n,\lambda}^{\text{post}}(\hat{\theta}_n|X^{(n)}, D^{(n)}) - \text{Bias}_{n,\lambda}^{\text{ante}}(\hat{\theta}_n|X^{(n)}))^2|X^{(n)}] \end{aligned}$$

$$\begin{aligned}
&= E_\lambda[\text{Bias}_{n,\lambda}^{\text{post}}(\hat{\theta}_n|X^{(n)}, D^{(n)})^2|X^{(n)}] - 2E_\lambda[\text{Bias}_{n,\lambda}^{\text{post}}(\hat{\theta}_n|X^{(n)}, D^{(n)})\text{Bias}_{n,\lambda}^{\text{ante}}(\hat{\theta}_n|X^{(n)})|X^{(n)}] + \text{Bias}_{n,\lambda}^{\text{ante}}(\hat{\theta}_n|X^{(n)})^2 \\
&= E_\lambda[\text{Bias}_{n,\lambda}^{\text{post}}(\hat{\theta}_n|X^{(n)}, D^{(n)})^2|X^{(n)}] - 2E_\lambda[\text{Bias}_{n,\lambda}^{\text{post}}(\hat{\theta}_n|X^{(n)}, D^{(n)})|X^{(n)}]\text{Bias}_{n,\lambda}^{\text{ante}}(\hat{\theta}_n|X^{(n)}) + \text{Bias}_{n,\lambda}^{\text{ante}}(\hat{\theta}_n|X^{(n)})^2 \\
&= E_\lambda[\text{Bias}_{n,\lambda}^{\text{post}}(\hat{\theta}_n|X^{(n)}, D^{(n)})^2|X^{(n)}] - 2\text{Bias}_{n,\lambda}^{\text{ante}}(\hat{\theta}_n|X^{(n)})^2 + \text{Bias}_{n,\lambda}^{\text{ante}}(\hat{\theta}_n|X^{(n)})^2 \\
&= E_\lambda[\text{Bias}_{n,\lambda}^{\text{post}}(\hat{\theta}_n|X^{(n)}, D^{(n)})^2|X^{(n)}] - \text{Bias}_{n,\lambda}^{\text{ante}}(\hat{\theta}_n|X^{(n)})^2, \tag{14}
\end{aligned}$$

where the first equality follows from definition, the second follows from (6) and (7), the third equality follows from expanding the square, the fourth equality follows since  $\text{Bias}_{n,\lambda}^{\text{ante}}(\hat{\theta}_n|X^{(n)})$  is constant conditional on  $X^{(n)}$ , and the fifth equality follows from (8). By (9),  $\text{Bias}_{n,\lambda}^{\text{ante}}(\hat{\theta}_n|X^{(n)})$  is the same across  $\lambda$ , and therefore it follows from (10)–(14) that (5) is equivalent to minimizing the first term in (14), i.e., the second moment of the ex-post bias. We summarize the results in the following lemma:

**Lemma 3.1.** *Suppose the treatment assignment scheme satisfies Assumption 2.1. Then, the set of solutions to (5) is the same as the set of solutions to*

$$\min_{\lambda \in \Lambda_n} E_\lambda[\text{Bias}_{n,\lambda}^{\text{post}}(\hat{\theta}_n|X^{(n)}, D^{(n)})^2|X^{(n)}], \tag{15}$$

and the set of solutions to

$$\min_{\lambda \in \Lambda_n} \text{Var}_\lambda[\text{Bias}_{n,\lambda}^{\text{post}}(\hat{\theta}_n|X^{(n)}, D^{(n)})|X^{(n)}]. \tag{16}$$

**Remark 3.1.** We have shown that minimizing the conditional MSE is equivalent to (15), i.e., minimizing the second moment of the ex-post bias, and (16), i.e., minimizing the variance of the ex-post bias conditional on the covariates. This equivalence holds since the mean of the ex-post bias is the ex-ante bias, which is the same across stratifications by (9). (15) is more convenient for intuition, while (16) is easier to solve. ■

The following theorem contains our main result on optimal stratification, which shows that (5) is solved by a matched-pair design, where units are ordered by their values of a scalar function of the covariates and paired adjacently. In particular, define the function

$$g(x) = E[Y_i(1) + Y_i(0)|X_i = x]. \tag{17}$$

For any measurable function  $h : \mathbf{R}^p \rightarrow \mathbf{R}$ , define  $h_i = h(X_i)$ . Let  $\pi^g \in \Pi_n$  be such that  $g_{\pi^g(1)} \leq \dots \leq g_{\pi^g(2n)}$ . Define the stratification

$$\lambda^g(X^{(n)}) = \{\{\pi^g(2s-1), \pi^g(2s)\} : 1 \leq s \leq n\}. \tag{18}$$

**Theorem 3.1.** *Suppose the treatment assignment scheme satisfies Assumption 2.1. Then,  $\lambda^g(X^{(n)})$  defined in (18) solves (5).*

**Remark 3.2.** Figure 3 illustrates the optimal stratification in (18). The outline of the proof of Theorem 3.1 is as follows. Lemma C.1 shows that each stratification is a convex combination of matched-pair designs. Therefore, one of the solutions to (5) must be a “vertex” of these convex combinations, i.e., a matched-pair design. Using the second part of Lemma 3.1, we show that the conditional MSEs of  $\hat{\theta}_n$  under matched-pair designs differ only in terms of the sum of squared distances in  $g$  within pairs. The sum is minimized by the

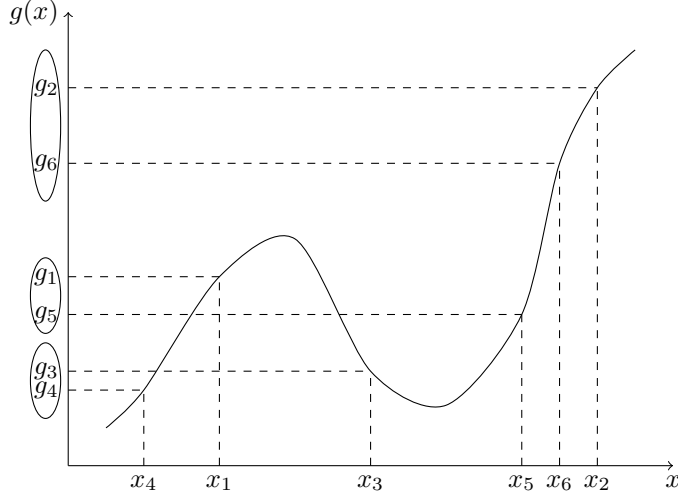


Figure 1: Illustration of the optimal stratification defined in (18). In the example,  $p = 1$ , i.e.,  $X_i$ 's are scalars. The optimal stratification is  $\{\{3, 4\}, \{1, 5\}, \{2, 6\}\}$ .

stratification defined in (18), according to a variant of the Hardy-Littlewood-Pólya rearrangement inequality for non-bipartite matching. ■

**Remark 3.3.** Note from (17) that  $g_i$  is a scalar regardless of the dimension  $p$  of  $X_i$ . Moreover, (18) depends not on the values but merely the ordering of  $g_i$ ,  $1 \leq i \leq 2n$ . For instance, if  $p = 1$  and we are certain that  $g(x)$  is monotonic in  $x$ , then it is optimal to order units by  $X_i$ ,  $1 \leq i \leq n$  and pair the units adjacently, regardless of the values of  $g_i$ ,  $1 \leq i \leq 2n$ . ■

**Remark 3.4.** Using similar arguments as those used to establish Theorem 3.1, it is possible to show that if MSE in (5) is replaced by any expected utility criterion, then one of the solutions is a matched-pair design. It is further possible to show the same conclusion holds for any criterion that is convex in the distribution of treatment status. Therefore, the optimality of matched-pair designs holds quite generally. That said, it is nontrivial to characterize the form of the optimal matched-pair design in those general settings and this is left for future work. ■

**Remark 3.5.** Theorem B.1 in the appendix examines the scenario where  $\tau_s \equiv \tau \in (0, 1)$ . Assume  $\tau = \frac{l}{k}$  where  $l, k \in \mathbb{Z}$ ,  $0 < l < k$ , and they are relatively prime, and that the sample size is  $kn$ . Define

$$g^\tau(X_i) = \frac{E[Y_i(1)|X_i]}{\tau} + \frac{E[Y_i(0)|X_i]}{1 - \tau}. \quad (19)$$

Let  $\pi^{\tau, g^\tau}$  be a permutation of  $\{1, \dots, kn\}$  such that  $g_{\pi^{\tau, g^\tau}(1)}^\tau \leq \dots \leq g_{\pi^{\tau, g^\tau}(kn)}^\tau$ . We show that (5) is solved by

$$\lambda^{\tau, g}(X^{(n)}) = \{\{\pi^{\tau, g^\tau}((s-1)k+1), \dots, \pi^{\tau, g^\tau}(sk)\} : 1 \leq s \leq n\}, \quad (20)$$

The scalar function  $g^\tau$  adjusts for treatment probabilities by inverse probability weighting. For a similar design, see Bold et al. (2018). ■

We illustrate Lemma 3.1, and in particular (15), in a small simulation study. In this example,  $2n = 100$ ;

$X_i = (X_{i,1}, X_{i,2})'$ ;  $X_{i,1}$  and  $X_{i,2}$  are both distributed as  $N(0, 1)$ , independent from each other, and i.i.d. across  $1 \leq i \leq 2n$ ; and  $E[Y_i(d)|X_i] = X_i'\beta(d)$  for  $\beta(0) = (0, 1.5)'$  and  $\beta(1) = (0.5, 2)'$ . As a result,  $\theta = 0$ . In Figure 2, we plot the densities of the distributions of  $\text{Bias}_{n,\lambda}^{\text{post}}(\hat{\theta}_n|X^{(n)}, D^{(n)})$  defined in (7) over 1000 draws of  $X^{(n)}$  and  $D^{(n)}$ , for different treatment assignment schemes:

**Oracle** stratified randomization using the infeasible optimal procedure defined by (18).

**by1** stratified randomization with two strata separated by the sample median of  $X_{i,1}$ .

**by2** stratified randomization with two strata separated by the sample median of  $X_{i,2}$ .

**SRS** Simple Random Sampling, i.e.,  $(D_i, 1 \leq i \leq 2n)$  are i.i.d. Bernoulli( $\frac{1}{2}$ ).

Note that the distribution of  $\text{Bias}_{n,\lambda}^{\text{post}}(\hat{\theta}_n|X^{(n)}, D^{(n)})$  under **Oracle** is much more concentrated than those under other treatment assignment schemes.

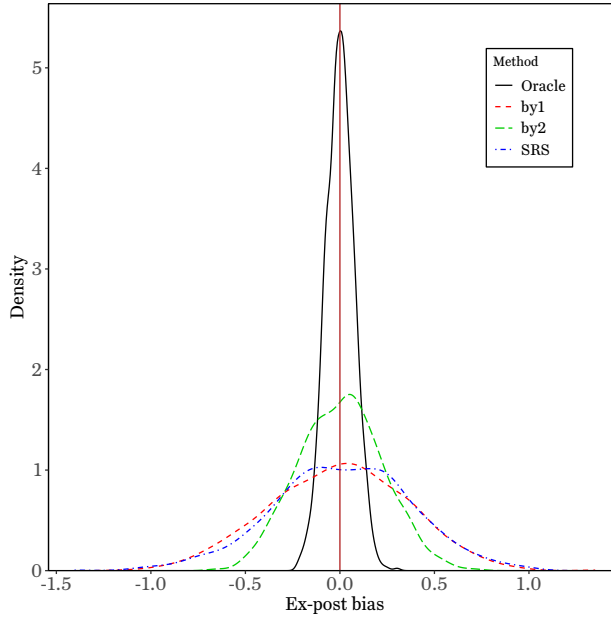


Figure 2: Densities of the distributions of the  $\text{Bias}_{n,\lambda}^{\text{post}}(\hat{\theta}_n|X^{(n)}, D^{(n)})$  over 1000 draws of  $X^{(n)}$  and  $D^{(n)}$  under all treatment assignment schemes.

## 4 Empirical counterparts

The optimal procedure in (18) depends on the function  $g$  defined in (17), which needs to be estimated in practice. Fortunately, pilot experiments are common in RCTs, and we could use data from pilot experiments to estimate  $g$ . In this section, we consider empirical counterparts to the optimal procedure defined by (18), when there is a pilot experiment. We describe the procedures in this section and comment on their asymptotic properties,

formally introducing asymptotic results in Section 5. For any random vector  $A$ , we denote by  $\tilde{A}_j$  the corresponding random vector of the  $j$ th unit in the pilot experiment. Suppose  $\tilde{W}^{(m)} = ((\tilde{Y}_j, \tilde{X}'_j, \tilde{D}_j)' : 1 \leq j \leq m)$  comes from the pilot experiment. We assume that  $((\tilde{Y}_j(1), \tilde{Y}_j(0), \tilde{X}_j) : 1 \leq j \leq m)$  is an i.i.d. sequence of random vectors with distribution  $Q$ , i.e., the units in the pilot are drawn from the same population as the units in the main experiment.

We first consider a plug-in procedure. Suppose  $\hat{g}_m$  is an estimator of  $g$  defined in (17). Concretely,  $\hat{g}_m$  is a random function from  $\mathbf{R}^p$  to  $\mathbf{R}$  that depends on  $\tilde{W}^{(m)}$ . We will abstract away from how  $\hat{g}_m$  is obtained but directly impose conditions on  $\hat{g}_m$  itself. Recall  $\Pi_n$  is the set of all permutations of  $\{1, \dots, 2n\}$  and let  $\pi^{\hat{g}_m} \in \Pi_n$  be such that  $\hat{g}_{m, \pi^{\hat{g}_m}(1)} \leq \dots \leq \hat{g}_{m, \pi^{\hat{g}_m}(2n)}$ . We define the following plug-in stratification for the main experiment:

$$\lambda^{\hat{g}_m}(X^{(n)}) = \{\{\pi^{\hat{g}_m}(2s-1), \pi^{\hat{g}_m}(2s)\} : 1 \leq s \leq n\}. \quad (21)$$

As Theorem 5.1 shows, the plug-in procedure enjoys the property that as the sample size of the pilot increases, the asymptotic variance of  $\hat{\theta}_n$  in (3) is that same as that under the optimal procedure defined by (18). The key condition for the property is that  $\hat{g}_m$  is consistent for  $g$  in a certain sense. See Assumption 5.3 below for more details. The assumption is satisfied by a large class of nonparametric estimation methods, including machine learning methods in high-dimensional settings, i.e., when the dimension of the covariates is large.

When the sample size of the pilot is small, the plug-in procedure generally does not have the efficiency property as in settings with large pilot. But even then, researchers may well be content with the plug-in procedure because it results in smaller limiting variance of  $\hat{\theta}_n$  than many alternatives. That said, we may be concerned that the plug-in estimator  $\hat{g}_m$  is a poor approximation for  $g$  in (17), and as a result, that under the plug-in stratification defined in (21), the conditional MSE and the asymptotic variance of  $\hat{\theta}_n$  is large. Therefore, we consider a penalized procedure under which, according to simulation studies in Section 6, the conditional MSE of  $\hat{\theta}_n$  is often smaller than that under the stratification defined in (21). The procedure is named so because it can be viewed as penalizing the plug-in procedure by the standard error of the plug-in estimate.

We will describe the procedure first and then explain the intuition why it is of this particular form. For  $d \in \{0, 1\}$ , define the least-square estimators based on the treated or control units as

$$\hat{\beta}_m(d) = \left( \sum_{1 \leq j \leq m: \tilde{D}_j = d} \tilde{X}_j \tilde{X}'_j \right)^{-1} \sum_{1 \leq j \leq m: \tilde{D}_j = d} \tilde{X}_j \tilde{Y}_j, \quad (22)$$

and the variance estimators assuming homoskedasticity as

$$\hat{\Sigma}_m(d) = \hat{\nu}_m^2(d) \left( \sum_{1 \leq j \leq m: \tilde{D}_j = d} \tilde{X}_j \tilde{X}'_j \right)^{-1}, \quad (23)$$

where

$$\hat{\nu}_m^2(d) = \frac{\sum_{1 \leq j \leq m} (\tilde{Y}_j - \tilde{X}'_j \hat{\beta}_m(d))^2 I\{\tilde{D}_j = d\}}{\sum_{1 \leq j \leq m} I\{\tilde{D}_j = d\}}.$$

Further define

$$\hat{\beta}_m = \hat{\beta}_m(1) + \hat{\beta}_m(0) \quad (24)$$

$$\hat{\Sigma}_m = \hat{\Sigma}_m(1) + \hat{\Sigma}_m(0). \quad (25)$$

Next, we define  $R_m$  as the result of the following Cholesky decomposition:

$$R'_m R_m = \hat{\beta}_m \hat{\beta}'_m + \hat{\Sigma}_m, \quad (26)$$

and the following transformation of the covariates:

$$Z_i = R_m X_i. \quad (27)$$

The penalized stratification matches units to minimize the sum of distances in terms of  $Z_i$  within pairs. Compared with  $\hat{g}_m(X_i)$ , the main difference is that  $Z_i$  is a vector of the same dimension  $p$  of  $X_i$ , instead of a scalar. Let  $\pi^{\text{pen}}$  denote the solution to the following problem:

$$\min_{\pi \in \Pi_n} \frac{1}{n} \sum_{1 \leq s \leq n} \|Z_{\pi(2s-1)} - Z_{\pi(2s)}\|. \quad (28)$$

When the dimension  $p$  of  $X_i$  is not too large, the problem could be solved quickly by the package `nbpMatching` in R. Finally, define the penalized stratification as

$$\lambda^{\text{pen}}(X^{(n)}) = \{\{\pi^{\text{pen}}(2s-1), \pi^{\text{pen}}(2s)\} : 1 \leq s \leq n\}. \quad (29)$$

(29) can be viewed as penalizing the plug-in procedure in (21) by the variance of the plug-in estimator.

We now briefly explain the intuition behind (28). For simplicity, suppose  $E[Y_i(d)|X_i] = X'_i \beta(d)$  for  $d \in \{0, 1\}$ . In addition, define  $\beta = \beta(1) + \beta(0)$ . (28) penalizes the the plug-in stratification by the standard error of the plug-in estimate. Indeed, the objective in (28) equals

$$\frac{1}{n} \sum_{1 \leq s \leq n} \hat{d}^{\frac{1}{2}}(X_{\pi(2s-1)}, X_{\pi(2s)}),$$

where for any  $x_1, x_2 \in \mathbf{R}^p$ ,

$$\hat{d}(x_1, x_2) = (x'_1 \hat{\beta}_m - x'_2 \hat{\beta}_m)^2 + (x_1 - x_2)' \hat{\Sigma}_m (x_1 - x_2). \quad (30)$$

If  $\hat{\Sigma}_m = 0$ , then (28) is solved by  $\pi^{\hat{g}_m}$  in the plug-in stratification in (21) with  $\hat{g}_m = X'_i \hat{\beta}_m$ . If on the other hand  $\hat{\Sigma}_m$  is large, which means that  $\hat{\beta}_m$  is a very noisy estimate for  $\beta$ , then the second term in (30) dominates, and  $\hat{g}_m$  contributes little to the solution to (28).

**Remark 4.1.** We now provide a further justification for (29) by discussing its optimality in a Bayesian framework. To begin with, note that the problem in (28) could also be defined with the squared norm  $\|Z_{\pi(2s-1)} - Z_{\pi(2s)}\|^2$ , and the two definitions are asymptotically equivalent. For more details, see Section 4 of Bai et al.

(2019). This asymptotically equivalent formulation is in fact optimal in the sense that it minimizes the integrated risk in a Bayesian framework with a diffuse normal prior, where the conditional expectations of potential outcomes are linear. With some abuse of notation, denote the conditional MSE in (4) by  $\text{MSE}(\lambda|g, X^{(n)})$ , where we make explicit the dependence on  $g$ . Suppose we have a prior distribution of  $g$ , denoted by  $F(dg)$ , which is normal. Let  $Q_X^n(dx^{(n)})$  denote the distribution of  $X^{(n)}$  and  $Q_{\tilde{W}}^m(d\tilde{w}^{(m)})$  denote the distribution of  $\tilde{W}^{(m)}$ . Consider the solution to following problem of minimizing the integrated risk across all measurable functions of the form  $u : (\tilde{w}^{(m)}, x^{(n)}) \mapsto \lambda \in \Lambda_n$ :

$$\min_u \iiint \text{MSE}(u(\tilde{w}^{(m)}, x^{(n)})|g, x^{(n)}) Q_X^n(dx^{(n)}) Q_{\tilde{W}}^m(d\tilde{w}^{(m)}) F(dg). \quad (31)$$

In Appendix D, we first show that the problem in (31) under any prior  $F$  is solved by a matched-pair design. To the best of our knowledge, this is the first result showing that matched-pair designs are optimal in general Bayesian frameworks. Next, we specialize the model by assuming  $E[Y_i(d)|X_i] = X_i'\beta(d)$ , define  $\beta = \beta(1) + \beta(0)$ , and show that  $F$  could be equivalently expressed as a distribution on  $\beta$ , which we further assume to be normal. One may be tempted to conjecture that the solution to (31) is to naïvely match units on the the value of  $X_i'\bar{\beta}$ , where  $\bar{\beta}$  is posterior mean of  $\beta$ , i.e.,  $\hat{\beta}_m$  in (24) shrunk towards the prior mean. We show, however, that the solution to (31) depends not only on the posterior mean of  $\beta$ , but also on the posterior variance of it. The posterior variance serves as a penalty to matching naively on the posterior mean of  $\beta$ : the larger the variance, the more it penalizes matching on the posterior mean. In the end, we show that when  $F$  diverges to the diffuse prior, the posterior mean converges to the OLS estimate, and the posterior variance converges to the variance estimate from OLS. As a result, the solution to (31) converges to the procedure defined by (28) with the squared norm  $\|Z_{\pi(2s-1)} - Z_{\pi(2s)}\|^2$ . ■

## 5 Asymptotic results and inference

Under matched-pair designs, it is challenging to derive asymptotic properties of the difference-in-means estimator and conduct inference for ATE, because of the heavy dependence of treatment status across units. Even if  $g$  in (17) is known, commonly-used inference procedures under matched-pair designs, including the two-sample  $t$ -test and the “matched pairs”  $t$ -test, are conservative in the sense that the limiting rejection probability under the null is equal to the nominal level. The issue is further complicated since  $g$  needs to be estimated, so that the stratifications in (21) and (29) depend on data from the pilot experiment. Extending results from Bai et al. (2019), we develop novel results of independent interest on the limiting behavior of the difference-in-means estimator under procedures involving a large number of strata, when the stratifications depend on data from the pilot experiment. These results enable us to establish the desired property of our proposed inference procedures. To begin with, we make the following mild moment restriction on the distributions of potential outcomes:

**Assumption 5.1.**  $E[Y_i^2(d)] < \infty$  for  $d \in \{0, 1\}$ .

## 5.1 Asymptotic results for plug-in with large pilot

In this subsection, we study the properties of  $\hat{\theta}_n$  defined in (3) under settings where the sample sizes of both the pilot and the main experiments increase. We henceforth refer to such a setting as an experiment with a large pilot. We first impose the following assumption on  $g$  defined in (17).

**Assumption 5.2.** The function  $g$  satisfies

- (a)  $0 < E[\text{Var}[Y_i(d)|g(X_i)]]$  for  $d \in \{0, 1\}$ .
- (b)  $\text{Var}[Y_i(d)|g(X_i) = z]$  is Lipschitz in  $z$ .
- (c)  $E[g^2(X_i)] < \infty$ .

Assumption 5.2(a)–(c) are conditions imposed on the target function  $g$  instead of the plug-in estimator  $\hat{g}_m$ . Assumption 5.2(a) is a mild restriction to rule out degenerate situations and to permit the application of suitable laws of large numbers and central limit theorems. Assumption 5.2(c) is another mild moment restriction to ensure the pairs are “close” in the limit. New sufficient conditions for Assumption 5.2(b) are provided in Appendix C.1. The results therein about the conditional expectation of a random variable given a manifold are new and may be of independent interest.

We additionally impose the following restriction on the estimator  $\hat{g}_m$ . In what follows, we use  $Q_X$  to denote the marginal distribution of  $X_i$  under  $Q$ .

**Assumption 5.3.** The sequence of estimators  $\{\hat{g}_m\}$  satisfies

$$\int_{\mathbf{R}^p} |\hat{g}_m(x) - g(x)|^2 Q_X(dx) \xrightarrow{P} 0$$

as  $m \rightarrow \infty$ .

Assumption 5.3 is commonly referred to as the  $L^2$ -consistency of the  $\hat{g}_m$  for  $g$ . When  $p$  is fixed and suitable smoothness conditions hold,  $L^2$ -consistency is satisfied by series and sieves estimators (Newey, 1997; Chen, 2007) and kernel estimators (Li and Racine, 2007). In high-dimensional settings, when  $p$  increases with  $n$  at suitable rates, it is satisfied by the LASSO estimator (Bühlmann and Van De Geer, 2011; Belloni et al., 2012, 2014; Chatterjee, 2013; Bellec et al., 2018), regression trees and random forests (Györfi et al., 2006; Biau, 2012; Denil et al., 2014; Scornet et al., 2015; Wager and Walther, 2015), neural nets (White, 1990; Chen and White, 1999; Chen, 2007; Farrell et al., 2018), and support vector machines (Steinwart and Christmann, 2008). The results therein are either exactly as stated in Assumption 5.3 or one of the following:

- (a)  $\sup_{x \in \mathbf{R}^p} |\hat{g}_m(x) - g(x)| \xrightarrow{P} 0$  as  $m \rightarrow \infty$ .
- (b)  $E[|\hat{g}_m(x) - g(x)|^2] \rightarrow 0$  as  $m \rightarrow \infty$ .

It is straightforward to see (a) implies Assumption 5.3. (b) also implies Assumption 5.3 by Markov’s inequality.



The next theorem reveals that under  $L^2$ -consistency of the estimator  $\hat{g}_m$ , the asymptotic variance of  $\hat{\theta}_n$  under the plug-in procedure is the same with that under the infeasible optimal procedure defined by (18).

**Theorem 5.1.** *Suppose the treatment assignment scheme satisfies Assumption 2.1,  $Q$  satisfies Assumption 5.1,  $g$  satisfies Assumption 5.2. Then, under  $\lambda^g(X^{(n)})$ , as  $n \rightarrow \infty$ ,*

$$\sqrt{n}(\hat{\theta}_n - \theta(Q)) \xrightarrow{d} N(0, \varsigma_g^2),$$

where

$$\varsigma_g^2 = \text{Var}[Y_i(1)] + \text{Var}[Y_i(0)] - \frac{1}{2}E[(g(X_i) - E[Y_i(1) + Y_i(0)])^2]. \quad (32)$$

In addition, suppose  $\hat{g}_m$  satisfies Assumption 5.3. Then, under  $\lambda^{\hat{g}_m}(X^{(n)})$  defined in (21), as  $m, n \rightarrow \infty$ ,

$$\sqrt{n}(\hat{\theta}_n - \theta(Q)) \xrightarrow{d} N(0, \varsigma_g^2).$$

**Remark 5.1.** Bai et al. (2019) studies the scenario where units are matched to minimize the sum of the Euclidean distance in terms of their covariates, and show that the limiting variance of  $\hat{\theta}_n$  is equal to  $\varsigma_g^2$  in (32). The results there, though, are derived assuming that the number of covariates  $p$  is fixed. Instead, we could allow for  $p$  to increase with the sample size  $n$ , as long as  $\hat{g}_m$  is  $L^2$ -consistent for  $g$ . ■

**Remark 5.2.** In analysis of randomized controlled trials researchers sometimes choose not to stratify but run regressions with covariate adjustments afterwards, or more precisely, regress outcomes on treatment and some observed covariates of the units. With the most flexible adjustment, the “efficiency bound” of  $\hat{\theta}_n$  is equal to  $\varsigma_g^2$  in (32). Unfortunately, in order to attain this bound, higher order smoothness conditions are usually required.  $L^2$ -consistency of the estimator seldom suffices, and the estimators are often required not only to be uniformly consistent, but also to have a sufficiently fast convergence rate. See Rothe (2020) for more details. As a result, in practice, researchers often opt to use only linear covariate adjustments. In contrast, the only assumption we require for Theorem 5.1 to hold is the  $L^2$ -consistency of  $\hat{g}_m$ . Furthermore, it is straightforward to combine stratification with covariate adjustments. By using a conventional argument in partitioned regression, one could show that with the same covariate adjustments, stratification will always lead to a weakly smaller limiting variance of  $\hat{\theta}_n$ . A further observation is that sometimes the optimal stratification  $\lambda^g(X^{(n)})$  in (18) is known without any need for estimation. For example, if  $X_i$  is a scalar and  $g(x)$  is known to be monotonic in  $x$ , then we could simply match units according to  $X_i$ ’s, and this would attain  $\varsigma_g^2$  in (32). On the other hand, with covariate adjustments we still need to nonparametrically estimate  $g$  at a sufficiently fast rate. ■

## 5.2 Inference under plug-in procedure

Next, we consider inference for the ATE. For any prespecified  $\theta_0 \in \mathbf{R}$ , we are interested in testing

$$H_0 : \theta(Q) = \theta_0 \text{ versus } H_1 : \theta(Q) \neq \theta_0 \quad (33)$$

at level  $\alpha \in (0, 1)$ . In order to do so, for  $d \in \{0, 1\}$ , define

$$\hat{\sigma}_n^2(d) = \frac{1}{n} \sum_{1 \leq i \leq 2n: D_i=d} (Y_i - \hat{\mu}_n(d))^2.$$

Define

$$\hat{\rho}_n = \frac{2}{n} \sum_{1 \leq j \leq \lfloor \frac{n}{2} \rfloor} (Y_{\pi^{\hat{g}_m}(4j-3)} + Y_{\pi^{\hat{g}_m}(4j-2)})(Y_{\pi^{\hat{g}_m}(4j-1)} + Y_{\pi^{\hat{g}_m}(4j)}) \quad (34)$$

and define  $\hat{\zeta}_n^{\hat{g}_m}$  such that

$$(\hat{\zeta}_n^{\hat{g}_m})^2 = \hat{\sigma}_n^2(1) + \hat{\sigma}_n^2(0) - \frac{1}{2}\hat{\rho}_n + \frac{1}{2}(\hat{\mu}_n(1) + \hat{\mu}_n(0))^2. \quad (35)$$

The test is

$$\phi_n^{\hat{g}_m}(W^{(n)}) = I\{|T_n^{\hat{g}_m}(W^{(n)})| > \Phi^{-1}(1 - \frac{\alpha}{2})\}, \quad (36)$$

where

$$T_n^{\hat{g}_m}(W^{(n)}) = \frac{\sqrt{n}(\hat{\theta}_n - \theta_0)}{\hat{\zeta}_n^{\hat{g}_m}}, \quad (37)$$

and  $\Phi^{-1}(1 - \frac{\alpha}{2})$  denotes the  $(1 - \frac{\alpha}{2})$ -th quantile of the standard normal distribution. Although the right-hand side of (35) is possibly negative, its limit in probability must be positive under assumptions imposed below. By Remark 5.5 below, we could always adjust it to be positive. Therefore, we assume all quantities like (35) are positive for the rest of the paper.

We start by studying the limiting behavior of the test defined in (36) with a large pilot. The following theorem shows that the test defined in (36) is asymptotically exact in the sense that when the sample sizes of both the pilot and the main experiments increase, the limiting rejection probability is equal to the nominal level.

**Theorem 5.2.** *Suppose the treatment assignment scheme satisfies Assumption 2.1,  $Q$  satisfies Assumption 5.1,  $g$  satisfies Assumption 5.2, and  $\hat{g}_m$  satisfies Assumption 5.3. Then, under  $\lambda^{\hat{g}_m}(X^{(n)})$  defined in (21), as  $m, n \rightarrow \infty$ ,*

$$(\hat{\zeta}_n^{\hat{g}_m})^2 \xrightarrow{P} \zeta_g^2.$$

Thus, for the problem of testing (33) at level  $\alpha \in (0, 1)$ ,  $\phi_n^{\hat{g}_m}(W^{(n)})$  defined in (36) satisfies

$$\lim_{m, n \rightarrow \infty} E[\phi_n^{\hat{g}_m}(W^{(n)})] = \alpha,$$

when  $Q$  additionally satisfies the null hypothesis, i.e.,  $\theta(Q) = \theta_0$ .

**Remark 5.3.** The studentization by (35) is crucial for the asymptotic exactness of (36). Commonly-used tests including the two-sample  $t$ -test (Riach and Rich, 2002; Gelman and Hill, 2006; Duflo et al., 2007) and the “matched pairs”  $t$ -test (Moses, 2006; Hsu and Lachenbruch, 2007; Armitage et al., 2008; Athey and Imbens, 2017) are asymptotically conservative in the sense that the limiting rejection probabilities under the null are no greater than and typically strictly less than the nominal level. See Bai et al. (2019) for more details. ■

**Remark 5.4.** In order for  $\hat{g}_m$  to satisfy Assumption 5.3, the following selection on observables condition is

usually required on the pilot experiment:

$$(\tilde{Y}^{(m)}(1), \tilde{Y}^{(m)}(0)) \perp\!\!\!\perp \tilde{D}^{(m)} | \tilde{X}^{(m)},$$

The condition is satisfied by a large class of treatment assignment schemes, including simple random sampling, covariate-adaptive randomization, re-randomization, etc. For more details, see [Bugni et al. \(2018\)](#) and [Bai et al. \(2019\)](#). ■

**Remark 5.5.** In finite sample one might be worried that the right hand side of (35) is negative. Furthermore, we always have access to an asymptotically conservative estimator for the limiting variance, for example,  $\hat{\zeta}_n^2 = \hat{\sigma}_n^2(1) + \hat{\sigma}_n^2(0)$ , whose probability limit is weakly greater than  $\zeta_g^2$ . So even though the right hand side of (35) is positive, it might be larger than  $\zeta_n^2$  in finite sample. To get over both problems, we could simply redefine the variance estimator to be  $\hat{\zeta}_n^2$  if the right hand side of (35) is less than or equal to 0, and the smaller one of the right hand side of (35) and  $\hat{\zeta}_n^2$  otherwise. ■

Next, we consider settings where the sample size of the main experiment increases while that of the pilot experiment is allowed to be fixed. We henceforth refer to such a setting as an experiment with a small pilot. We show that test defined in (36) is again asymptotically exact in the sense that the limiting rejection probability under the null is equal to the nominal level when the sample size of the main experiment increases, regardless of the sample size of the pilot. The restrictions that we put on  $\hat{g}_m$ , however, are more likely to be satisfied when  $\hat{g}_m$  is constructed using simple methods such as least squares. We impose the following restriction in addition to Assumption 5.1:

**Assumption 5.4.** The estimator  $\hat{g}_m$  satisfies

$$Q\{\hat{g}_m \in \mathbf{H}\} = 1,$$

where  $\mathbf{H}$  is the set of all measurable functions  $h : \mathbf{R}^p \rightarrow \mathbf{R}$  such that

- (a)  $0 < E[\text{Var}[Y_i(d)|h(X_i)]]$  for  $d \in \{0, 1\}$ .
- (b)  $E[Y_i^r(d)|h(X_i) = z]$  is Lipschitz in  $z$  for  $r = 1, 2$  and  $d = 0, 1$ .
- (c)  $E[h^2(X_i)] < \infty$ .

Assumption 5.4 is imposed on the distributions of potential outcomes conditional on  $\hat{g}_m$ , where  $\hat{g}_m$  is viewed as a fixed function given data from the pilot experiment. In fact, with small pilots, Assumption 5.4 contains the same set of conditions as those in Assumption 5.2, the only difference being that they are imposed on  $\hat{g}_m$  instead of  $g$ . In the definition of  $\mathbf{H}$ , (a) is a mild restriction to rule out degenerate situations and to permit the application of suitable laws of large numbers and central limit theorems, and (c) is another mild moment restriction to ensure the pairs are “close” in the limit. New sufficient conditions for (b) are provided in Appendix C.1. Note, in particular, that (b) is more likely to be satisfied when  $\hat{g}_m$  is constructed using simple estimation methods such as least squares.

The following theorem shows that the test defined in (36) is asymptotically exact in the sense that as the sample size of the main experiment increases, the limiting rejection probability under the null is equal to the nominal level. Note, in particular, that the sample size of the pilot is allowed to be fixed.

**Theorem 5.3.** *Suppose the treatment assignment scheme satisfies Assumption 2.1,  $Q$  satisfies Assumption 5.1, and  $\hat{g}_m$  satisfies Assumption 5.4. Suppose  $Q$  additionally satisfies the null hypothesis, i.e.,  $\theta(Q) = \theta_0$ . Then, under  $\lambda^{\hat{g}_m}(X^{(n)})$  defined in (21), for the problem of testing (33) at level  $\alpha \in (0, 1)$ ,  $\phi_n^{\hat{g}_m}(W^{(n)})$  defined in (36) satisfies*

$$\lim_{n \rightarrow \infty} E[\phi_n^{\hat{g}_m}(W^{(n)})] = \alpha.$$

**Remark 5.6.** Note that we use the same test  $\phi_n^{\hat{g}_m}$  with large (Theorem 5.2) and small (Theorem 5.3) pilots, and it is asymptotically exact either way. When  $m$  increases at a rate such that Assumption 5.3 is satisfied, the asymptotic variance of  $\hat{\theta}_n$  as  $m, n \rightarrow \infty$  is  $\varsigma_g^2$ , which equals the asymptotic variance under the infeasible optimal procedure defined by (18). Yet when  $m$  is fixed, the asymptotic variance of  $\hat{\theta}_n$  as  $n \rightarrow \infty$  is generally larger than  $\varsigma_g^2$ . Moreover, as previously commented, the assumptions in the two settings are non-nested. Assumption 5.4 is more likely to be satisfied when the plug-in estimator  $\hat{g}_m$  is constructed using simple estimation methods, but does not require  $\hat{g}_m$  to be consistent for  $g$  in any sense. On the other hand, Assumptions 5.2 and Assumption 5.3 could potentially allow for more complicated estimation methods but require  $\hat{g}_m$  to be  $L^2$ -consistent for  $g$ .

**Remark 5.7.** In fact, the asymptotic exactness of  $\phi_n^{\hat{g}_m}(W^{(n)})$  holds conditional on data from the pilot experiment, i.e.,

$$\lim_{n \rightarrow \infty} E[\phi_n^{\hat{g}_m}(W^{(n)}) | \tilde{W}^{(m)}] = \alpha \quad (38)$$

with probability one for  $\tilde{W}^{(m)}$ . See the proof of Theorem 5.3 in the appendix for more details. Furthermore, it follows from the proof that the test is also asymptotically exact under

$$\lambda^h(X^{(n)}) = \{\{\pi^h(2s-1), \pi^h(2s)\} : 1 \leq s \leq n\}, \quad (39)$$

where  $h_{\pi^h(1)} \leq \dots \leq h_{\pi^h(2n)}$  and  $h$  is a fixed function satisfying  $h \in \mathbf{H}$  for  $\mathbf{H}$  defined in (5.4). ■

**Remark 5.8.** As an intermediate step in the proof of Theorem 5.3, we derive the limiting variance of  $\hat{\theta}_n$  under  $\lambda^h(X^{(n)})$  defined in (39), where  $h$  is a fixed function satisfying  $h \in \mathbf{H}$ . The limiting variance equals

$$\varsigma_h^2 = \text{Var}[Y_i(1)] + \text{Var}[Y_i(0)] - \frac{1}{2} E[(E[Y_i(1) + Y_i(0) | h(X_i)] - E[Y_i(1) + Y_i(0)])^2]. \quad (40)$$

Comparing (40) with (32), we could show the minimum of  $\varsigma_h^2$  over  $h \in \mathbf{H}$  occurs when  $h = g$ , and the minimum is unique unless there exists and  $h \in \mathbf{H}$  for which  $E[Y_i(1) + Y_i(0) | h(X_i)] = E[Y_i(1) + Y_i(0) | X_i]$  with probability one. This result enables us to compare the limiting variance of  $\hat{\theta}_n$  across a large class of stratifications, and in particular, all stratifications with a fixed number of large strata. Indeed, all such stratifications could be defined by a discrete-valued function  $h : \mathbf{R}^p \rightarrow \{1, \dots, R\}$  for a fixed integer  $R$ , and therefore  $\varsigma_h^2 \geq \varsigma_g^2$  unless  $E[Y_i(1) + Y_i(0) | h(X_i) = r] = E[Y_i(1) + Y_i(0) | X_i]$  with probability one, i.e., when  $E[Y_i(1) + Y_i(0) | X_i]$  is the same within each stratum. Another corollary is that if  $h \in \mathbf{H}$  and  $h_c$  is a constant function, then the stratification  $\lambda^{h_c}(X^{(n)}) = \{\{1, \dots, 2n\}\}$  with all units in one stratum satisfies  $\varsigma_{h_c}^2 \geq \varsigma_h^2$ , unless again the degeneracy

condition holds, this time requiring  $E[Y_i(1) + Y_i(0)|h(X_i)]$  to be a constant. Any  $\hat{g}_m$  with  $Q\{\hat{g}_m \in \mathbf{H}\} = 1$  is a constant function in  $\mathbf{H}$  conditional on the pilot data  $\tilde{W}^{(m)}$ , so in this sense, almost all stratifications are better than not stratifying at all, because it results in a weakly smaller and typically strictly smaller limiting variance of  $\hat{\theta}_n$ . See Theorem B.2 for more details. By direct calculation we could also show that for any  $h \in \mathbf{H}$ ,  $\varsigma_h^2$  is weakly less than and typically strictly less than the limiting variance of  $\hat{\theta}_n$  under simple random sampling, i.e., when treatment status is determined by i.i.d. coin flips. ■

**Remark 5.9.** Sometimes political or logistical considerations or estimation of subpopulation treatment effects require researchers to prespecify different treated fractions across subpopulations. In those settings, as discussed in Appendix B,  $\hat{\theta}_n$  is no longer consistent for  $\theta$  in (1). Instead, it is natural to use the estimator from the fully saturated regression with all interaction terms of treatment status and strata indicators, i.e.,  $\hat{\theta}_n^{\text{sat}}$  defined in (62). Appendix B discusses straightforward extensions of the optimality result in Theorem 3.1 and empirical counterparts including that in (21). These results are closely related to Tabord-Meehan (2020), who considers stratification trees which lead to a small number of large strata. In particular, Remark B.1 discusses a way to combine his procedure and procedures in this paper, under which the asymptotic variance of  $\hat{\theta}_n^{\text{sat}}$  is no greater than and typically strictly less than that under his procedure alone. ■

### 5.3 Inference under penalized procedure

We now consider inference under the penalized procedure defined by (29) with a small pilot. This subsection follows closely the exposition in Section 4 of Bai et al. (2019). Since in general  $Z$  defined in (27) is not a scalar, the correction term in (34) could no longer be defined as before since it relies on  $\pi^{\hat{g}_m}$ , where  $\hat{g}_m$  is a scalar. Instead, we need to match the pairs to ensure that the two pairs matched are close in terms of  $Z$ . Define

$$\bar{Z}_s = \frac{Z_{\pi^{\text{pen}}(2s-1)} + Z_{\pi^{\text{pen}}(2s)}}{2},$$

and  $\bar{\pi}$  as the solution of the following problem:

$$\min_{\pi \in \Pi_n} \frac{1}{n} \sum_{1 \leq j \leq \lfloor \frac{n}{2} \rfloor} \|\bar{Z}_{\pi(2j-1)} - \bar{Z}_{\pi(2j)}\|.$$

Let  $\tilde{\pi}^{\text{pen}} \in \Pi_n$  be such that for  $1 \leq s \leq n$ ,

$$\tilde{\pi}^{\text{pen}}(2s-1) = \pi^{\text{pen}}(2\bar{\pi}(s)-1) \text{ and } \tilde{\pi}^{\text{pen}}(2s) = \pi^{\text{pen}}(2\bar{\pi}(s)).$$

In other words,  $\tilde{\pi}^{\text{pen}}$  matches the pairs defined by  $\pi^{\text{pen}}$  based on the midpoints of pairs. Since  $\tilde{\pi}^{\text{pen}}$  rearranges  $\pi^{\text{pen}}$  in (29) while preserving the units in each stratum, it follows that for  $\lambda^{\text{pen}}(X^{(n)})$  defined in (29), we have

$$\lambda^{\text{pen}}(X^{(n)}) = \{ \{ \tilde{\pi}^{\text{pen}}(2s-1), \tilde{\pi}^{\text{pen}}(2s) \} : 1 \leq s \leq n \}.$$

We then define the test similarly to (36), with  $\pi^{\hat{g}_m}$  replaced by  $\tilde{\pi}^{\text{pen}}$ . In particular, define

$$\hat{\rho}_n^{\text{pen}} = \frac{2}{n} \sum_{1 \leq j \leq \lfloor \frac{n}{2} \rfloor} (Y_{\tilde{\pi}^{\text{pen}}(4j-3)} + Y_{\tilde{\pi}^{\text{pen}}(4j-2)})(Y_{\tilde{\pi}^{\text{pen}}(4j-1)} + Y_{\tilde{\pi}^{\text{pen}}(4j)})$$

and let  $\hat{\zeta}_n^{\text{pen}}$  be such that

$$(\hat{\zeta}_n^{\text{pen}})^2 = \hat{\sigma}_n^2(1) + \hat{\sigma}_n^2(0) - \frac{1}{2}\hat{\rho}_n^{\text{pen}} + \frac{1}{2}(\hat{\mu}_n(1) + \hat{\mu}_n(0))^2.$$

The test is

$$\phi_n^{\text{pen}}(W^{(n)}) = I\{|T_n^{\text{pen}}(W^{(n)})| > \Phi^{-1}(1 - \frac{\alpha}{2})\}, \quad (41)$$

where

$$T_n^{\text{pen}}(W^{(n)}) = \frac{\sqrt{n}(\hat{\theta}_n - \theta_0)}{\hat{\zeta}_n^{\text{pen}}}, \quad (42)$$

and  $\Phi^{-1}(1 - \frac{\alpha}{2})$  denotes the  $(1 - \frac{\alpha}{2})$ -th quantile of the standard normal distribution.

Under the penalized procedure, we impose the following assumption on  $Q$ :

**Assumption 5.5.** (a)  $0 < E[\text{Var}[Y_i(d)|R_m X_i]]$  for  $d \in \{0, 1\}$ .

(b)  $E[Y_i^r(d)|R_m X_i = z]$  is Lipschitz in  $z$  for  $r \in \{1, 2\}$  and  $d \in \{0, 1\}$ .

(c) The support of  $R_m X_i$  is compact.

Assumption 5.5(a)–(b) are the counterparts to Assumption 2.1(a) and (c) of Bai et al. (2019). Assumption 5.5(c) is also imposed in Section 4 of Bai et al. (2019). The following theorem establishes the asymptotic exactness of the test defined in (41), in the sense that the limiting rejection probability under the null equals the nominal level. Note, in particular, that the sample size of the pilot is allowed to be fixed.

**Theorem 5.4.** *Suppose the treatment assignment scheme satisfies Assumption 2.1 and  $Q$  satisfies Assumptions 5.1 and 5.5. Suppose  $Q$  additionally satisfies the null hypothesis, i.e.,  $\theta(Q) = \theta_0$ . Then, under  $\lambda^{\text{pen}}(X^{(n)})$  defined in (29), for the problem of testing (33) at level  $\alpha \in (0, 1)$ ,  $\phi_n^{\text{pen}}(W^{(n)})$  defined in (36) satisfies*

$$\lim_{n \rightarrow \infty} E[\phi_n^{\text{pen}}(W^{(n)})] = \alpha.$$

**Remark 5.10.** In some setups, it may be possible to improve the estimator  $\hat{g}_m$  by imposing shape restrictions on  $g$ . See, for instance, Chernozhukov et al. (2015) and Chetverikov et al. (2018). ■

## 5.4 Inference with pooled data

So far we have disregarded data from the pilot experiment in the test defined in (36) except when computing  $\hat{g}_m$ . We end this section by describing a test that combines data from the pilot and the main experiments. Define

$$\tilde{\theta}_m = \tilde{\mu}_m(1) - \tilde{\mu}_m(0),$$

where

$$\tilde{\mu}_m(d) = \frac{\sum_{1 \leq j \leq m} \tilde{Y}_j I\{\tilde{D}_j = d\}}{\sum_{1 \leq j \leq m} I\{\tilde{D}_j = d\}}$$

for  $d \in \{0, 1\}$ . We define the new estimator for  $\theta(Q)$  as

$$\hat{\theta}_n^{\text{combined}} = \frac{m}{m+2n} \tilde{\theta}_m + \frac{2n}{2n+m} \hat{\theta}_n.$$

We define the test as

$$\phi_n^{\text{combined}}(W^{(n)}, \tilde{W}^{(m)}) = I\{|T_n^{\text{combined}}(W^{(n)}, \tilde{W}^{(m)})| > \Phi^{-1}(1 - \frac{\alpha}{2})\}, \quad (43)$$

where

$$T_n^{\text{combined}}(W^{(n)}, \tilde{W}^{(m)}) = \frac{\sqrt{m+2n}(\hat{\theta}_n^{\text{combined}} - \theta_0)}{\sqrt{\frac{m}{m+2n} \zeta_{\text{pilot},m}^2 + \frac{2n}{m+2n} 2(\zeta_n^{\hat{g}_m})^2}}, \quad (44)$$

and  $\Phi^{-1}(1 - \frac{\alpha}{2})$  denotes the  $(1 - \frac{\alpha}{2})$ -th quantile of the standard normal distribution.

The following theorem shows that the test defined in (43) is asymptotically exact as the sample sizes of both the pilot and the main experiments increase. The main additional requirement is that as  $m \rightarrow \infty$ ,  $\sqrt{m}(\tilde{\theta}_m - \theta(Q))$  converges in distribution to a normal distribution whose variance is consistently estimable. The assumption is satisfied by many treatment assignment schemes, including simple random sampling and covariate-adaptive randomization. See [Bugni et al. \(2018\)](#) and [Bugni et al. \(2019\)](#) for more details.

**Theorem 5.5.** *Suppose the treatment assignment scheme satisfies Assumption 2.1,  $Q$  satisfies Assumptions 5.1,  $g$  satisfies Assumption 5.2, and  $\hat{g}_m$  satisfies Assumption 5.3. Suppose in addition that as  $m \rightarrow \infty$ ,  $\sqrt{m}(\tilde{\theta}_m - \theta(Q)) \xrightarrow{d} N(0, \zeta_{\text{pilot}}^2)$ ,  $\zeta_{\text{pilot},m}^2 \xrightarrow{P} \zeta_{\text{pilot}}^2$ , and that as  $m, n \rightarrow \infty$ ,*

$$\frac{m}{m+2n} \rightarrow \nu \in [0, 1].$$

Then, under  $\lambda^{\hat{g}_m}(X^{(n)})$  defined in (21), as  $m, n \rightarrow \infty$ ,

$$\frac{\sqrt{m+2n}(\hat{\theta}_n^{\text{combined}} - \theta(Q))}{\sqrt{\frac{m}{m+2n} \zeta_{\text{pilot},m}^2 + \frac{2n}{m+2n} 2(\zeta_n^{\hat{g}_m})^2}} \xrightarrow{d} N(0, 1).$$

Thus, for the problem of testing (33) at level  $\alpha \in (0, 1)$ ,  $\phi_n^{\text{combined}}(W^{(n)}, \tilde{W}^{(m)})$  in (43) satisfies

$$\lim_{m,n \rightarrow \infty} E[\phi_n^{\text{combined}}(W^{(n)}, \tilde{W}^{(m)})] = \alpha,$$

whenever  $Q$  additionally satisfies the null hypothesis, i.e.  $\theta(Q) = \theta_0$ .

**Remark 5.11.** Although Theorem 5.5 is stated under  $\lambda^{\hat{g}_m}(X^{(n)})$  in (21), it is straightforward to establish similar results when  $\lambda^{\hat{g}_m}(X^{(n)})$  in the main experiment is replaced by other stratifications, e.g., (29). ■

## 6 Simulation

In this section, we examine the properties of the procedures discussed in Section 4 in a small simulation study. For  $d \in \{0, 1\}$  and  $1 \leq i \leq 2n$ , potential outcomes are generated according to the equation:

$$Y_i(d) = \mu(d) + m_d(X_i) + \sigma_d(X_i)\epsilon_i(d),$$

where  $\mu(d)$ ,  $m_d(X_i)$ ,  $\sigma_d(X_i)$ , and  $\epsilon_i(d)$  are specified in each model as follows. In each of the following specifications,  $2n = 200$ ;  $((X_i, \epsilon_i(0), \epsilon_i(1)) : 1 \leq i \leq 2n)$  are i.i.d.;  $X_i, \epsilon_i(0), \epsilon_i(1)$  are independent; and  $\mu(0) = 0$ . For each model, we generate data from a very small pilot experiment of sample size  $m = 20$ , in which half of the units are treated.

**Model 1**  $p = 2$ ;  $X_{i,1} \sim \text{Beta}(2, 2)$ ,  $X_{i,2} \sim \text{Beta}(2, 2)$ ;  $m_d(X_i) = X_i' \beta(d)$  and  $\epsilon_i(d) \sim N(0, 1)$  for  $d \in \{0, 1\}$ ;  $\beta(1) = \beta(0) = (1, 1)'$ ;  $\sigma_0(X_i) = \sigma_1(X_i) = 0.1$ .

**Model 2** as in Model 1, but  $\beta(1) = \beta(0) = (3, 0.1)'$ .

**Model 3** as in Model 1, but  $\sigma_0(X_i) = \sigma_1(X_i) = 1$  and  $\epsilon_i(d) \sim \text{Unif}[-\frac{1}{2}, \frac{1}{2}]$  for  $d \in \{0, 1\}$ .

**Model 4** as in Model 2, but  $\sigma_0(X_i) = \sigma_1(X_i) = 1$  and  $\epsilon_i(d) \sim \text{Unif}[-\frac{1}{2}, \frac{1}{2}]$  for  $d \in \{0, 1\}$ .

**Model 5** as in Model 1, but  $m_1(X_i) = m_0(X_i) = X_{i,1}^2$ ,  $\sigma_0(X_i) = \sigma_1(X_i) = 0.1$ , and  $\epsilon_i(d) \sim N(0, 1)$  for  $d \in \{0, 1\}$ .

**Model 6** as in Model 5, but  $m_1(X_i) = m_0(X_i) = X_{i,1}^2 + X_{i,2}^2$ .

Model 1 is a symmetric model with small variances in error terms. Model 2 differs from Model 1 in that  $X_{i,1}$  is the predominant component in potential outcomes. Models 3 and 4 are similar to Models 1 and 2, the only difference being that the error terms have larger variances. Models 5 and 6 are non-linear and are designed to study properties of the plug-in and the penalized procedures under misspecification. In Model 5, only  $X_{i,1}$  affects the potential outcomes, while  $X_{i,1}$  and  $X_{i,2}$  are symmetric in Model 6.

We consider the following procedures:

**Oracle** matched-pair design with the infeasible optimal stratification in (18).

**Plug-in** matched-pair design with the plug-in stratification in (21) with  $\hat{g}_m(x) = x' \hat{\beta}_m$  for  $\hat{\beta}_m$  in (24).

**Pen** matched-pair design with the penalized stratification in (29).

**MPeuc** matched-pair design minimizing the sum of Euclidean distances within pairs.

**by1** stratified randomization with two strata separated by the sample median of  $X_{i,1}$ .

**by2** stratified randomization with two strata separated by the sample median of  $X_{i,2}$ .

**MP1** matched-pair design using  $X_{i,1}$  only, i.e., stratification in (21) with  $\hat{g}_m(x) = x_1$ .



**MP2** matched-pair design using  $X_{i,2}$  only, i.e., stratification in (21) with  $\hat{g}_m(x) = x_2$ .

Stratifications in **Pen** and **MPeuc** are computed using the package `nbpMatching` in R.

We first present results on the conditional MSE of  $\hat{\theta}_n$  defined in (4). In these results, we set  $\mu(1) = \mu(0) = 0$ , so that  $\theta(Q) = 0$  as well. By Lemma 3.1 and in particular (16), the conditional MSEs of  $\hat{\theta}_n$  under stratifications differ only in terms of the variance of the ex-post bias conditional on the covariates. Therefore, for a given stratification  $\lambda$ , a set of covariates  $X^{(n)}$ , and the function  $g$  defined in (17), we define a constant multiple of the objective in (16) as the loss:

$$L(\lambda|g, X^{(n)}) = 4n^2 \text{Var}_\lambda[E[\hat{\theta}_n|X^{(n)}, D^{(n)}]|X^{(n)}]. \quad (45)$$

Table 1 displays the summary statistics of the values of the loss defined in (45) for different stratifications across 1000 draws of  $X^{(n)}$ . We label the columns according to the procedures. In each model, we calculate ratios of values of the loss for each procedure against those for **Oracle**, and present the quartiles and means of the ratios across the 1000 draws of  $X^{(n)}$ .

Model		<b>Oracle</b>	<b>Plug-in</b>	<b>Pen</b>	<b>MPeuc</b>	<b>by1</b>	<b>by2</b>	<b>MP1</b>	<b>MP2</b>
1	25%	1.00	2.50	3.69	22.51	2344.62	2353.34	885.77	903.36
	50%	1.00	8.46	5.76	35.86	3852.52	3848.06	1455.54	1435.83
	75%	1.00	28.03	9.93	55.50	5853.40	5866.36	2238.42	2183.49
	Mean	1.00	25.07	8.22	40.76	4281.19	4293.87	1653.90	1641.51
2	25%	1.00	2.08	4.33	67.39	3238.83	10723.31	6.89	5192.29
	50%	1.00	5.34	5.96	86.24	4211.93	14112.38	8.48	6954.55
	75%	1.00	15.21	9.53	108.13	5239.90	17414.65	10.57	8640.57
	Mean	1.00	12.85	8.14	89.26	4305.93	14377.14	8.90	7169.01
3	25%	1.00	16.28	8.57	22.52	2329.58	2340.74	894.50	902.57
	50%	1.00	68.97	14.04	35.55	3835.03	3850.74	1455.64	1466.20
	75%	1.00	230.33	25.64	54.02	5734.63	5783.08	2230.34	2226.22
	Mean	1.00	205.52	21.42	40.65	4288.67	4299.91	1650.97	1662.17
4	25%	1.00	8.86	10.27	67.58	3266.09	10924.10	6.91	5440.39
	50%	1.00	43.88	15.49	87.50	4125.96	13824.46	8.57	6847.59
	75%	1.00	131.81	26.43	109.16	5197.76	17364.76	10.65	8744.17
	Mean	1.00	104.72	22.07	89.41	4291.05	14343.10	8.97	7168.34
5	25%	1.00	27.39	71.83	415.34	19128.24	57595.61	1.00	27631.81
	50%	1.00	116.62	103.72	501.70	22248.95	66572.16	1.00	32579.67
	75%	1.00	333.13	176.04	599.89	26430.16	77215.74	1.00	38871.75
	Mean	1.00	318.20	150.85	520.67	23158.28	68653.31	1.00	34162.98
6	25%	1.00	244.36	115.27	214.18	27727.82	11878.77	13124.19	1424.60
	50%	1.00	342.09	150.88	265.06	32936.14	14190.12	15817.15	1726.21
	75%	1.00	517.14	197.98	328.09	39810.35	17243.41	18864.44	2118.38
	Mean	1.00	424.81	168.61	276.24	34031.92	14659.61	16327.06	1798.22

Table 1: Summary statistics for ratios of the values of the loss in (45) under all stratifications against those under the infeasible optimal stratifications (**Oracle**), over 1000 draws of  $X^{(n)}$ , in Models 1-6.

Unsurprisingly, **Oracle** always has the smallest values of the loss. Ad-hoc procedures including **by1**, **by2**, **MP1**, **MP2** perform miserably most of the time. Although **MP1** performs well under Models 2, 4, and 5, it is because there  $X_{i,1}$  is a predominant element of potential outcomes. In particular, Model 5 is an example where

$g$  defined in (17) is a monotonic function of the first covariate, so that **MP1** solves (5) and has the same values of loss with **Oracle**. We separately discuss the remaining three procedures, **Plug-in**, **Pen**, and **MPeuc**:

**Plug-in**: In most models, **Plug-in** outperforms ad-hoc procedures including **by1**, **by2**, **MP1**, **MP2**, which is somewhat surprising since the sample size of pilot is only  $m = 20$ . In Models 1–2, where the variances of  $\epsilon_i(d)$ 's are small, **Plug-in** also improves upon **MPeuc**, and the improvement is pronounced in Model 2. But when the variances of  $\epsilon_i(d)$ 's are large, it performs worse than **Pen** and **MPeuc**, as could be seen from Models 3–6.

**Pen**: In Models 1–4, **Pen** is the best among all procedures. In all models, it performs better than **Plug-in** and **MPeuc**, remarkably so than **Plug-in** in Models 3–6. The improvement upon **MPeuc** is most pronounced in Models 2 and 4, where  $X_{i,2}$  contributes little to potential outcomes. These are examples in which **MPeuc** assigns equal weights to two covariates while regression-based methods could detect that one of them dominates. Even when potential outcomes are non-linear (Models 5–6), the values of its loss are smaller than those under **MPeuc**.

**MPeuc**: In all models, it is not as poor as the ad-hoc procedures including **by1**, **by2**, **MP1**, **MP2**, but is obviously worse than **Pen**. In Models 2 and 4, where only  $X_{i,1}$  matters, it is obviously worse than **Pen** and **Plug-in**, because the pilot informs us that  $X_{i,1}$  is much more important than  $X_{i,2}$ , which is not taken into account by Euclidean matching.

Next, for  $\theta_0 = 0$ , we consider the problem of testing (33) at level  $\alpha = 0.05$ . For Models 1–6, we compute the rejection probabilities of suitable tests under stratifications mentioned previously, when  $\mu(0) = 0$  and  $\theta = \mu(1) = 0, 0.0.1, 0.02, 0.04$ . In particular, we use the following tests under each stratification:

**Oracle**: test in (36) with  $\hat{g}_m = g$  for  $g$  defined in (17).

**Plug-in**: test in (36) with  $\hat{g}_m(x) = x' \hat{\beta}_m$  for  $\hat{\beta}_m$  defined in (24).

**Pen**: test in (41).

**MPeuc**: test in (41) with  $Z$  replaced by  $X$ .

**by1**: test in (36) with  $\hat{g}_m(x) = I\{x_1 > \text{med}(X_{i,1} : 1 \leq i \leq 2n)\}$ .

**by2**: test in (36) with  $\hat{g}_m(x) = I\{x_2 > \text{med}(X_{i,2} : 1 \leq i \leq 2n)\}$ .

**MP1**: test in (36) with  $\hat{g}_m(x) = x_1$ .

**MP2**: test in (36) with  $\hat{g}_m(x) = x_2$ .

Table 2 displays the rejection probabilities for Models 1–6 under all stratifications using tests described above. Note that loss properties in Table 1 translate into power properties in Table 2. Indeed, while all tests under all stratifications have correct sizes, the test in (41) under the penalized stratification in (29) has higher

power than most other tests under other stratifications, except that under **Oracle**. In Models 1–2, the corresponding tests under **Plug-in** and **Pen** have higher power than that under **MPeuc**, while being comparable in other models, except in Model 6, where potential outcomes are highly non-linear. The comparison is most pronounced in Model 2, where  $g$  in (17) depends mostly on  $x_1$ , because **Plug-in** and **Pen** incorporate information from the pilot while **MPeuc** doesn't. The test under **Pen** performs better than that under **Plug-in** in Models 1–5. Finally, note that tests under matched-pair designs, including **Plug-in**, **Pen**, and **MPeuc** usually perform much better than tests under stratifications with a small number of large strata, including **by1** and **by2**.

Model		Oracle	Plug-in	Pen	MPeuc	by1	by2	MP1	MP2
1	$\theta = 0$	5.63	5.15	5.61	5.48	5.02	5.27	5.44	5.45
	$\theta = 0.01$	11.21	10.63	11.2	11	6.34	6.41	6.15	6.24
	$\theta = 0.02$	30.26	28.32	29.76	27.31	8.02	8.19	9.83	9.6
	$\theta = 0.04$	79.44	76.86	79.98	75.4	17.71	18.12	20.87	23.19
2	$\theta = 0$	5.43	5.05	5.12	5.24	5.37	5.47	5.32	5.88
	$\theta = 0.01$	11.72	10.84	11.06	9.68	5.54	5.57	10.96	5.53
	$\theta = 0.02$	28.52	27.45	27.88	20.5	7.35	5.6	27.14	5.81
	$\theta = 0.04$	79.82	76.23	78.6	62.6	11.98	6.79	77.77	7.19
3	$\theta = 0$	5.08	5.61	5.32	5.34	5.51	5.7	5.37	5.26
	$\theta = 0.01$	5.69	6.11	6.33	5.58	5.93	5.46	5.51	5.57
	$\theta = 0.02$	8.22	7.49	8.18	8.43	6.92	6.92	7.27	7.67
	$\theta = 0.04$	17.52	16.66	16.94	16.84	11.82	12.31	12.67	12.84
4	$\theta = 0$	5.69	5.55	5.7	5.31	5.43	5.16	5.2	5.14
	$\theta = 0.01$	6.31	6.2	6.69	5.98	5.72	5.49	6.32	5.72
	$\theta = 0.02$	8.1	7.98	8.13	7.87	6.97	5.91	8.05	5.88
	$\theta = 0.04$	16.73	16.77	17.02	16.75	9.69	7.28	16.81	7.28
5	$\theta = 0$	5.33	5.26	5.66	5.5	5.47	5.38	5.6	5.16
	$\theta = 0.01$	11.44	10.93	11.57	11.5	7.78	6.56	11.64	6.5
	$\theta = 0.02$	30.34	28.2	30.02	28.44	14.36	9.28	30.02	9.23
	$\theta = 0.04$	80.81	77.12	79.89	77.46	40.39	20.83	80.52	21.93
6	$\theta = 0$	5.15	5.47	3.51	4.94	5.57	5.78	5.78	5.72
	$\theta = 0.01$	6.77	6.84	4.44	6.46	5.72	5.7	5.62	6.52
	$\theta = 0.02$	12.41	11.49	8.72	11.22	6.79	7.91	6.69	10.55
	$\theta = 0.04$	31.94	29.34	24.37	29.18	10.45	16.31	10.94	25.43

Table 2: Rejection probabilities for Models 1–6 under all stratifications using tests in Section 4.

## 7 Empirical application

To illustrate our procedures in practice, we replicate part of the experiment in DellaVigna and Pope (2018) on Amazon Mechanical Turk (MTurk) and the TurkPrime Prime Panels, using the penalized procedure defined by (29). MTurk is an online crowdsourcing platform widely used to conduct economic and behavioral experiments. For more information about running experiments on Amazon MTurk, see Horton et al. (2011), Mason and Suri (2012), Paolacci and Chandler (2014), Kuziemko et al. (2015), and Litman et al. (2017). Prime Panels is another online platform with over 30 million participants and their reliable demographics.

DellaVigna and Pope (2018) run a large-scale experiment to compare the effectiveness of multiple incentives for efforts in one setting, as well as compare experimental results with expert forecasts. The 18 treat-

ments include various monetary and behavioral incentives. We focus on one of the treatments, which is a monetary incentive. In the experiment, subjects are asked to alternately press the “a” and “b” buttons on their keyboard as quickly as possible in 10 minutes. One alternate press counts as 1 point. All subjects are paid some base rate upon finishing the experiment. In the treatment we replicate, subjects in the treated group are paid an extra \$0.01 for every 100 points they score, while subjects in the control group receive no extra payment. In [DellaVigna and Pope \(2018\)](#), the base payment is \$1, but we use about \$1.25 in the pilot and \$2 in the main experiment to minimize attrition. In our notation, the outcome  $Y$  is the points scored, the treatment  $D$  indicates whether the subject receives extra payment ( $D = 1$ ) or not ( $D = 0$ ). The covariates  $X$  include a constant term, age, gender, ethnicity, education, and income. We re-index gender and ethnicity as binary variables and regard the rest as continuous.

The sample size in the original experiment in [DellaVigna and Pope \(2018\)](#) is 1098. In the original experiment, all the units are in one stratum and the treated fraction is approximately  $\frac{1}{2}$ . There is a pilot experiment in the preregistration stage but the results used in neither designing the main experiment nor analysis in their paper. In our replication, we perform the pilot experiment on Prime Panels and the main experiment on MTurk. The sample size of the pilot experiment is  $m = 44$ , and that of the main experiment is  $2n = 176$ . We could not replicate the original experiment with 1098 units because of the budget constraint.

After collecting data from the pilot experiment, we calculate the penalized stratification defined in (29), and conduct inference on the ATE in two ways: disregarding data from the pilot experiment as in (41), and combining data from the pilot and main experiments as in (43). We compare the results with the original ones in [DellaVigna and Pope \(2018\)](#). For a meaningful comparison, we also present the scaled-up version of the original standard errors in [DellaVigna and Pope \(2018\)](#) to match the sample size in our replication. Table 3 lists the sample sizes and difference-in-means estimates, standard errors, and  $t$ -statistics. Since there is only one stratum in [DellaVigna and Pope \(2018\)](#), the two-sample  $t$ -test is asymptotically exact in their setup. The columns correspond to the following:

**Pen** penalized stratification in (29) and the test statistic in (42).

**Combined** penalized stratification in (29) and the test statistic in (44).

**Original (scaled)** results in [DellaVigna and Pope \(2018\)](#), with sample size scaled down to  $2n + m$  and standard error scaled up accordingly.

**Original** results in [DellaVigna and Pope \(2018\)](#) and the two-sample  $t$ -statistic.

We see that the standard error under **Combined** is 29% smaller than that under **Original (scaled)**. Equivalently, to attain the same standard error, **Combined** requires only about half the sample size of that under the stratification in [DellaVigna and Pope \(2018\)](#).

	Pen	Combined	Original (scaled)	Original
sample size	176	220	220	1098
$\hat{\theta}_n$	644	624	-	499
s.e.	108.16	<b>92.05</b>	<b>129.95</b>	58.70
$t$ -statistic	5.95	6.78	-	8.50

Table 3: Summary statistics from [DellaVigna and Pope \(2018\)](#) and our replication.

## 8 Minimax procedure

Finally, we discuss alternative procedures without reliable pilot data. In some experiments pilot data is not available, or even if there is a pilot experiment, the units might not be drawn from the same population as the main experimental units. On the other hand, the procedure in [Theorem 3.1](#) is optimal in population, which translates into optimality with large pilots in [Theorem 5.1](#), while the penalized procedure in [\(29\)](#) is based on optimality in integrated risk in a Bayesian framework, assuming linearity and normality. It is then natural to ask about finite sample optimality without linearity and normality. To answer the question, we introduce a minimax problem. We briefly highlight the results and leave all details to [Appendix E](#). By [Lemma 3.1](#) and in particular [\(16\)](#), the conditional MSEs of  $\hat{\theta}_n$  under stratifications differ only in terms of the variance of the ex-post bias conditional on the covariates, and hence we define a constant multiple of it as the loss in [\(45\)](#). Moreover, we have

$$L(\lambda|g, X^{(n)}) = 4n^2 \text{Var}_\lambda[E[\hat{\theta}_n|X^{(n)}, D^{(n)}]|X^{(n)}] = \sum_{1 \leq s \leq S} \frac{1}{n_s - 1} \sum_{i,j \in \lambda_s, i < j} (g_i - g_j)^2. \quad (46)$$

Consider the following minimax problem to find the stratification  $\lambda$  that has the best worst-case performance in terms of the loss in [\(46\)](#), where the worst-case is among a class of functions  $\mathcal{G}$ .

$$\min_{\lambda \in \Lambda} \max_{h \in \mathcal{G}} L(\lambda|h, X^{(n)}). \quad (47)$$

Our framework requires  $\mathcal{G}$  to have a bounded polyhedron structure, in the sense made precise by [Assumption E.1](#). The assumption is satisfied by a large class of shape restrictions on  $\mathcal{G}$ , including Lipschitz continuity, monotonicity, and convexity.

Our first result shows that when  $p = 1$ , under a Lipschitz model, [\(47\)](#) is solved by matching on  $X$  directly. It reflects the intuition to match on the covariate itself when little information is available on how the covariate affects potential outcomes. For more details, see [Theorem E.1](#). Unfortunately, such a result no longer holds when  $p > 1$ . Indeed, [Example E.7](#) shows that matched-pair designs may not even be minimax-optimal. We show, however, that under [Assumption E.1](#) it is possible to reformulate [\(47\)](#) into a mixed-integer linear program. The reformulation is based on the special structure in [\(46\)](#), which enables us to rewrite [\(47\)](#) into a problem in graph theory, related to but more complicated than what is known in the literature as the clique partitioning problem. The program is computationally intensive, and therefore we consider a relaxation which replaces  $\lambda \in \Lambda$  in the minimization in [\(47\)](#) with  $\lambda \in \Lambda^{\text{pair}}$ . The resulting program, related to what

is known in the literature as the minimum-weight perfect matching problem, is computationally much easier and could be computed using modern solvers such as Gurobi. In Appendix E, we compute the solutions in a simulation study. Simulation evidence suggests that although the minimax matched-pair design is in general not minimax-optimal among all stratifications, it is often close to optimal in a sense we make precise in the appendix.

## 9 Conclusion and recommendations for empirical practice

This paper provides a framework under which a certain matched-pair design is optimal among all stratified randomization procedures. To the best of our knowledge, this is the first formal justification in the literature on the use of matched-pair designs based on optimality results. We show it is optimal to match units according to the sum of expectations of potential outcomes if treated and untreated conditional on the covariates. We then provide empirical counterparts to the optimal stratification and study their properties. In particular, we provide different procedures under large and small pilots, as well as inference procedures under each of them. From the theoretical point of view, stratifying impacts the estimation efficiency of RCTs in terms of the ex-ante MSE, i.e., before treatment status is assigned, and the ex-post bias, i.e., after treatment status is assigned. Lemma 3.1 shows that ex-post bias translates into ex-ante MSE, and hence impacts the estimation of treatment effects in an RCT. From a practical point of view, matched-pair designs weakly improve estimation and typically strictly do so, as long as the function used in matching satisfies the regularity conditions laid out in Assumption 5.4. Therefore, we recommend researchers to consider using matched-pair designs, or corresponding procedures in Appendix B, when treated fractions are identical across strata but not  $\frac{1}{2}$  and when they are in addition allowed to vary across subpopulations.

Both our theoretical and simulation results suggest that the efficiency for estimation of ATE could be improved, often notably, by incorporating information from pilot data. Therefore, we recommend researchers to perform pilot studies, on the same population as the main experiment. Based on Theorem 5.2, we recommend researchers to use flexible nonparametric estimation methods to estimate the target function in (17) when the pilot is large. When the pilot is small, researchers could still use the plug-in procedure with simple estimators such as least squares, but could also consider the penalized procedure.

## A Proof of main results

For the rest of the appendix we introduce the following definition for the convex combination of matched-pair designs.

**Definition A.1.** For  $\lambda, \lambda' \in \Lambda_n^{\text{pair}}$  and  $\delta \in [0, 1]$ , define  $\delta\lambda \oplus (1 - \delta)\lambda'$  as the randomization between  $\lambda$  and  $\lambda'$  such that  $\lambda$  is implemented with probability  $\delta$ . Define the convex hull formed by all convex combinations of any matched-pair designs as

$$\text{co}(\Lambda_n^{\text{pair}}) = \left\{ \bigoplus_{1 \leq j \leq J} \delta_j \lambda^j : \lambda^j \in \Lambda_n^{\text{pair}} \text{ and } \delta_j \geq 0 \text{ for } 1 \leq j \leq J, \sum_{1 \leq j \leq J} \delta_j = 1, 1 \leq J < \infty \right\}. \quad (48)$$

### A.1 Proof of Theorem 3.1

Define  $V(\lambda)$  as the objective in (16) multiplied by  $n^2$ . We have

$$\begin{aligned} V(\lambda) &= n^2 \text{Var}_\lambda[E[\hat{\theta}_n | X^{(n)}, D^{(n)}] | X^{(n)}] \\ &= \text{Var}_\lambda \left[ \sum_{1 \leq i \leq 2n} [D_i E[Y_i(1) | X_i] - (1 - D_i) E[Y_i(0) | X_i]] \middle| X^{(n)} \right] \\ &= \text{Var}_\lambda \left[ \sum_{1 \leq i \leq 2n} D_i (E[Y_i(0) | X_i] + E[Y_i(1) | X_i]) \middle| X^{(n)} \right] \\ &= (g^{(n)})' \text{Var}_\lambda[D^{(n)}] g^{(n)}. \end{aligned}$$

Recall from Section 2 that  $\Lambda_n^{\text{pair}}$  is the set of all matched-pair designs. For any  $\lambda = \{\{\pi(1), \pi(2)\}, \dots, \{\pi(2n-1), \pi(2n)\}\} \in \Lambda_n^{\text{pair}}$ ,

$$V(\lambda) = \frac{1}{4} \sum_{1 \leq s \leq n} (g_{\pi(2s-1)} - g_{\pi(2s)})^2. \quad (49)$$

By Lemma C.2, we have  $V(\lambda^g(X^{(n)})) \leq V(\lambda)$ .

Recall the definition of convex combinations of matched-pair designs from Definition A.1. To conclude the proof, note that by Lemma C.1, for any  $\lambda \in \Lambda$  we have

$$\lambda = \bigoplus_{1 \leq j \leq J} \delta_j \lambda^j,$$

where  $\lambda^j \in \Lambda_n^{\text{pair}}$  and  $\delta_j \geq 0$  for  $1 \leq j \leq J$ ,  $\sum_{1 \leq j \leq J} \delta_j = 1$ , and  $1 \leq J < \infty$ . Then,

$$\text{MSE}(\lambda | X^{(n)}) = \sum_{1 \leq j \leq J} \delta_j \text{MSE}(\lambda^j | X^{(n)}) \geq \min_{1 \leq j \leq J} \text{MSE}(\lambda^j | X^{(n)}) \geq \text{MSE}(\lambda^g(X^{(n)}) | X^{(n)}),$$

where the last inequality follows because  $\lambda^g(X^{(n)})$  minimizes  $\text{MSE}(\lambda | X^{(n)})$  over  $\Lambda_n^{\text{pair}}$ . The theorem therefore follows. ■

### A.2 Proof of Theorem 5.3

First, note that the assumptions in Lemma C.4 hold because of Lemma C.7 and Assumption 5.4, and that  $\hat{g}_m$  is a fixed function conditional on  $\tilde{W}^{(m)}$ . Hence, by Lemma C.4 with  $\tau = \frac{1}{2}$ , with probability one for  $\tilde{W}^{(m)}$ , as  $n \rightarrow \infty$ ,

$$\sup_{t \in \mathbf{R}} \left| Q\{\sqrt{n}(\hat{\theta}_n - \theta(Q)) \leq t | \tilde{W}^{(m)}\} - \Phi(z/s_{\hat{g}_m}) \right| \rightarrow 0, \quad (50)$$

where

$$\varsigma_{\hat{g}_m}^2 = \text{Var}[Y_i(1)] + \text{Var}[Y_i(0)] - \frac{1}{2}E[(E[Y_i(1) + Y_i(0)|\hat{g}_m(X_i), \tilde{W}^{(m)}] - E[Y_i(1) + Y_i(0)])^2]. \quad (51)$$

On the other hand, note that the assumptions in Lemma C.5 hold because of Lemma C.7 and Assumption 5.4, and that  $\hat{g}_m$  is a fixed function conditional on  $\tilde{W}^{(m)}$ . Hence, by Lemma C.5 with  $\tau = \frac{1}{2}$ , with probability one for  $\tilde{W}^{(m)}$ , for all  $\epsilon > 0$ , as  $n \rightarrow \infty$ ,

$$Q\{(|\hat{\varsigma}_n^{\hat{g}_m}|^2 - \varsigma_{\hat{g}_m}^2) > \epsilon|\tilde{W}^{(m)}\} \rightarrow 0. \quad (52)$$

The conditional convergence in (38) follows immediately from (50) and (52). Since the conditional convergence holds with probability one for  $\tilde{W}^{(m)}$ , and  $\phi_n^{\hat{g}_m}(W^{(n)}) \in [0, 1]$ , the unconditional convergence follows from the dominated convergence theorem. ■

### A.3 Proof of Theorem 5.1

The first assertion follows from Lemma C.4 with  $h = g$  and  $\tau = \frac{1}{2}$ . We now show that under  $\lambda^{\hat{g}_m}(X^{(n)})$  defined in (21),  $\sqrt{n}(\hat{\theta}_n - \theta(Q)) \xrightarrow{d} N(0, \varsigma_g^2)$  for  $\varsigma_g^2$  defined in (32) as  $m, n \rightarrow \infty$ . By repeating arguments in the proof of Lemma C.4, we write

$$\sqrt{n}(\hat{\theta}_n - \theta(Q)) = A_n - B_n + C_n - D_n,$$

where

$$\begin{aligned} A_n &= \frac{1}{\sqrt{n}} \sum_{1 \leq i \leq 2n} (Y_i(1)D_i - E[Y_i(1)D_i|g^{(n)}, D^{(n)}]) \\ B_n &= \frac{1}{\sqrt{n}} \sum_{1 \leq i \leq 2n} (Y_i(0)(1 - D_i) - E[Y_i(0)(1 - D_i)|g^{(n)}, D^{(n)}]) \\ C_n &= \frac{1}{\sqrt{n}} \sum_{1 \leq i \leq 2n} (E[(Y_i(1) + Y_i(0))D_i|g^{(n)}, D^{(n)}] - D_i E[Y_i(1) + Y_i(0)]) \\ D_n &= \frac{1}{\sqrt{n}} \sum_{1 \leq i \leq 2n} (E[Y_i(0)|g^{(n)}, D^{(n)}] - E[Y_i(0)]). \end{aligned}$$

Note that unlike in Lemma C.4, the quantities above are conditioned on  $g^{(n)}$  for  $g$  defined in (17), instead of  $\hat{g}_m^{(n)}$ . Note that by Assumptions 5.2(c), 5.3, and Lemma C.8,

$$\frac{1}{n} \sum_{1 \leq s \leq n} (g_{\pi \hat{g}_m(2s-1)} - g_{\pi \hat{g}_m(2s)})^2 \xrightarrow{P} 0. \quad (53)$$

Since Assumption 5.2(a)–(b) and (53) hold, by repeating arguments in the proof of Lemma C.4 with  $\tau = \frac{1}{2}$ , it is straightforward to establish that as  $m, n \rightarrow \infty$ ,

$$\sqrt{n}(\hat{\theta}_n - \theta(Q)) \xrightarrow{d} N(0, \varsigma_g^2), \quad (54)$$

Note that (53) is enough to derive the asymptotic representation for  $C_n$  so that we need not impose Lipschitz conditions on  $E[Y_i(d)|g(X_i)]$ . ■



## A.4 Proof of Theorem 5.2

In light of Theorem 5.1, we only need to show that  $(\hat{\zeta}_n^{\hat{g}^m})^2 \xrightarrow{P} \zeta_g^2$  as  $m, n \rightarrow \infty$ . Similar arguments as those used in Lemma C.5 go through if (53) holds and

$$\frac{1}{n} \sum_{1 \leq j \leq \lfloor \frac{n}{2} \rfloor} |g_{\pi \hat{g}^m(4j-k)} - g_{\pi \hat{g}^m(4j-l)}|^2 \xrightarrow{P} 0 \quad (55)$$

for  $k \in \{2, 3\}$  and  $l \in \{0, 1\}$ . Since (55) follows from Assumptions 5.3 by Lemma C.8, the proof is concluded.

## A.5 Proof of Theorem 5.5

To begin with, note that we need only establish that as  $m, n \rightarrow \infty$ ,

$$\sqrt{m+2n}(\hat{\theta}_n^{\text{combined}} - \theta(Q)) \xrightarrow{d} N(0, \nu \zeta_{\text{pilot}}^2 + (1-\nu)2\zeta^2), \quad (56)$$

and the rest follows from Slutsky's lemma. We prove (56) by contradiction. Suppose (56) does not hold. Then, there exists a subsequence still denoted by  $\{m, n\}$  for notational simplicity, along which as  $m, n \rightarrow \infty$ ,

$$\sup_{t \in \mathbf{R}} \left| \sqrt{m+2n}(\hat{\theta}_n^{\text{combined}} - \theta(Q)) - \Phi(z/\sqrt{\nu \zeta_{\text{pilot}}^2 + (1-\nu)2\zeta^2}) \right| \rightarrow c, \quad (57)$$

where  $c > 0$ , and

$$\frac{m}{m+2n} \rightarrow \nu \in [0, 1].$$

Now consider this subsequence. Since the two convergences in the Lemma C.8 hold in probability, there exists a further subsequence along which they hold with probability one. By repeating the proof of Theorem 5.2, we could see that along this subsequence, as  $m, n \rightarrow \infty$ , with probability one for  $\tilde{W}^{(m)}$ ,

$$\sup_{t \in \mathbf{R}} \left| Q\{\sqrt{n}(\hat{\theta}_n - \theta(Q)) \leq t | \tilde{W}^{(m)}\} - \Phi(z/\zeta_g) \right| \rightarrow 0. \quad (58)$$

Along the subsequence we construct, since  $\frac{m}{m+2n} \rightarrow \nu$ , by (58), Slutsky's lemma and Lemma C.3,

$$\sqrt{m+2n}(\hat{\theta}_n^{\text{combined}} - \theta(Q)) \xrightarrow{d} N(0, \nu \zeta_{\text{pilot}}^2 + (1-\nu)2\zeta^2),$$

which is a contradiction to (57). The theorem therefore holds. ■

## A.6 Proof of Theorem 5.4

Follows from Theorem 4.2 in Bai et al. (2019) and by repeating arguments in the proof of Lemma C.4. ■

## B Supplementary results

The next theorem shows that the infeasible optimal stratification has a similar structure to (18) when  $\tau \neq \frac{1}{2}$ .

**Theorem B.1.** *Suppose the sample size is  $kn$  for  $k \in \mathbb{Z}$  and the treatment assignment scheme satisfies  $\tau_s \equiv \tau = \frac{l}{k}$ , where  $l \in \mathbb{Z}$ ,  $0 < l < k$ , and  $k$  and  $l$  are relatively prime. Then, (5) is solved by  $\lambda^{\tau, g}$  defined in (20), where  $g_{\pi^\tau, g^\tau(1)}^\tau \leq \dots \leq g_{\pi^\tau, g^\tau(kn)}^\tau$  for  $g^\tau$  defined in (19).*

**PROOF OF THEOREM B.1.** First, note that

$$\hat{\theta}_n = \frac{1}{kn} \sum_{1 \leq i \leq kn} \left( \frac{1}{\tau} Y_i(1) D_i - \frac{1}{1-\tau} Y_i(0) (1 - D_i) \right).$$

Next,

$$\text{MSE}(\lambda | X^{(n)}) = (E_\lambda[\hat{\theta}_n | X^{(n)}] - \theta(Q))^2 + \text{Var}_\lambda[\hat{\theta}_n | X^{(n)}].$$

By repeating arguments in the proof of Lemma 3.1,

$$E_\lambda[\hat{\theta}_n | X^{(n)}] - \theta(Q) = \frac{1}{kn} \sum_{1 \leq i \leq kn} (E[Y_i(1) | X_i] - E[Y_i(0) | X_i] - \theta(Q)),$$

identical across all  $\lambda \in \Lambda_n$ , so that we need only consider conditional variances of  $\hat{\theta}$  given  $X^{(n)}$  which could be decomposed as in (11). By repeating arguments in the proof of Lemma 3.1, for any  $\lambda \in \Lambda_n$ , the first term of the right-hand side of (11) equals

$$\frac{1}{k^2 n^2} \sum_{1 \leq i \leq kn} \left( \frac{\text{Var}[Y_i(1) | X_i]}{\tau} + \frac{\text{Var}[Y_i(0) | X_i]}{1-\tau} \right),$$

again identical across all  $\lambda \in \Lambda_n$ . Therefore, we need only consider

$$\text{Var}_\lambda[E[\hat{\theta}_n | X^{(n)}, D^{(n)}] | X^{(n)}].$$

By repeating arguments in the proof of Lemma C.1, a stratum of size  $kl$  where  $l > 1$  is a convex combination of stratifications with strata only of size  $k$ . We could therefore focus on the case where each stratum is of size  $k$ . For any stratification of the form  $\lambda = \{\{\pi((s-1)k+1, \dots, \pi(sk))\} : 1 \leq s \leq n\}$ ,

$$\text{Var}_\lambda[E[\hat{\theta}_n | X^{(n)}, D^{(n)}] | X^{(n)}] \propto \sum_{1 \leq s \leq n} \sum_{(s-1)k+1 \leq j \leq sk} (g_{\pi(j)}^\tau - \bar{g}_s^\tau)^2,$$

where  $g_i^\tau$  is defined in (19) and

$$\bar{g}_s^\tau = \frac{1}{k} \sum_{(s-1)k+1 \leq j \leq sk} g_{\pi(j)}^\tau.$$

To see this, first note that units are independent across strata, so that by repeating arguments in the proof of Lemma 3.1,

$$\text{Var}_\lambda[E[\hat{\theta}_n | X^{(n)}, D^{(n)}] | X^{(n)}] \propto \sum_{1 \leq s \leq n} \text{Var}_\lambda \left[ \sum_{(s-1)k+1 \leq j \leq sk} g_{\pi(j)}^\tau D_{\pi(j)} \right].$$

Next,

$$\begin{aligned} & \text{Var}_\lambda \left[ \sum_{(s-1)k+1 \leq j \leq sk} g_{\pi(j)}^\tau D_{\pi(j)} \right] \\ &= \frac{1}{\binom{k}{l}} \sum_{(s-1)k+1 \leq j_1 < \dots < j_l \leq sk} \left( \sum_{1 \leq \iota \leq l} g_{\pi(j_\iota)}^\tau - l \bar{g}_s^\tau \right)^2 \\ &= \frac{l}{k} \sum_{(s-1)k+1 \leq j \leq sk} (g_{\pi(j)}^\tau - \bar{g}_s^\tau)^2 + \frac{1}{\binom{k}{l}} \sum_{(s-1)k+1 \leq j_1 < \dots < j_l \leq sk} \sum_{1 \leq \iota_1 \neq \iota_2 \leq l} (g_{\pi(j_{\iota_1})}^\tau - \bar{g}_s^\tau)(g_{\pi(j_{\iota_2})}^\tau - \bar{g}_s^\tau) \\ &= \frac{l}{k} \sum_{(s-1)k+1 \leq j \leq sk} (g_{\pi(j)}^\tau - \bar{g}_s^\tau)^2 + \frac{\binom{k-2}{l-2}}{\binom{k}{l}} \left[ \left( \sum_{(s-1)k+1 \leq j \leq sk} g_{\pi(j)}^\tau - k \bar{g}_s^\tau \right)^2 - \sum_{(s-1)k+1 \leq j \leq sk} (g_{\pi(j)}^\tau - \bar{g}_s^\tau)^2 \right] \end{aligned}$$

$$\propto \sum_{(s-1)k+1 \leq j \leq sk} (g_{\pi(j)}^\tau - \bar{g}_s^\tau)^2,$$

where the first equality holds by definition, the second holds by expanding the square, the third holds by accounting for cross product terms, and the fourth holds because the first term inside the square bracket on the fourth line is 0. Therefore, the problem is reduced to optimal univariate clustering of  $kn$  units on the real line where each cluster is of size  $k$ , and the conclusion follows by arguing similarly to in the proof of Lemma C.2. ■

For a measurable function  $h : \mathbf{R}^p \rightarrow \mathbf{R}$ , let  $\pi^h$  be a permutation of  $\{1, \dots, kn\}$  such that  $h_{\pi^h(1)} \leq \dots \leq h_{\pi^h(kn)}$ . Define

$$\lambda^{\tau, h}(X^{(n)}) = \{\{\pi^h((s-1)k+1), \dots, \pi^h(sk)\} : 1 \leq s \leq n\}. \quad (59)$$

Further define  $\bar{h}_s^\tau = \frac{1}{k} \sum_{(s-1)k+1 \leq j \leq sk} h_{\pi^h(j)}$ .

**Assumption B.1.**  $h$  satisfies

- (a)  $0 < E[\text{Var}[Y_i(d)|h(X_i)]]$  for  $d \in \{0, 1\}$ .
- (b)  $E[Y_i^r(d)|h(X_i) = z]$  is Lipschitz for  $r = 1, 2$  and  $d = 0, 1$ .
- (c)  $\frac{1}{n} \sum_{1 \leq s \leq n} \sum_{(s-1)k+1 \leq j \leq sk} |h_{\pi^h(j)} - \bar{h}_s^\tau|^2 \xrightarrow{P} 0$ .

The next theorem is the limiting counterpart to Theorems 3.1 and B.1. It shows that across all stratifications defined by (59) for  $h$  satisfying Assumption B.1, the asymptotic variance of  $\hat{\theta}_n$  is minimized by choosing  $h = g^\tau$  defined in (19).

**Theorem B.2.** Suppose  $h : \mathbf{R}^p \rightarrow \mathbf{R}$  be a measurable function that satisfies Assumption B.1. Then,

$$\varsigma_{\tau, g^\tau}^2 \leq \varsigma_{\tau, h}^2,$$

for  $\varsigma_{\tau, g^\tau}^2$  and  $\varsigma_{\tau, h}^2$  defined in (64) and  $g^\tau$  defined in (19). Moreover, the inequality is strict unless  $E\left[\frac{Y_i(1)}{\tau} + \frac{Y_i(0)}{1-\tau} \middle| h(X_i)\right] = g^\tau(X_i)$  with probability one under  $Q$ .

**PROOF OF THEOREM B.2.** By the definition of  $\varsigma_{\tau, h}^2$  in (64), minimizing  $\varsigma_{\tau, h}^2$  with respect to  $h$  is equivalent to maximizing

$$E\left[\left(E\left[\frac{Y_i(1)}{\tau} + \frac{Y_i(0)}{1-\tau} \middle| h(X_i)\right] - \left(\frac{E[Y_i(1)]}{\tau} + \frac{E[Y_i(0)]}{1-\tau}\right)\right)^2\right].$$

Next, note that

$$\begin{aligned} & E\left[\left(g^\tau(X_i) - \left(\frac{E[Y_i(1)]}{\tau} + \frac{E[Y_i(0)]}{1-\tau}\right)\right)^2\right] \\ &= E\left[\left(g^\tau(X_i) - E\left[\frac{Y_i(1)}{\tau} + \frac{Y_i(0)}{1-\tau} \middle| h(X_i)\right] + E\left[\frac{Y_i(1)}{\tau} + \frac{Y_i(0)}{1-\tau} \middle| h(X_i)\right] - \left(\frac{E[Y_i(1)]}{\tau} + \frac{E[Y_i(0)]}{1-\tau}\right)\right)^2\right] \\ &= E\left[\left(g^\tau(X_i) - E\left[\frac{Y_i(1)}{\tau} + \frac{Y_i(0)}{1-\tau} \middle| h(X_i)\right]\right)^2\right] + E\left[\left(E\left[\frac{Y_i(1)}{\tau} + \frac{Y_i(0)}{1-\tau} \middle| h(X_i)\right] - \left(\frac{E[Y_i(1)]}{\tau} + \frac{E[Y_i(0)]}{1-\tau}\right)\right)^2\right], \end{aligned} \quad (60)$$

$$\geq E\left[\left(E\left[\frac{Y_i(1)}{\tau} + \frac{Y_i(0)}{1-\tau} \middle| h(X_i)\right] - \left(\frac{E[Y_i(1)]}{\tau} + \frac{E[Y_i(0)]}{1-\tau}\right)\right)^2\right]. \quad (61)$$

where the last inequality is strict except unless  $E \left[ \frac{Y_i(1)}{\tau} + \frac{Y_i(0)}{1-\tau} \middle| h(X_i) \right] = g^\tau(X_i)$  with probability one under  $Q$ . To show (60), note that

$$\begin{aligned} & E \left[ \left( g^\tau(X_i) - E \left[ \frac{Y_i(1)}{\tau} + \frac{Y_i(0)}{1-\tau} \middle| h(X_i) \right] \right) \left( E \left[ \frac{Y_i(1)}{\tau} + \frac{Y_i(0)}{1-\tau} \middle| h(X_i) \right] - \left( \frac{E[Y_i(1)]}{\tau} + \frac{E[Y_i(0)]}{1-\tau} \right) \right) \right] \\ &= E \left[ E \left[ g^\tau(X_i) - E \left[ \frac{Y_i(1)}{\tau} + \frac{Y_i(0)}{1-\tau} \middle| h(X_i) \right] \middle| h(X_i) \right] \left( E \left[ \frac{Y_i(1)}{\tau} + \frac{Y_i(0)}{1-\tau} \middle| h(X_i) \right] - \left( \frac{E[Y_i(1)]}{\tau} + \frac{E[Y_i(0)]}{1-\tau} \right) \right) \right] \\ &= 0, \end{aligned}$$

where the second equality holds because

$$E[g^\tau(X_i)|h(X_i)] = E \left[ \frac{Y_i(1)}{\tau} + \frac{Y_i(0)}{1-\tau} \middle| h(X_i) \right]$$

by the law of iterated expectation. The lemma is thus proved. ■

If  $\tau_s$ 's are allowed to differ across  $s$ , then  $\hat{\theta}_n$  is generally inconsistent for  $\theta$ . In such settings researchers often use the estimator from the fully saturated regression in Bugni et al. (2019). For  $1 \leq s \leq S$  and  $d \in \{0, 1\}$ , define

$$\hat{\mu}_{n,s}(1) = \frac{1}{n_s \tau_s} \sum_{i \in \lambda_s: D_i=1} Y_i$$

and

$$\hat{\mu}_{n,s}(0) = \frac{1}{n_s(1-\tau_s)} \sum_{i \in \lambda_s: D_i=0} Y_i.$$

The estimator is

$$\hat{\theta}_n^{\text{sat}} = \sum_{1 \leq s \leq S} \frac{n_s}{n} (\hat{\mu}_{n,s}(1) - \hat{\mu}_{n,s}(0)). \quad (62)$$

Note that  $\hat{\theta}_n^{\text{sat}}$  and  $\hat{\theta}_n$  coincide whenever  $\tau_s \equiv \tau \in (0, 1)$ . See Bugni et al. (2018), Tabord-Meehan (2020), and Bugni et al. (2019) for more details. By repeating arguments used in the proof of Theorem 3.1 and Theorem B.1, we could find the stratification that minimizes  $\text{MSE}(\hat{\theta}_n^{\text{sat}} | X^{(n)})$ , which is defined as in (4) with  $\hat{\theta}_n$  replaced by  $\hat{\theta}_n^{\text{sat}}$ . The solution is as follows: we first calculate the stratification defined in (20) with  $\tau$ ,  $g$ , and  $X^{(n)}$  defined separately for each subpopulation, and then take the union of those stratifications. Moreover, the next theorem enables us to derive feasible procedures similar to (21) when treated fractions are allowed to vary across subpopulations. In particular, it reveals any plug-in estimator that satisfies the regularity conditions in Assumption B.1 leads to a procedure under which the asymptotic variance of  $\hat{\theta}_n^{\text{sat}}$  is no greater than and typically strictly less than that under procedures with each subpopulation as a stratum.

**Theorem B.3.** *Suppose the sample size is  $n$ . Define a function  $f : \mathbf{R}^p \rightarrow \{1, \dots, R\}$  where  $R \geq 1$  is an integer. Define  $N_r = \{i : f(X_i) = r\}$ ,  $X^{N_r} = (X_i : i \in N_r)$ ,  $n_r = |N_r|$ , and  $p(r) = Q\{f(X_i) = r\}$ . Define  $\lambda^{\text{large}} = \bigcup_{1 \leq r \leq R} N_r$ . For  $1 \leq r \leq R$ , let  $\tau_r$  be the treated fraction in  $N_r$ . Define functions  $h^r : \mathbf{R}^p \rightarrow \mathbf{R}$  for  $1 \leq r \leq R$ . Define  $\lambda^{\text{small}} = \bigcup_{1 \leq r \leq R} \lambda^{\tau_r, h^r}(X^{N_r})$ , where  $\lambda^{\tau_r, h^r}(X^{N_r})$  is defined in (59). Suppose  $Q$  satisfies Assumption 5.1. Then, under  $\lambda^{\text{large}}$ , for  $\hat{\theta}_n^{\text{sat}}$  defined in (62), as  $n \rightarrow \infty$ ,*

$$\sqrt{n}(\hat{\theta}_n^{\text{sat}} - \theta(Q)) \xrightarrow{d} N(0, \varsigma_{\text{large}}^2),$$

where

$$\varsigma_{\text{large}}^2 = E \left[ \frac{\text{Var}[Y_i(1)]}{\tau_{f_i}} + \frac{\text{Var}[Y_i(0)]}{1-\tau_{f_i}} - \tau_{f_i}(1-\tau_{f_i}) E \left[ \left( E \left[ \frac{Y_i(1)}{\tau_{f_i}} + \frac{Y_i(0)}{1-\tau_{f_i}} \middle| f(X_i) \right] - \left( \frac{E[Y_i(1)]}{\tau_{f_i}} + \frac{E[Y_i(0)]}{1-\tau_{f_i}} \right) \right)^2 \right] \right].$$

Suppose in addition that  $h^r, 1 \leq r \leq R$  satisfy Assumption [B.1](#), under  $Q$  restricted to  $\{x \in \mathbf{R}^p : f(x) = r\}$ . Then, under  $\lambda^{\text{small}}$ , for  $\hat{\theta}_n^{\text{sat}}$  defined in [\(62\)](#), as  $n \rightarrow \infty$ ,

$$\sqrt{n}(\hat{\theta}_n^{\text{sat}} - \theta(Q)) \xrightarrow{d} N(0, \varsigma_{\text{small}}^2),$$

where

$$\varsigma_{\text{small}}^2 = E \left[ \frac{\text{Var}[Y_i(1)]}{\tau_{f_i}} + \frac{\text{Var}[Y_i(0)]}{1 - \tau_{f_i}} - \tau_{f_i}(1 - \tau_{f_i}) E \left[ \left( E \left[ \frac{Y_i(1)}{\tau_{f_i}} + \frac{Y_i(0)}{1 - \tau_{f_i}} \middle| h^{f_i}(X_i) \right] - \left( \frac{E[Y_i(1)]}{\tau_{f_i}} + \frac{E[Y_i(0)]}{1 - \tau_{f_i}} \right) \right)^2 \right] \right].$$

In addition,  $\varsigma_{\text{small}}^2 \leq \varsigma_{\text{large}}^2$  where the inequality is strict unless

$$E \left[ \frac{Y_i(1)}{\tau_{f_i}} + \frac{Y_i(0)}{1 - \tau_{f_i}} \middle| h^{f_i}(X_i) \right] = E \left[ \frac{Y_i(1)}{\tau_{f_i}} + \frac{Y_i(0)}{1 - \tau_{f_i}} \middle| f(X_i) \right]$$

with probability one under  $Q$ . Moreover, among all choices of  $(h^r : 1 \leq r \leq R)$ ,  $\varsigma_{\text{small}}^2$  is minimized by setting  $h^r = g^{r^*}$ , where  $g^{r^*}$  is defined in [\(19\)](#).

**Remark B.1.** [Tabord-Meehan \(2020\)](#) considers stratification trees, which leads to a small number of large strata, with different treated fractions in each stratum. Using results from [Theorem B.3](#), it is straightforward to combine his procedure with procedures in this paper. The asymptotic variance of  $\hat{\theta}_n^{\text{sat}}$  under the combined procedure is no greater than and typically strictly less than that under his procedure alone. The combined procedure is as follows: First, perform the procedure in [Tabord-Meehan \(2020\)](#), which produces a finite number of strata with a target treated fraction for each stratum. Second, we view each stratum as a subpopulation and calculate the stratification in [\(59\)](#) either with a fixed function  $h$  or some plug-in estimate, with  $\tau$  equal the target treated fraction. Finally, we take the union of these stratifications. The desired property now follows from [Theorem B.3](#). ■

**PROOF OF THEOREM B.3.** The first convergence holds by [Theorem 3.1 of Bugni et al. \(2019\)](#). For the second convergence, note that

$$\begin{pmatrix} \sqrt{n_1}(\hat{\mu}_{n,1}(1) - \hat{\mu}_{n,1}(0)) \\ \vdots \\ \sqrt{n_R}(\hat{\mu}_{n,R}(1) - \hat{\mu}_{n,R}(0)) \end{pmatrix} \xrightarrow{d} N(0, \text{diag}(\varsigma_{\tau_r, h^r}^2 : 1 \leq r \leq R)).$$

Meanwhile, note that  $\frac{n_r}{n} \xrightarrow{P} p(r)$  for  $1 \leq r \leq R$ . The convergence then follows by the Slutsky's lemma. The last two results could be shown similarly to [Theorem B.2](#). ■

## C Auxiliary Lemmas

In the rest of the appendix, we use  $a \lesssim b$  to denote that there exists  $c \geq 0$  such that  $a \leq cb$ .

**Lemma C.1.** *If the treatment assignment scheme satisfies [Assumption 2.1](#), then  $\Lambda_n \subseteq \text{co}(\Lambda_n^{\text{pair}})$ .*

**PROOF OF LEMMA C.1.** We first prove that  $\lambda_0 = \{\{X_1, \dots, X_{2n}\}\}$  is a convex combination of matched-pair designs. Indeed,

$$\lambda_0 = \frac{1}{|\Lambda_n^{\text{pair}}|} \bigoplus_{\lambda \in \Lambda_n^{\text{pair}}} \lambda,$$

where

$$|\Lambda_n^{\text{pair}}| = \frac{\binom{2n}{n} n!}{2^n}.$$

Next, consider  $\lambda = \{\lambda_1, \dots, \lambda_S\}$ . Let  $\Lambda_n^{\text{pair}}(\lambda_s)$  denote the set of all matched-pair designs of units in  $\lambda_s$ . Then,

$$\lambda = \frac{1}{\prod_{1 \leq s \leq S} |\Lambda_n^{\text{pair}}(\lambda_s)|} \bigoplus_{\xi^s \in \Lambda_n^{\text{pair}}(\lambda_s): 1 \leq s \leq S} \bigcup_{1 \leq s \leq S} \xi^s,$$

and the conclusion follows. ■

**Example C.1.** Let  $n = 4$  and define

$$\begin{aligned} \lambda^0 &= \{\{1, 2, 3, 4\}\} \\ \lambda^1 &= \{\{1, 2\}, \{3, 4\}\} \\ \lambda^2 &= \{\{1, 3\}, \{2, 4\}\} \\ \lambda^3 &= \{\{1, 4\}, \{2, 3\}\}. \end{aligned}$$

We have  $\lambda^0 = \frac{1}{3}\lambda^1 \oplus \frac{1}{3}\lambda^2 \oplus \frac{1}{3}\lambda^3$ . ■

**Lemma C.2.** Suppose  $m \geq 2$ , and  $x_1, \dots, x_{2m}$  are real number such that  $x_1 \leq \dots \leq x_{2m}$ . Then, for any  $\pi \in \Pi_n$ ,

$$\sum_{k=1}^m x_{\pi(2k-1)} x_{\pi(2k)} \leq \sum_{k=1}^m x_{2k-1} x_{2k}. \quad (63)$$

**PROOF OF LEMMA C.2.** We need only consider the case where there exists  $k_1 < k_2 < k_3 < k_4$  such that at least one of  $\pi(k_1), \pi(k_2)$  is greater than at least one of  $\pi(k_3), \pi(k_4)$  because the lemma trivially holds otherwise. Suppose without loss of generality that  $\pi(k_2) < \pi(k_3) < \pi(k_4) < \pi(k_1)$ , then it is easy to verify that

$$x_{\pi(k_1)} x_{\pi(k_2)} + x_{\pi(k_3)} x_{\pi(k_4)} \leq x_{\pi(k_2)} x_{\pi(k_3)} + x_{\pi(k_1)} x_{\pi(k_4)}$$

so that by interchanging two indices we decrease the sum weakly. A finite number of those interchanges maps  $\pi$  back to the identity operator, and hence (63) holds. ■

**Lemma C.3.** Let  $X_n, Y_n, Z_n$  be random variables. Suppose  $Y_n = g(Z_n) \xrightarrow{d} Y$  as  $n \rightarrow \infty$ , where  $g : \mathbf{R} \rightarrow \mathbf{R}$  is measurable and  $X_n \xrightarrow{d} X$  conditional on  $Z_n$ , with probability one for  $Z_n$ . Furthermore, suppose the distributions of both  $X$  and  $Y$  are continuous everywhere. Then

$$(X_n, Y_n) \xrightarrow{d} (X, Y),$$

where  $X \perp\!\!\!\perp Y$ .

**PROOF OF LEMMA C.3.** Since  $X$  and  $Y$  both have continuous distribution function, we need only show for any  $x, y \in \mathbf{R}$ ,

$$P\{X_n \leq x, Y_n \leq y\} \rightarrow P\{X \leq x\}P\{Y \leq y\}.$$

To this end, note that

$$\begin{aligned} &P\{X_n \leq x, Y_n \leq y\} - P\{X \leq x\}P\{Y \leq y\} \\ &= E[E[I\{X_n \leq x\}I\{Y_n \leq y\}|Z_n]] - P\{X \leq x\}P\{Y \leq y\} \end{aligned}$$

$$\begin{aligned}
&= E[E\{I\{X_n \leq x\}|Z_n\}I\{Y_n \leq y\}] - P\{X \leq x\}P\{Y \leq y\} \\
&= E[(E\{I\{X_n \leq x\}|Z_n\} - P\{X \leq x\})I\{Y_n \leq y\}] + E[P\{X \leq x\}(I\{Y_n \leq y\} - P\{Y \leq y\})] \\
&= E[(P\{X_n \leq x|Z_n\} - P\{X \leq x\})I\{Y_n \leq y\}] + (P\{Y_n \leq y\} - P\{Y \leq y\})P\{X \leq x\}
\end{aligned}$$

For the first term on the right-hand side, note that

$$P\{X_n \leq x|Z_n\} - P\{X \leq x\} \rightarrow 0$$

with probability one for  $Z_n$ , and hence the product inside the expectation converges to 0 with probability one as well, which in turn implies the expectation converges to 0 by the dominated convergence theorem since probabilities are bounded. The second term converges to 0 because of the definition of convergence in distribution and the fact that the distribution of  $Y$  has no discontinuity. ■

**Lemma C.4.** *Suppose the sample size is  $kn$  for  $k \in \mathbb{Z}$  and the treatment assignment scheme satisfies  $\tau_s \equiv \tau = \frac{l}{k}$ , where  $l \in \mathbb{Z}$ ,  $0 < l < k$ , and they are relatively prime. Suppose  $Q$  satisfies Assumption 5.1 and  $h$  satisfies Assumption B.1. Then, under  $\lambda^{\tau, h}(X^{(n)})$  defined in (59), as  $n \rightarrow \infty$ ,*

$$\sqrt{kn}(\hat{\theta}_n - \theta(Q)) \xrightarrow{d} N(0, \varsigma_{\tau, h}^2),$$

where

$$\varsigma_{\tau, h}^2 = \frac{\text{Var}[Y_i(1)]}{\tau} + \frac{\text{Var}[Y_i(0)]}{1-\tau} - \tau(1-\tau)E\left[\left(E\left[\frac{Y_i(1)}{\tau} + \frac{Y_i(0)}{1-\tau} \middle| h(X_i)\right] - \left(\frac{E[Y_i(1)]}{\tau} + \frac{E[Y_i(0)]}{1-\tau}\right)\right)^2\right]. \quad (64)$$

**PROOF OF LEMMA C.4.** To begin with, note that

$$\sqrt{kn}(\hat{\theta}_n - \theta(Q)) = A_n - B_n + C_n - D_n,$$

where

$$\begin{aligned}
A_n &= \frac{1}{\sqrt{kn}} \sum_{1 \leq i \leq kn} \left( \frac{Y_i(1)D_i}{\tau} - E\left[\frac{Y_i(1)D_i}{\tau} \middle| h^{(n)}, D^{(n)}\right] \right) \\
B_n &= \frac{1}{\sqrt{kn}} \sum_{1 \leq i \leq kn} \left( \frac{Y_i(0)(1-D_i)}{1-\tau} - E\left[\frac{Y_i(1)D_i}{\tau} \middle| h^{(n)}, D^{(n)}\right] \right) \\
C_n &= \frac{1}{\sqrt{kn}} \sum_{1 \leq i \leq kn} \left( E\left[\frac{Y_i(1)D_i}{\tau} \middle| h^{(n)}, D^{(n)}\right] - E[Y_i(1)] \right) \\
D_n &= \frac{1}{\sqrt{kn}} \sum_{1 \leq i \leq kn} \left( E\left[\frac{Y_i(0)(1-D_i)}{1-\tau} \middle| h^{(n)}, D^{(n)}\right] - E[Y_i(0)] \right).
\end{aligned}$$

Note that, conditional on  $h^{(n)}$  and  $D^{(n)}$ ,  $A_n$  and  $B_n$  are independent and  $C_n$  and  $D_n$  are constant.

We first study the limiting behavior of  $A_n$ . Conditional on  $h^{(n)}$  and  $D^{(n)}$ , the terms in the sum are independent but not identically distributed. Therefore, we proceed to verify that the Lindeberg condition holds in probability conditional on  $h^{(n)}$  and  $D^{(n)}$ . To that end, define

$$s_n^2 = s_n^2(h^{(n)}, D^{(n)}) = \sum_{1 \leq i \leq kn} \text{Var}\left[\frac{Y_i(1)D_i}{\tau} \middle| h^{(n)}, D^{(n)}\right]$$

and note that

$$\begin{aligned}
s_n^2 &= \sum_{1 \leq i \leq kn} \text{Var} \left[ \frac{Y_i(1)D_i}{\tau} \middle| h^{(n)}, D^{(n)} \right] \\
&= \frac{1}{\tau^2} \sum_{1 \leq i \leq kn} D_i \text{Var}[Y_i(1)|h^{(n)}] \\
&= \frac{1}{\tau^2} \sum_{1 \leq i \leq kn: D_i=1} \text{Var}[Y_i(1)|h(X_i)],
\end{aligned}$$

where the second equality follows from (2) and the third follows from the fact that units are i.i.d. It follows that

$$\tau \frac{s_n^2}{kn} = \frac{1}{kn} \sum_{1 \leq i \leq kn} \text{Var}[Y_i(1)|h(X_i)] + \left( \frac{1-\tau}{\tau kn} \sum_{1 \leq i \leq kn: D_i=1} \text{Var}[Y_i(1)|h(X_i)] - \frac{1}{kn} \sum_{1 \leq i \leq kn: D_i=0} \text{Var}[Y_i(1)|h(X_i)] \right). \quad (65)$$

By Assumption 5.1,

$$\frac{1}{kn} \sum_{1 \leq i \leq kn} \text{Var}[Y_i(1)|h(X_i)] \xrightarrow{P} E[\text{Var}[Y_i(1)|h(X_i)]] < E[Y_i(1)] < \infty. \quad (66)$$

Meanwhile,

$$\begin{aligned}
&\left| \frac{1-\tau}{\tau kn} \sum_{1 \leq i \leq kn: D_i=1} \text{Var}[Y_i(1)|h(X_i)] - \frac{1}{kn} \sum_{1 \leq i \leq kn: D_i=0} \text{Var}[Y_i(1)|h(X_i)] \right| \\
&\lesssim \left| \frac{1-\tau}{\tau kn} \sum_{1 \leq i \leq kn: D_i=1} h_i - \frac{1}{kn} \sum_{1 \leq i \leq kn: D_i=0} h_i \right| \\
&= \frac{1}{\tau kn} \left| \sum_{1 \leq s \leq n} \sum_{(s-1)k+1 \leq j \leq sk: D_{\pi\tau, h(j)}=1} (h_{\pi\tau, h(j)} - \bar{h}_s^\tau) \right| \\
&\leq \frac{1}{\tau kn} \sum_{1 \leq s \leq n} \sum_{(s-1)k+1 \leq j \leq sk: D_{\pi\tau, h(j)}=1} |h_{\pi\tau, h(j)} - \bar{h}_s^\tau| \\
&\lesssim \frac{1}{n} \sum_{1 \leq s \leq n} \sum_{(s-1)k+1 \leq j \leq sk} |h_{\pi\tau, h(j)} - \bar{h}_s^\tau| \\
&\leq \left( \frac{1}{n} \sum_{1 \leq s \leq n} \sum_{(s-1)k+1 \leq j \leq sk} |h_{\pi\tau, h(j)} - \bar{h}_s^\tau|^2 \right)^{1/2} \xrightarrow{P} 0, \quad (67)
\end{aligned}$$

where the first inequality holds by Assumption B.1(b), the second holds by using Assumption B.1(c), the third holds by inspection, the last holds by the Cauchy-Schwarz inequality, and the equality holds by inspection. Combining (65), (66), and (67), we have

$$\frac{s_n^2}{kn} \xrightarrow{P} \frac{E[\text{Var}[Y_i(1)|h(X_i)]]}{\tau} > 0, \quad (68)$$

where the inequality holds by Assumption B.1(a).

We now argue that the Lindeberg condition holds in probability conditional on  $h^{(n)}$  and  $D^{(n)}$ , i.e., for any  $\epsilon > 0$ ,

$$E_n = \frac{1}{s_n^2 \tau^2} \sum_{1 \leq i \leq kn} E[|Y_i(1)D_i - E[Y_i(1)D_i|h^{(n)}, D^{(n)}]|^2 I\{|Y_i(1)D_i - E[Y_i(1)D_i|h^{(n)}, D^{(n)}]| > \epsilon \tau s_n\} | h^{(n)}, D^{(n)}] \xrightarrow{P} 0.$$

To this end, first note that for any  $M > 0$ ,

$$P\{\epsilon \tau s_n > M\} \rightarrow 1 \quad (69)$$



because of (68). Next, note that

$$E[Y_i(1)D_i|h^{(n)}, D^{(n)}] = E[Y_i(1)|h(X_i)]D_i$$

because of (2). As a result, for any  $M > 0$

$$\begin{aligned} E_n &= \frac{1}{s_n^2 \tau^2} \sum_{1 \leq i \leq kn: D_i=1} E[|Y_i(1) - E[Y_i(1)|h^{(n)}, D^{(n)}]|^2 I\{|Y_i(1) - E[Y_i(1)|h^{(n)}, D^{(n)}]| > \epsilon \tau s_n\} |h^{(n)}, D^{(n)}] \\ &\leq \frac{1}{s_n^2 \tau^2} \sum_{1 \leq i \leq kn} E[|Y_i(1) - E[Y_i(1)|h^{(n)}, D^{(n)}]|^2 I\{|Y_i(1) - E[Y_i(1)|h^{(n)}, D^{(n)}]| > \epsilon \tau s_n\} |h^{(n)}, D^{(n)}] \\ &\leq \frac{1}{s_n^2 \tau^2} \sum_{1 \leq i \leq kn} E[|Y_i(1) - E[Y_i(1)|h(X_i)]|^2 I\{|Y_i(1) - E[Y_i(1)|h(X_i)]| > M\} |h^{(n)}, D^{(n)}] + o_p(1) \\ &= \frac{kn}{s_n^2 \tau^2} \frac{1}{kn} \sum_{1 \leq i \leq kn} E[|Y_i(1) - E[Y_i(1)|h(X_i)]|^2 I\{|Y_i(1) - E[Y_i(1)|h(X_i)]| > M\} |h^{(n)}, D^{(n)}] + o_p(1) \end{aligned} \quad (70)$$

$$\xrightarrow{P} (E[\text{Var}[Y_i(1)|h(X_i)]])^{-1} E[|Y_i(1) - E[Y_i(1)|h(X_i)]|^2 I\{|Y_i(1) - E[Y_i(1)|h(X_i)]| > M\}], \quad (71)$$

where the first inequality holds by inspection, the second holds because of (69) and the equality follows because (2) and  $Q_n = Q^{kn}$ , and the convergence in probability follows from (68) and the fact that Assumption B.1(a) implies

$$\begin{aligned} E[|Y_i(1) - E[Y_i(1)|h(X_i)]|^2 I\{|Y_i(1) - E[Y_i(1)|h(X_i)]| > M\}] \\ \leq E[|Y_i(1) - E[Y_i(1)|h(X_i)]|^2] = E[\text{Var}[Y_i(1)|h(X_i)]] \leq E[Y_i^2(1)] < \infty. \end{aligned}$$

In addition, by the dominated convergence theorem,

$$\lim_{M \rightarrow \infty} E[|Y_i(1) - E[Y_i(1)|h(X_i)]|^2 I\{|Y_i(1) - E[Y_i(1)|h(X_i)]| > M\}] = 0.$$

To show that  $E_n \xrightarrow{P} 0$ , fix any subsequence, which we still call  $\{n\}$  with some abuse of notation, and we argue that there is a further subsequence along which  $E_n$  converges to 0 almost surely. Indeed, for the subsequence  $\{n\}$ , for any fixed  $M$ , the preceding display is bounded by (70), which we define as  $U_n(M)$ . We know from above that  $U_n(M) \xrightarrow{P} U(M)$ , where  $U(M)$  is defined as (71). Hence, there exists a further subsequence  $\{n\}$  along which  $U_n(M) \rightarrow U(M)$  almost surely. We then choose a sequence  $\{M_n\}_{n \geq 1}$  such that  $M_n \rightarrow \infty$ . By the dominated convergence theorem,  $\lim_{n \rightarrow \infty} U(M_n) = 0$ . By a diagonalizing argument, we could construct a further subsequence  $\{n\}$  along which  $U_n(M_n) \rightarrow 0$ . Along this subsequence, since  $E_n \leq U_n(M_n)$  for each  $n$ , the almost sure limit of  $E_n$  must be zero because it is non-negative.

We now argue that

$$\sup_{t \in \mathbf{R}} \left| P\{A_n \leq t | h^{(n)}, D^{(n)}\} - \Phi\left(t / \sqrt{E[\text{Var}[Y_i(1)|h(X_i)]] / \tau}\right) \right| \xrightarrow{P} 0.$$

Fix any subsequence. Since  $E_n \xrightarrow{P} 0$ , there exists a further subsequence along which  $E_n \rightarrow 0$  with probability one for  $h^{(n)}, D^{(n)}$ . Because of the Lindeberg condition and (68), it follows that with probability one for  $h^{(n)}, D^{(n)}$ ,  $A_n \xrightarrow{d} N(0, E[\text{Var}[Y_i(1)|h(X_i)]] / \tau)$  conditional on  $h^{(n)}, D^{(n)}$ . But then the left-hand side of the preceding display must converge almost surely to 0 by Pólya's theorem. Since for any subsequence there exists a further subsequence along which it converges to 0 almost surely, it must converge to 0 in probability.

A similar argument establishes that

$$\sup_{t \in \mathbf{R}} \left| P\{B_n \leq t | h^{(n)}, D^{(n)}\} - \Phi\left(t / \sqrt{E[\text{Var}[Y_i(0)|h(X_i)]] / (1 - \tau)}\right) \right| \xrightarrow{P} 0.$$

Since  $A_n$  and  $B_n$  are independent conditional on  $h^{(n)}$  and  $D^{(n)}$ , it follows by a similar subsequencing argument as above that

$$\sup_{t \in \mathbf{R}} \left| P\{A_n - B_n \leq t | h^{(n)}, D^{(n)}\} - \Phi\left(\frac{t}{\sqrt{E[\text{Var}[Y_i(1)|h(X_i)]/\tau + E[\text{Var}[Y_i(0)|h(X_i)]/(1-\tau)]}}\right) \right| \xrightarrow{P} 0. \quad (72)$$

To study  $C_n$ , note that by (2),

$$C_n = \frac{1}{\sqrt{kn}} \sum_{1 \leq i \leq kn} \left( E \left[ \frac{Y_i(1)}{\tau} \middle| h(X_i) \right] D_i - E[Y_i(1)] \right).$$

So we have

$$E[C_n | h^{(n)}] = \frac{1}{\sqrt{kn}} \sum_{1 \leq i \leq kn} (E[Y_i(1)|h(X_i)] - E[Y_i(1)]).$$

Furthermore, by Assumptions B.1(b)-(c),

$$\text{Var}[C_n | h^{(n)}] \propto \frac{1}{kn} \sum_{1 \leq s \leq n} (h_{\pi\tau, h(i)} - \bar{h}_\tau^s)^2 \xrightarrow{P} 0,$$

where the first relation could be established by repeating the arguments used in the last step of establishing Theorem B.1. It therefore follows by Markov's inequality that for any  $\epsilon > 0$ ,

$$P\{|C_n - E[C_n | h^{(n)}]| > \epsilon | h^{(n)}\} \xrightarrow{P} 0,$$

and since probabilities are bounded and hence uniformly integrable,

$$P\{|C_n - E[C_n | h^{(n)}]| > \epsilon\} \xrightarrow{P} 0,$$

and hence

$$C_n = \frac{1}{\sqrt{kn}} \sum_{1 \leq i \leq kn} (E[Y_i(1)|h(X_i)] - E[Y_i(1)]) + o_p(1).$$

A similar proof shows that

$$D_n = \frac{1}{\sqrt{kn}} \sum_{1 \leq i \leq kn} (E[Y_i(0)|h(X_i)] - E[Y_i(0)]) + o_p(1).$$

and therefore

$$\begin{aligned} C_n - D_n &= \frac{1}{\sqrt{kn}} \sum_{1 \leq i \leq kn} (E[Y_i(1)|h(X_i)] - E[Y_i(1)] - (E[Y_i(0)|h(X_i)] - E[Y_i(0)])) + o_p(1) \\ &\xrightarrow{d} N\left(0, E[(E[Y_i(1)|h(X_i)] - E[Y_i(1)] - (E[Y_i(0)|h(X_i)] - E[Y_i(0)]))^2]\right). \end{aligned}$$

We now show by contradiction that

$$\sup_{t \in \mathbf{R}} |P\{\sqrt{n}(\hat{\theta}_n - \theta(Q)) \leq t\} - \Phi(t/\varsigma_h)| \rightarrow 0.$$

Suppose not, then there must exist a subsequence along which the left-hand side of the above display converges to some  $\delta > 0$ . Along this subsequence, we could find a further subsequence along which the left-hand side of (72) converges to 0 with probability one for  $h^{(n)}$  and  $D^{(n)}$ , i.e.,

$$A_n - B_n \xrightarrow{d} N\left(0, \frac{E[\text{Var}[Y_i(1)|h(X_i)]]}{\tau} + \frac{E[\text{Var}[Y_i(0)|h(X_i)]]}{1-\tau}\right)$$

with probability one for  $h^{(n)}$  and  $D^{(n)}$ . Since  $C_n - D_n$  is constant for each  $h^{(n)}$  and  $D^{(n)}$ , Lemma C.3 establishes that

$$A_n - B_n + C_n - D_n \xrightarrow{d} N\left(0, \frac{E[\text{Var}[Y_i(1)|h(X_i)]]}{\tau} + \frac{E[\text{Var}[Y_i(0)|h(X_i)]]}{1-\tau} + E[(E[Y_i(1)|h(X_i)] - E[Y_i(1)] - (E[Y_i(0)|h(X_i)] - E[Y_i(0)]))^2]\right),$$

which, by Pólya's Theorem, implies a contradiction.

Finally, note that

$$\begin{aligned} & \frac{E[\text{Var}[Y_i(1)|h(X_i)]]}{\tau} + \frac{E[\text{Var}[Y_i(0)|h(X_i)]]}{1-\tau} + E[(E[Y_i(1)|h(X_i)] - E[Y_i(1)] - (E[Y_i(0)|h(X_i)] - E[Y_i(0)]))^2] \\ &= \frac{\text{Var}[Y_i(1)]}{\tau} + \frac{\text{Var}[Y_i(0)]}{1-\tau} - \frac{\text{Var}[E[Y_i(1)|h(X_i)]]}{\tau} - \frac{\text{Var}[E[Y_i(0)|h(X_i)]]}{1-\tau} + \\ & \quad E[(E[Y_i(1)|h(X_i)] - E[Y_i(1)] - (E[Y_i(0)|h(X_i)] - E[Y_i(0)]))^2] \\ &= \frac{\text{Var}[Y_i(1)]}{\tau} + \frac{\text{Var}[Y_i(0)]}{1-\tau} - \frac{1-\tau}{\tau} E[(E[Y_i(1)|h(X_i)] - E[Y_i(1)])^2] - \frac{\tau}{1-\tau} E[(E[Y_i(0)|h(X_i)] - E[Y_i(0)])^2] \\ & \quad - 2E[(E[Y_i(1)|h(X_i)] - E[Y_i(1)])(E[Y_i(0)|h(X_i)] - E[Y_i(0)])] \\ &= \frac{\text{Var}[Y_i(1)]}{\tau} + \frac{\text{Var}[Y_i(0)]}{1-\tau} - \tau(1-\tau)E\left[\left(E\left[\frac{Y_i(1)}{\tau} + \frac{Y_i(0)}{1-\tau} \middle| h(X_i)\right] - \left(\frac{E[Y_i(1)]}{\tau} + \frac{E[Y_i(0)]}{1-\tau}\right)\right)^2\right], \end{aligned}$$

and the result follows. ■

**Assumption C.1.**  $h$  satisfies

$$\frac{1}{n} \sum_{1 \leq j \leq \lfloor \frac{n}{2} \rfloor} |h_{\pi^h(4j-k)} - h_{\pi^h(4j-l)}|^2 \xrightarrow{P} 0$$

for  $k \in \{2, 3\}$  and  $l \in \{0, 1\}$ .

**Lemma C.5.** Define

$$\hat{\rho}_n = \frac{2}{n} \sum_{1 \leq j \leq \lfloor \frac{n}{2} \rfloor} (Y_{\pi^h(4j-3)} + Y_{\pi^h(4j-2)})(Y_{\pi^h(4j-1)} + Y_{\pi^h(4j)})$$

and

$$(\hat{\zeta}_n^h)^2 = \hat{\sigma}_n^2(1) + \hat{\sigma}_n^2(0) - \frac{1}{2}\hat{\rho}_n + \frac{1}{2}(\hat{\mu}_n(1) + \hat{\mu}_n(0))^2.$$

Suppose the treatment assignment scheme satisfies Assumption 2.1,  $Q$  satisfies Assumption 5.1, and  $h$  satisfies Assumptions B.1 and C.1. Then, under  $\lambda^{\frac{1}{2}, h}$  defined in (59),

$$(\hat{\zeta}_n^h)^2 \xrightarrow{P} \zeta_{\frac{1}{2}, h}^2.$$

**PROOF OF LEMMA C.5.** To begin with, note that  $\hat{\mu}_n(d) \xrightarrow{P} E[Y_i(d)]$  and  $\hat{\sigma}_n^2(d) \xrightarrow{P} \text{Var}[Y_i(d)]$  for  $d \in \{0, 1\}$ , by Lemma 6.5 in Bai et al. (2019). Next, we show that

$$E[\hat{\rho}_n | h^{(n)}] \xrightarrow{P} \rho^2. \tag{73}$$

For notational simplicity, we define  $\mu_d(h_i) = E[Y_i(d) | h(X_i) = h_i]$  for  $d \in \{0, 1\}$ . To see this, note that

$$\begin{aligned} & E[(Y_{\pi^h(4j-3)} + Y_{\pi^h(4j-2)})(Y_{\pi^h(4j-1)} + Y_{\pi^h(4j)}) | h^{(n)}] \\ &= \frac{1}{4}(\mu_1(h_{\pi^h(4j-3)}) + \mu_0(h_{\pi^h(4j-2)}))(\mu_1(h_{\pi^h(4j-1)}) + \mu_0(h_{\pi^h(4j)})) \\ & \quad + \frac{1}{4}(\mu_1(h_{\pi^h(4j-3)}) + \mu_0(h_{\pi^h(4j-2)}))(\mu_1(h_{\pi^h(4j)}) + \mu_0(h_{\pi^h(4j-1)})) \\ & \quad + \frac{1}{4}(\mu_1(h_{\pi^h(4j-2)}) + \mu_0(h_{\pi^h(4j-3)}))(\mu_1(h_{\pi^h(4j-1)}) + \mu_0(h_{\pi^h(4j)})) \end{aligned}$$

$$\begin{aligned}
& + \frac{1}{4}(\mu_1(h_{\pi^h(4j-2)}) + \mu_0(h_{\pi^h(4j-3)}))(\mu_1(h_{\pi^h(4j)}) + \mu_0(h_{\pi^h(4j-1)})) \\
& = \frac{1}{4}(g_h(h_{\pi^h(4j-3)}) + g_h(h_{\pi^h(4j-2)}))(g_h(h_{\pi^h(4j-1)}) + g_h(h_{\pi^h(4j)})) \\
& = \frac{1}{4} \sum_{k \in \{2,3\}, l \in \{0,1\}} g_h^2(h_{\pi^h(4j-k)}) + g_h^2(h_{\pi^h(4j-l)}) - (g_h(h_{\pi^h(4j-k)}) - g_h(h_{\pi^h(4j-l)}))^2.
\end{aligned}$$

As a result,

$$\begin{aligned}
E[\hat{\rho}_n | h^{(n)}] & = \sum_{1 \leq j \leq \lfloor \frac{n}{2} \rfloor} E[(Y_{\pi^h(4j-3)} + Y_{\pi^h(4j-2)})(Y_{\pi^h(4j-1)} + Y_{\pi^h(4j)}) | h^{(n)}] \\
& = \frac{1}{2n} \sum_{1 \leq i \leq 2n} g_h^2(h(X_i)) - \frac{1}{4n} \sum_{1 \leq j \leq \lfloor \frac{n}{2} \rfloor} \sum_{k \in \{2,3\}, l \in \{0,1\}} (g_h(h_{\pi^h(4j-k)}) - g_h(h_{\pi^h(4j-l)}))^2.
\end{aligned}$$

(73) then follows from Assumption B.1(b), C.1, the fact that

$$\begin{aligned}
E[g_h^2(h(X_i))] & \lesssim E[E[Y_i(1) | h(X_i)]^2] + E[E[Y_i(0) | h(X_i)]^2] \\
& \leq E[E[Y_i^2(1) | h(X_i)]] + E[E[Y_i^2(0) | h(X_i)]] = E[Y_i^2(1) + Y_i^2(0)] < \infty
\end{aligned}$$

because of Assumption 5.1, and an application of the WLLN.

It remains to show  $\hat{\rho}_n - E[\hat{\rho}_n | h^{(n)}] \xrightarrow{P} 0$ . We will prove

$$\frac{2}{n} \sum_{1 \leq j \leq \lfloor \frac{n}{2} \rfloor} (Y_{\pi^h(4j-2)} Y_{\pi^h(4j)} - E[Y_{\pi^h(4j-2)} Y_{\pi^h(4j)} | h^{(n)}]) \xrightarrow{P} 0,$$

and the others follow similarly. We will repeatedly use the following elementary inequalities for any  $a, b \in \mathbf{R}$  and  $\lambda > 0$ :

$$\begin{aligned}
|a + b| I\{|a + b| > \lambda\} & \leq 2|a| I\{|a| > \lambda/2\} + 2|b| I\{|b| > \lambda/2\} \\
|ab| I\{|ab| > \lambda\} & \leq |a|^2 I\{|a| > \sqrt{\lambda}\} + |b|^2 I\{|b| > \sqrt{\lambda}\}.
\end{aligned}$$

To begin with,

$$E[Y_{\pi^h(4j-2)} Y_{\pi^h(4j)} | h^{(n)}] = \frac{1}{2} \mu_1(h_{\pi^h(4j-2)}) \mu_0(h_{\pi^h(4j)}) + \frac{1}{2} \mu_1(h_{\pi^h(4j)}) \mu_0(h_{\pi^h(4j-2)})$$

Next, note that

$$\begin{aligned}
& \frac{2}{n} \sum_{1 \leq j \leq \lfloor \frac{n}{2} \rfloor} E[|Y_{\pi^h(4j-2)} Y_{\pi^h(4j)} - E[Y_{\pi^h(4j-2)} Y_{\pi^h(4j)} | h^{(n)}]| I\{|Y_{\pi^h(4j-2)} Y_{\pi^h(4j)} - E[Y_{\pi^h(4j-2)} Y_{\pi^h(4j)} | h^{(n)}]| > \lambda\} | h^{(n)}] \\
& \leq \frac{2}{n} \sum_{1 \leq j \leq \lfloor \frac{n}{2} \rfloor} E[|Y_{\pi^h(4j-2)} Y_{\pi^h(4j)}| I\{|Y_{\pi^h(4j-2)} Y_{\pi^h(4j)}| > \sqrt{\lambda/2}\} | h^{(n)}] \\
& \quad + E[|E[Y_{\pi^h(4j-2)} Y_{\pi^h(4j)} | h^{(n)}]| I\{|E[Y_{\pi^h(4j-2)} Y_{\pi^h(4j)} | h^{(n)}]| > \sqrt{\lambda/2}\} | h^{(n)}] \\
& \leq \frac{2}{n} \sum_{1 \leq j \leq \lfloor \frac{n}{2} \rfloor} E[Y_{\pi^h(4j-2)}^2 I\{|Y_{\pi^h(4j-2)}| > \sqrt{\lambda/2}\} | h^{(n)}] + E[Y_{\pi^h(4j)}^2 I\{|Y_{\pi^h(4j)}| > \sqrt{\lambda/2}\} | h^{(n)}] \\
& \quad + |\mu_1(h_{\pi^h(4j-2)}) \mu_0(h_{\pi^h(4j)})| I\{|\mu_1(h_{\pi^h(4j-2)}) \mu_0(h_{\pi^h(4j)})| > \lambda/2\} \\
& \quad + |\mu_1(h_{\pi^h(4j)}) \mu_0(h_{\pi^h(4j-2)})| I\{|\mu_1(h_{\pi^h(4j)}) \mu_0(h_{\pi^h(4j-2)})| > \lambda/2\} \\
& \leq \frac{2}{n} \sum_{1 \leq j \leq \lfloor \frac{n}{2} \rfloor} E[Y_{\pi^h(4j-2)}^2(1) I\{|Y_{\pi^h(4j-2)}(1)| > \sqrt{\lambda/2}\} | h^{(n)}] + E[Y_{\pi^h(4j-2)}^2(0) I\{|Y_{\pi^h(4j-2)}(0)| > \sqrt{\lambda/2}\} | h^{(n)}]
\end{aligned}$$

$$\begin{aligned}
& + E[Y_{\pi^h(4j)}^2(1)I\{|Y_{\pi^h(4j)}(1)| > \sqrt{\lambda/2}\}|h^{(n)}] + E[Y_{\pi^h(4j)}^2(0)I\{|Y_{\pi^h(4j)}(0)| > \sqrt{\lambda/2}\}|h^{(n)}] \\
& + \mu_1^2(h_{\pi^h(4j-2)})I\{|\mu_1(h_{\pi^h(4j-2)})| > \sqrt{\lambda/2}\} + \mu_0^2(h_{\pi^h(4j)})I\{|\mu_0(h_{\pi^h(4j)})| > \sqrt{\lambda/2}\} \\
& + \mu_1^2(h_{\pi^h(4j)})I\{|\mu_1(h_{\pi^h(4j)})| > \sqrt{\lambda/2}\} + \mu_0^2(h_{\pi^h(4j-2)})I\{|\mu_0(h_{\pi^h(4j-2)})| > \sqrt{\lambda/2}\} \\
\lesssim & \frac{1}{2n} \sum_{1 \leq i \leq 2n} E[Y_i^2(1)I\{|Y_i(1)| > \sqrt{\lambda/2}\}|h(X_i)] + E[Y_i^2(0)I\{|Y_i(1)| > \sqrt{\lambda/2}\}|h(X_i)] \\
& + E[Y_i^2(1)|h(X_i)]I\{E[Y_i^2(1)|h(X_i)] > \sqrt{\lambda/2}\} + E[Y_i^2(0)|h(X_i)]I\{E[Y_i^2(0)|h(X_i)] > \sqrt{\lambda/2}\} \\
\stackrel{P}{\rightarrow} & E[Y_i^2(1)I\{|Y_i(1)| > \sqrt{\lambda/2}\}] + E[Y_i^2(0)I\{|Y_i(1)| > \sqrt{\lambda/2}\}] + E[E[Y_i^2(1)|h(X_i)]I\{E[Y_i^2(1)|h(X_i)] > \sqrt{\lambda/2}\}] \\
& + E[E[Y_i^2(0)|h(X_i)]I\{E[Y_i^2(0)|h(X_i)] > \sqrt{\lambda/2}\}], \tag{74}
\end{aligned}$$

where the last line follows from WLLN and the law of iterated expectation. Since by Assumption 5.1 we have  $E[Y_i^2(d)] < \infty$  and hence  $E[E[Y_i(d)|h(X_i)]^2] < E[Y_i^2(d)]$  by Jensen's inequality, the limit as  $\lambda \rightarrow \infty$  of the last line is 0, by the dominated convergence theorem. We finish the proof by arguing by contradiction. Suppose

$$\hat{\rho}_n - E[\hat{\rho}_n|h^{(n)}]$$

does not converge in probability to 0. There must then exist  $\epsilon > 0$  and  $\delta > 0$  and a subsequence, which for simplicity we again denote by  $\{n\}$ , such that

$$P\{|\hat{\rho}_n - E[\hat{\rho}_n|h^{(n)}]| > \epsilon\} \rightarrow \delta \tag{75}$$

along this subsequence. But because of (74), there exists a further subsequence along which the condition in Lemma 6.3 of Bai et al. (2019) holds with probability one for  $h^{(n)}$ , but then along this subsequence  $\hat{\rho}_n - E[\hat{\rho}_n|h^{(n)}] \xrightarrow{P} 0$  conditional on  $h^{(n)}$  with probability one for  $h^{(n)}$ , i.e., for any  $\epsilon > 0$ , with probability one for  $h^{(n)}$ ,

$$P\{|\hat{\rho}_n - E[\hat{\rho}_n|h^{(n)}]| > \epsilon|h^{(n)}\} \rightarrow 0.$$

Since probabilities are bounded and hence uniformly integrable,

$$P\{|\hat{\rho}_n - E[\hat{\rho}_n|h^{(n)}]| > \epsilon\} \rightarrow 0$$

along the chosen subsequence, which implies a contradiction to (75). ■

**Lemma C.6.** Suppose  $U_i$ ,  $1 \leq i \leq n$  are i.i.d. random variables where  $E|U_i|^r < \infty$ . Then

$$n^{-1/r} \max_{1 \leq i \leq n} |U_i| \xrightarrow{P} 0.$$

**PROOF OF LEMMA C.6.** Note that for all  $\epsilon > 0$ ,

$$\begin{aligned}
P\left\{n^{-1/r} \max_{1 \leq i \leq n} |U_i| > \epsilon\right\} &= P\left\{\max_{1 \leq i \leq n} |U_i|^r > n\epsilon^r\right\} \\
&\leq nP\{|U_i|^r > n\epsilon^r\} \leq \frac{n}{n\epsilon^r} E[|U_i|^r I\{|U_i|^r > n\epsilon^r\}] = \frac{1}{\epsilon^r} E[|U_i|^r I\{|U_i|^r > n\epsilon^r\}] \rightarrow 0,
\end{aligned}$$

where the convergence follows because of the dominated convergence theorem and that  $E|U_i|^r < \infty$ . ■

**Lemma C.7.** Suppose  $E[h^2(X_i)] < \infty$ . Then Assumptions B.1(c) and C.1 hold.

**PROOF OF LEMMA C.7.** We prove the case where  $\tau = \frac{1}{2}$  and the results follow similarly for any  $\tau \in (0, 1)$ . Note that

$$\sum_{1 \leq s \leq n} |h_{\pi^h(2s-1)} - h_{\pi^h(2s)}|^2 \leq |h_{\pi^h(2n)} - h_{\pi^h(1)}|^2 \leq 4 \max_{1 \leq i \leq 2n} h^2(X_i),$$

where the first inequality follows from the definition of  $\pi^h$  and the second inequality follows by inspection, and therefore it follows from Lemma C.6 that

$$\frac{1}{n} \sum_{1 \leq s \leq n} |h_{\pi^h(2s-1)} - h_{\pi^h(2s)}|^2 \leq \frac{4}{n} \max_{1 \leq i \leq 2n} h^2(X_i) \xrightarrow{P} 0.$$

Assumption B.1(c) thus holds. To see Assumption C.1 holds, note that

$$\frac{1}{n} \sum_{1 \leq j \leq \lfloor \frac{n}{2} \rfloor} |h_{\pi^h(4j-k)} - h_{\pi^h(4j-l)}|^2 \lesssim \frac{1}{n} |h_{\pi^h(2n)} - h_{\pi^h(1)}|^2,$$

and the result follows similarly as above. ■

**Lemma C.8.** Suppose  $g$  satisfies Assumption 5.2(c) and  $\hat{g}_m$  satisfies Assumption 5.3. Then, as  $m, n \rightarrow \infty$ ,

$$\frac{1}{n} \sum_{1 \leq s \leq n} |g_{\pi^{\hat{g}_m}(2s-1)} - g_{\pi^{\hat{g}_m}(2s)}|^2 \xrightarrow{P} 0,$$

and

$$\frac{1}{n} \sum_{1 \leq j \leq \lfloor \frac{n}{2} \rfloor} |g_{\pi^{\hat{g}_m}(4j-k)} - g_{\pi^{\hat{g}_m}(4j-l)}|^2 \xrightarrow{P} 0$$

for  $k \in \{2, 3\}$  and  $l \in \{0, 1\}$ .

**PROOF OF LEMMA C.8.** We only prove the first conclusion as the second could be shown by similar arguments. We first show that Assumption 5.3 implies

$$\frac{1}{n} \sum_{1 \leq i \leq 2n} |\hat{g}_i - g_i|^2 \xrightarrow{P} 0. \quad (76)$$

Suppose Assumption 5.3 holds. For any  $\epsilon > 0, \delta > 0$ , there exists  $M > 0$  such that for  $m > M$ ,

$$P \left\{ \int_{\mathbf{R}^p} |\hat{g}_m(x) - g(x)|^2 Q_X(dx) > \frac{\epsilon \delta}{2} \right\} \leq \frac{\delta}{2}. \quad (77)$$

By Markov's inequality again, if

$$\int_{\mathbf{R}^p} |\hat{g}_m(x) - g(x)|^2 Q_X(dx) \leq \frac{\epsilon \delta}{2},$$

then by the independence of  $\tilde{W}^{(m)}$  and  $W^{(n)}$ ,

$$P \left\{ \frac{1}{2n} \sum_{1 \leq i \leq 2n} |\hat{g}_i - g_i|^2 > \epsilon \mid \tilde{W}^{(m)} \right\} \leq \frac{E \left[ \frac{1}{2n} \sum_{1 \leq i \leq 2n} |\hat{g}_i - g_i|^2 \mid \tilde{W}^{(m)} \right]}{\epsilon} = \frac{\int_{\mathbf{R}^p} |\hat{g}_m(x) - g(x)|^2 Q_X(dx)}{\epsilon} \leq \frac{\delta}{2}. \quad (78)$$

Then,

$$P \left\{ \frac{1}{2n} \sum_{1 \leq i \leq 2n} |\hat{g}_i - g_i|^2 > \epsilon \right\} \leq P \left\{ \frac{1}{2n} \sum_{1 \leq i \leq 2n} |\hat{g}_i - g_i|^2 > \epsilon \mid \tilde{W}^{(m)} \right\} P \left\{ \int_{\mathbf{R}^p} |\hat{g}_m(x) - g(x)|^2 Q_X(dx) \leq \frac{\epsilon \delta}{2} \right\}$$

$$\begin{aligned}
& + P \left\{ \int_{\mathbf{R}^p} |\hat{g}_m(x) - g(x)|^2 Q_X(dx) > \frac{\epsilon\delta}{2} \right\} \\
& \leq \frac{\delta}{2} \left( 1 - \frac{\delta}{2} \right) + \frac{\delta}{2} \leq \delta,
\end{aligned}$$

where the first inequality follows by definition, and the second inequality follows from (77) and (78).

Next, note that since  $|a + b|^2 \leq 2(a^2 + b^2)$  for any  $a, b \in \mathbf{R}$ ,

$$\begin{aligned}
& \frac{1}{n} \sum_{1 \leq s \leq n} |g_{\pi^{\hat{g}_m}(2s-1)} - g_{\pi^{\hat{g}_m}(2s)}|^2 \\
& \lesssim \frac{1}{n} \sum_{1 \leq s \leq n} |\hat{g}_{\pi^{\hat{g}_m}(2s-1)} - \hat{g}_{\pi^{\hat{g}_m}(2s)}|^2 + \frac{1}{n} \sum_{1 \leq i \leq 2n} |\hat{g}_i - g_i|^2.
\end{aligned} \tag{79}$$

Next, note that

$$\begin{aligned}
& \frac{1}{n} \sum_{1 \leq s \leq n} |\hat{g}_{\pi^{\hat{g}_m}(2s-1)} - \hat{g}_{\pi^{\hat{g}_m}(2s)}|^2 \\
& \leq \frac{1}{n} \max_{1 \leq i \leq 2n} |\hat{g}_i|^2 \\
& \lesssim \frac{1}{n} \max_{1 \leq i \leq 2n} |g_i|^2 + \frac{1}{n} \max_{1 \leq i \leq 2n} |\hat{g}_i - g_i|^2 \\
& \lesssim \frac{1}{n} \max_{1 \leq i \leq 2n} |g_i|^2 + \frac{1}{n} \sum_{1 \leq i \leq 2n} |\hat{g}_i - g_i|^2.
\end{aligned} \tag{80}$$

The conclusion then follows from (76), (79), (80), Assumption 5.2(c) and an application of Lemma C.6. ■

## C.1 Sufficient conditions for Lipschitz continuity

Let  $f$  denote the density function of  $X$ . Recall that  $C^{(r)}$  is the class of functions which are  $r$ th continuously differentiable. We impose the following assumption on  $h$  in Assumption B.1 and  $f$ .

**Assumption C.2.** The function  $h$  and density function  $f$  satisfy the following conditions.

- (a)  $h \in C^{(2)}$ .
- (b)  $\frac{\partial h(x)}{\partial x_p} \neq 0$  Lebesgue a.e.
- (c)  $f \in C^{(2)}$ .

**Lemma C.9** (Theorem 24.4 of Munkres (1997)). *Let  $O$  be open in  $\mathbf{R}^p$  and  $f : O \rightarrow \mathbf{R}$  be of class  $C^{(r)}$  for  $r \geq 1$ . Let  $M$  be the set of points  $x$  for which  $f(x) = 0$  and  $N$  be the set of points  $x$  for which  $f(x) \geq 0$ . Suppose  $M$  is non-empty and  $Df(x)$  has rank 1 at each point of  $M$ . Then  $N$  is a  $p$ -manifold in  $\mathbf{R}^p$  and  $\partial N = M$ .*

**Lemma C.10.** *Suppose Assumption C.2(a)-(b) hold. Then  $M = \{x : h(x) = z\}$  is a  $(p-1)$ -manifold in  $\mathbf{R}^p$ .*

**PROOF OF LEMMA C.10.** For each  $x \in M$ , we aim at providing a coordinate patch on  $M$  about  $x$ . Indeed, by Assumption C.2(a)-(b) and Theorem 9.2 (implicit function theorem) of Munkres (1997), there exists an open set  $U$  containing  $u = (x_1, \dots, x_{p-1})$ , an open ball  $B(z)$  containing  $z$  and an open set  $O$  in  $\mathbf{R}$  containing  $x_p$ , and a function  $k : U \times B(z) \rightarrow \mathbf{R}^p$  of class  $C^{(2)}$  such that  $h(u, k(u, z')) = z'$  for all  $u \in U$ ,  $z' \in B(z)$  and  $x \in O$ . Moreover,  $k(U \times B(z)) = O$ . Define the coordinate patch  $\alpha(u; z) = (u, k(u, z))$ . The conclusion follows by Theorem 5-2 of Spivak (1965). ■

Note that  $M = \{x : h(x) = z\}$  is a  $(p - 1)$ -manifold by Lemmas C.9 and C.10. In what follows, we will need the definition of the integral of a function  $g$  over the manifold  $M$ . In order to do so, note that there exists a coordinate patch as  $\{\alpha_j : U_j \subseteq \mathbf{R}^{p-1} \rightarrow V_j \subseteq M, j \in \mathcal{J}\}$ , where  $\alpha_j(u) = \alpha_j(u, z)$ , and each  $\alpha_j(u) = (u, k_j(u))$  for some function  $k_j : U \rightarrow \mathbf{R}$  which is of class  $C^2$ , as shown in the proof of Lemma C.10, and  $\alpha_j(U_j) = V_j$ . Next, there exists a partition of unity  $\{\phi_i : i \in \mathcal{I}\}$  dominated by the  $\{V_j : j \in \mathcal{J}\}$ . Moreover, both  $\mathcal{I}$  and  $\mathcal{J}$  could be chosen to be countable, according to Section 25 of Munkres (1997). The integral of a scalar function  $g$  over the manifold is written as

$$\int_M g \, dV = \sum_{j \in \mathcal{J}} \sum_{i \in \mathcal{I}} \int_{U_j} [(g\phi_i) \circ \alpha_j] V(D\alpha_j),$$

where  $V(A) = \sqrt{\det(A'A)}$  is the volume. We have

$$D\alpha_j = \begin{bmatrix} I_{p-1} & \frac{\partial k_j(u, z)}{\partial u} \end{bmatrix},$$

so that

$$V(D\alpha_j) = \sqrt{1 + \frac{\partial k_j(u, z)}{\partial u'} \frac{\partial k_j(u, z)}{\partial u}} = \frac{\|\nabla h(u, k_j(u, z))\|}{|D_p h(u, k_j(u, z))|},$$

where  $D_p = \frac{\partial}{\partial x_p}$ , by the implicit function theorem and matrix determinant lemma. Note that on one hand, for each  $j \in \mathcal{J}$ , only a finite number of  $\phi_i$  is positive, and on the other hand,  $\{\phi_i : i \in \mathcal{I}\}$  is dominated by the coordinate patch, which means that each  $\phi_i$  is supported on a compact set inside a single  $V_j$ . As a result, the order of the above double sum could be interchanged.

By p.345 of Bogachev (2007), the conditional expectation of a function  $g$  on the manifold  $M$  is defined as

$$E[g(X)|M] = \lim_{t \rightarrow 0} \frac{E[g(X)I\{z \leq h(X) \leq z + t\}]}{P\{z \leq h(X) \leq z + t\}}.$$

**Lemma C.11.** *Suppose Assumption C.2(a)–(c) hold. Then*

$$E[g(X)|M] = \frac{\int_M \frac{fg}{\|\nabla h\|} \, dV}{\int_M \frac{f}{\|\nabla h\|} \, dV}. \quad (81)$$

For a continuously differentiable function  $h : \mathbf{R}^p \rightarrow \mathbf{R}$ ,  $x \in \mathbf{R}^p$  is a critical point of  $h$  if  $\nabla h(x) = 0$ , where  $\nabla h(x)$  is the gradient of  $h$  at  $x$ ; otherwise  $x$  is a regular point of  $h$ . A value  $z$  is a critical value of  $h$  if the set  $\{x : h(x) = z\}$  contains at least one critical point; otherwise  $z$  is a regular value of  $h$ .

**PROOF OF LEMMA C.11.** By L'Hospital's rule,

$$E[g(X)|M] = \frac{\lim_{t \rightarrow 0} \frac{E[g(X)I\{z \leq h(X) \leq z + t\}]}{t}}{\lim_{t \rightarrow 0} \frac{P\{z \leq h(X) \leq z + t\}}{t}},$$

and the lemma follows from Lemma A.1 of Chernozhukov et al. (2018). In particular, the denominator equals the one in (81) directly by that lemma, while for the numerator we merely need to redefine the 'density' function as  $fg$  and the same proof goes through. ■

**Lemma C.12.** *Suppose Assumption C.2(a)–(b) hold. Let  $M = \{x : h(x) = z\}$ , where  $z$  is a regular value of  $h$  on  $\mathbf{R}^p$ . Then*



for any  $g \in C^{(2)}$ ,

$$\frac{\partial}{\partial z} \int_M g \, dV = \int_M \frac{D_p g}{D_p h} \, dV + \int_M g \frac{1}{\|\nabla h\|^2} \sum_{1 \leq i \leq p} \frac{D_i h D_{i_p} h}{D_p h} \, dV - \int_M g \frac{D_{pp} h}{D_p^2 h} \, dV. \quad (82)$$

**PROOF OF LEMMA C.12.** To begin with, note that

$$\begin{aligned} & \frac{\partial}{\partial z} \int_{U_j} [(g\phi_i) \circ \alpha_j] V(D\alpha_j) \\ &= \int_{U_j} D_p(g\phi_i) \frac{\partial k_j(u, z)}{\partial z} \frac{\|\nabla h\|}{|D_p h|} \\ & \quad + \int_{U_j} g\phi_i \frac{|D_p h|}{\|\nabla h\|} \frac{\partial k_j(u, z)}{\partial z} \frac{1}{D_p^4 h} \left( D_p^2 h \sum_{1 \leq i \leq p} D_i h D_{i_p} h - D_p h D_{pp} h \sum_{1 \leq i \leq p} D_i^2 h \right), \end{aligned} \quad (83)$$

where  $D_{ij}h = \partial_i \partial_j h$  for any function  $h \in C^{(2)}$ . We have suppressed the arguments of  $h$ , being  $(u, k_j(u, z))$ . Note that it is legitimate to pass differentiation inside the integral by the dominated convergence theorem. By the Implicit Function Theorem again,

$$\frac{\partial k_j(u, z)}{\partial z} = \frac{1}{D_p h(u, k_j(u, z))}. \quad (84)$$

By Theorem 7.17 of Rudin (1976), we know that  $\frac{\partial}{\partial z} \int_M g(x) \, dV$  is the sum over  $i \in \mathcal{I}, j \in \mathcal{J}$  of the two terms in (83). Using (84), the sum of the first term is

$$\begin{aligned} & \sum_{j \in \mathcal{J}} \sum_{i \in \mathcal{I}} \int_{U_j} (\phi_i D_p g + g D_p \phi_i) \frac{1}{D_p h} \frac{\|\nabla h\|}{|D_p h|} \\ &= \sum_j \int_{U_j} \frac{D_p g}{D_p h} V(D\alpha_j) \\ &= \int_M \frac{D_p g}{D_p h} \, dV, \end{aligned} \quad (85)$$

because  $\sum_{i \in \mathcal{I}} \phi_i = 1$  and hence  $\sum_{i \in \mathcal{I}} D_p \phi_i = D_p \sum_{i \in \mathcal{I}} \phi_i = 0$ . Again, the interchange of differentiation and sum is allowed because the sum is actually over a finite number of terms, by definition of a partition of unity. The sum of the second term is

$$\begin{aligned} & \sum_{j \in \mathcal{J}} \sum_{i \in \mathcal{I}} \int_{U_j} g\phi_i \frac{|D_p h|}{\|\nabla h\|} \frac{1}{D_p^4 h} \sum_{1 \leq i \leq p} (D_i h D_p h D_{i_p} h - D_i^2 h D_{pp} h) \\ &= \sum_{j \in \mathcal{J}} \int_{U_j} g \frac{D_p^2 h}{\|\nabla h\|^2} \frac{1}{D_p^4 h} \sum_{1 \leq i < p} (D_i h D_p h D_{i_p} h - D_i^2 h D_{pp} h) V(D\alpha) \\ &= \int_M g \frac{1}{\|\nabla h\|^2 D_p^2 h} \sum_{1 \leq i \leq p} (D_i h D_p h D_{i_p} h - D_i^2 h D_{pp} h) \, dV \\ &= \int_M g \frac{1}{\|\nabla h\|^2} \sum_{1 \leq i \leq p} \frac{D_i h D_{i_p} h}{D_p h} \, dV - \int_M g \frac{D_{pp} h}{D_p^2 h} \, dV. \end{aligned} \quad (86)$$

(82) now follows from (85) and (86). ■

**Theorem C.1.** Suppose Assumption C.2 holds. If  $z$  is a regular value of  $h$ , then

$$\frac{\partial}{\partial z} E[g(X)|M] = \frac{\int_M \frac{D_p(fg/D_p h)}{\|\nabla h\|} \, dV \int_M \frac{f}{\|\nabla h\|} \, dV - \int_M \frac{D_p(f/D_p h)}{\|\nabla h\|} \, dV \int_M \frac{fg}{\|\nabla h\|} \, dV}{\left[ \int_M \frac{f}{\|\nabla h\|} \, dV \right]^2}. \quad (87)$$

**PROOF OF THEOREM C.1.** To begin with, replace  $g$  in Lemma C.12 with  $\frac{f}{\|\nabla h\|}$ . We then have

$$\begin{aligned}
& \frac{\partial}{\partial z} \int_M \frac{f}{\|\nabla h\|} dV \\
&= \int_M \frac{\|\nabla h\| D_p f - \frac{f \sum_{1 \leq i \leq p} D_i h D_{ip} h}{\|\nabla h\|}}{\|\nabla h\|^2 D_p h} dV \\
&\quad + \int_M \frac{f}{\|\nabla h\|^3} \sum_{1 \leq i \leq p} \frac{D_i h D_{ip} h}{D_p h} dV - \int_M \frac{f D_{pp} h}{\|\nabla h\| D_p^2 h} dV \\
&= \int_M \frac{D_p f D_p h - f D_{pp} h}{\|\nabla h\| D_p^2 h} dV \\
&= \int_M \frac{D_p (f/D_p h)}{\|\nabla h\|} dV. \tag{88}
\end{aligned}$$

By the same arguments,

$$\frac{\partial}{\partial z} \int_M \frac{fg}{\|\nabla h\|} dV = \int_M \frac{D_p (fg/D_p h)}{\|\nabla h\|} dV. \tag{89}$$

(87) now follows from (88) and (89) together with the quotient rule. ■

In general, by the Law of Iterated Expectation

$$E[Y_i^r(d)|h(X) = z] = E[E[Y_i^r(d)|X]|h(X) = z].$$

Suppose  $h$  and the density function of  $X$ ,  $f(X)$  satisfy the smoothness conditions in Assumption C.2, the derivative

$$\frac{\partial}{\partial z} E[g(X)|h(X) = z]$$

is given in Theorem C.1, where  $g(x) = E[Y_i^r(d)|X = x]$  for  $r = 1, 2$  and  $d = 0, 1$ . In particular, it is equal to

$$\begin{aligned}
& E \left[ \frac{D_p g}{D_p h} + \frac{g D_p f}{f D_p h} - \frac{g D_{pp} h}{D_p^2 h} \middle| h(X) = z \right] - E \left[ \frac{D_p f}{f D_p h} - \frac{D_{pp} h}{D_p^2 h} \middle| h(X) = z \right] E \left[ g \middle| h(X) = z \right] \\
&= E \left[ \frac{D_p g}{D_p h} \middle| h(X) = z \right] + \text{Cov} \left[ \frac{D_p f}{f D_p h} - \frac{D_{pp} h}{D_p^2 h}, g \middle| h(X) = z \right]. \tag{90}
\end{aligned}$$

**Lemma C.13.** Each of the following conditions imply the boundedness of (90).

1.  $h$  is linear,  $\|D_p g\|_\infty < \infty$ ,  $\|g\|_\infty < \infty$  and  $\|D_p(\ln f)\|_\infty < \infty$ .
2.  $h$  is linear,  $\sup_{z \in \mathbf{R}} |E[D_p g|h(X) = z]| < \infty$ ,  $\sup_{z \in \mathbf{R}} |E[g^2|h(X) = z]| < \infty$  and  $\sup_{z \in \mathbf{R}} |E[D_p^2(\ln f)|h(X) = z]| < \infty$ .
3.  $h$  includes linear and interaction terms,  $\left\| \frac{D_p g}{D_p h} \right\|_\infty < \infty$ ,  $\|g\|_\infty < \infty$  and  $\left\| \frac{D_p(\ln f)}{D_p h} \right\|_\infty < \infty$ .

**PROOF OF LEMMA C.13.** Follows from inspection. ■

## D Details of penalized matching

In this section, we consider the solution to the Bayesian problem in (31) a particular example that motivates the penalized matching procedure defined by (29). For simplicity, we focus on the special case under which and  $Y_i(d) \sim N(X_i' \beta(d), \sigma^2)$

for  $d \in \{0, 1\}$ . Note that the potential outcomes are homoskedastic conditional on the covariates. Define  $\beta = \beta(1) + \beta(0)$ , and we have  $g(x) = x'\beta$ . As before, we suppose  $\tilde{W}^{(m)} = ((\tilde{Y}_j, \tilde{X}_j', \tilde{D}_j)' : 1 \leq j \leq m)$  is available from a pilot experiment. Suppose the prior on  $\beta(d)$  is  $G_d \stackrel{d}{=} N(\eta(d), \Sigma(d))$  for  $d \in \{0, 1\}$ , being independent across  $d \in \{0, 1\}$ . The prior distribution of  $\beta$  is then  $G(d\beta) \stackrel{d}{=} N(\eta(1) + \eta(0), \Sigma(1) + \Sigma(0))$ . We could show that the posterior distribution of  $\beta(d)$  conditional on  $\tilde{W}^{(m)}$  is

$$\bar{G}_d(d\beta|\tilde{W}^{(m)}) \stackrel{d}{=} N(\bar{\eta}, \bar{\Sigma}),$$

where for  $d \in \{0, 1\}$ ,

$$\begin{aligned} \bar{\eta}(d) &= \left( (\sigma^2)^{-1} \sum_{j:\tilde{D}_j=d} \tilde{X}_j \tilde{X}_j' + \Sigma^{-1}(d) \right)^{-1} \left( (\sigma^2)^{-1} \sum_{j:\tilde{D}_j=d} \tilde{X}_j \tilde{Y}_j + \Sigma^{-1}(d)\eta(d) \right) \\ \bar{\Sigma}(d) &= \left( (\sigma^2)^{-1} \sum_{j:\tilde{D}_j=d} \tilde{X}_j \tilde{X}_j' + \Sigma^{-1}(d) \right)^{-1}. \end{aligned}$$

Define  $\bar{\eta} = \bar{\eta}(1) + \bar{\eta}(0)$  and  $\bar{\Sigma} = \bar{\Sigma}(1) + \bar{\Sigma}(0)$ . The posterior distribution for  $\beta$  is

$$\bar{G}(d\beta|\tilde{W}^{(m)}) \stackrel{d}{=} (\bar{\eta}, \bar{\Sigma}),$$

since  $G_d(d\beta)$ 's are independent across  $d \in \{0, 1\}$ .

The next lemma provides the solution to the Bayesian problem in (31), where the choice set is over all measurable functions  $u : (\tilde{w}^{(m)}, x^{(n)}) \mapsto \lambda \in \Lambda_n$ .

**Lemma D.1.** *The solution to (31) maps each  $(\tilde{w}^{(m)}, x^{(n)})$  to  $\lambda = \{\{\pi(2s-1), \pi(2s)\} : 1 \leq s \leq n/2\}$ , where  $\pi$  solves*

$$\min_{\pi \in \Pi_n} \sum_{1 \leq s \leq n} \bar{d}(x_{\pi(2s-1)}, x_{\pi(2s)}),$$

where

$$\bar{d}(x_1, x_2) = (x_1' \bar{\eta} - x_2' \bar{\eta})^2 + (x_1 - x_2)' \bar{\Sigma} (x_1 - x_2). \quad (91)$$

**PROOF.** First note that by (9) and (12), (31) is equivalent to

$$\min_u \iiint L(u(\tilde{w}^{(m)}, x^{(n)})|\beta, x^{(n)}) Q_X^n(dx^{(n)}) Q_{\tilde{W}}^m(d\tilde{w}^{(m)}) G(d\beta). \quad (92)$$

Next, note that we could solve the problem pointwise for  $\tilde{w}^{(m)}$  and  $x^{(n)}$  since (92) is equivalent to

$$\min_u \bar{R}(u|\tilde{W}^{(m)}), \quad (93)$$

where

$$\bar{R}(u|\tilde{W}^{(m)}) = \int L(u(\tilde{W}^{(m)}, x^{(n)})|\beta, x^{(n)}) \bar{G}(d\beta|\tilde{W}^{(m)}).$$

To solve (93), first note that since  $\bar{R}(u|\tilde{W}^{(m)})$  is linear in  $u$ , by Lemma C.1, it is solved by a matched-pair design. Next,

$$\bar{R}(u|\tilde{W}^{(m)}) = \sum_{1 \leq s \leq n} ((x_{\pi(2s-1)}' \bar{\eta} - x_{\pi(2s)}' \bar{\eta})^2 + (x_{\pi(2s-1)} - x_{\pi(2s)})' \bar{\Sigma} (x_{\pi(2s-1)} - x_{\pi(2s)})).$$

As a result, minimizing it is equivalent to minimizing the sum of the distances defined in (91). ■

Finally, we want the prior to be irrelevant. For the purpose, suppose that  $\Sigma = cI$  where  $I$  is an identity matrix. We let the constant  $c \rightarrow \infty$ , so that the prior diverges to a diffuse (uninformative) one. Then,  $\bar{\eta}(d)$  converges to  $\hat{\beta}_m(d)$  in (22) and  $\bar{\Sigma}(d)$  converges to  $\hat{\Sigma}_m(d)$  defined in (23). Therefore, we define  $\hat{\beta}_m$  as in (24) and  $\hat{\Sigma}_m$  as in (25). The metric (91) converges to the metric defined in (30).

## E Minimax matching

This section describes the minimax procedure in detail. First note that  $L(\lambda|h, X^{(n)})$  depends on  $h$  only through  $h^{(n)}$ , and hence (47) is equivalent to

$$\min_{\lambda \in \Lambda} \max_{h^{(n)} \in G} L(\lambda|h^{(n)}), \quad (94)$$

where

$$L(\lambda|h^{(n)}) = L(\lambda|h, X^{(n)})$$

and

$$G = \{h^{(n)} : h \in \mathcal{G}, h_1 = 0\}.$$

The restriction  $h_1 = 0$  is a location normalization, since  $L(\lambda|h^{(n)})$  only depends on  $h^{(n)}$  through pairwise differences and is therefore shift-invariant. In order to solve (94) computationally, we impose the following requirement on  $G$ :

**Assumption E.1.**  $G$  is a bounded polyhedron in  $\mathbf{R}^n$ .

We now provide examples of  $G$  that satisfy Assumption E.1.

**Example E.1.** Consider the class of Lipschitz functions:

$$G = \{h^{(n)} : |h_i - h_j| \leq M \|X_i - X_j\| \text{ for } i \neq j, h_1 = 0\}. \quad (95)$$

$G$  satisfies Assumption E.1. ■

**Example E.2.** When  $p > 2$ , i.e.,  $X_i$  is multivariate, consider the class of functions which are Lipschitz along each dimension:

$$G = \left\{ h^{(n)} : |h_i - h_j| \leq \sum_{1 \leq l \leq p} M_l |X_{il} - X_{jl}| \text{ for } i \neq j, h_1 = 0 \right\}.$$

$G$  satisfies Assumption E.1. ■

**Example E.3.** Consider the class of functions Lipschitz in a known index. For a known function  $w$ , define

$$G = \left\{ h^{(n)} : |h_i - h_j| \leq M |\nu(X_i) - \nu(X_j)| \text{ for } i \neq j, h_1 = 0 \right\}. \quad (96)$$

$G$  satisfies Assumption E.1. ■

**Example E.4.** Consider the class of linear functions with coefficients in a bounded polyhedron. For a bounded polyhedron  $\mathcal{B}$  in  $\mathbf{R}^p$ , define

$$G = \{X^{(n)}\beta - X_1'\beta\mathbf{1}_n : \beta \in \mathcal{B}\}.$$

$G$  satisfies Assumption E.1. ■

**Example E.5.** Consider the class of monotonically increasing functions. Without loss of generality assume that  $X_1 \leq \dots \leq X_n$ . For  $M > 0$ , define

$$G = \{h^{(n)} : h_i \leq h_j \text{ for } i < j, h_n \leq M, h_1 = 0\}.$$

Since  $G$  is bounded and defined by linear inequalities, it satisfies Assumption E.1. ■

**Example E.6.** Consider the class of convex functions. Without loss of generality assume that  $X_1 \leq \dots \leq X_n$ . For  $M > 0$ , define

$$G = \left\{ h^{(n)} : h_i \leq \frac{X_{i+1} - X_i}{X_{i+1} - X_{i-1}} h_{i-1} + \frac{X_i - X_{i-1}}{X_{i+1} - X_{i-1}} h_{i+1}, 2 \leq i \leq 2n-1, |h_n| \leq M, h_1 = 0 \right\}.$$

Since  $G$  is bounded and defined by linear inequalities, it satisfies Assumption E.1. ■

Consider the minimax problem (94) with  $G$  defined in (96). The following theorem shows that without any information of how the covariate affects potential outcomes beyond the index, the best we could do is to match on the index itself.

**Theorem E.1.** *The solution to (94) with  $G$  defined in (96) is  $\lambda^\nu = \{\{\pi^\nu(2s-1), \pi^\nu(2s)\} : 1 \leq s \leq n\}$  where  $\nu_{\pi^\nu(1)} \leq \dots \leq \nu_{\pi^\nu(2n)}$ .*

**PROOF OF THEOREM E.1.** Without loss of generality, consider  $p = 1$  and  $\nu(x) = x$ . The general case is proved in exactly the same way. We use another expression of (46). Define  $\Delta_i = g_{\pi(i+1)} - g_{\pi(i)}$  for  $i = 1, \dots, 2n-1$ . For  $\lambda^0 = \{\{1, \dots, 2n\}\}$ ,

$$\begin{aligned} L(\lambda^0 | g, X^{(n)}) &= \frac{1}{2n(2n-1)} \sum_{1 \leq i \leq 2n} \left[ (2n-1)g_i - \sum_{j \neq i} g_j \right]^2 \\ &= \frac{1}{2n(2n-1)} \sum_{1 \leq i \leq 2n} \left[ - \sum_{1 \leq j \leq i-1} j \Delta_j + \sum_{i \leq j \leq 2n-1} (2n-j) \Delta_j \right]^2 \\ &= \frac{1}{2n(2n-1)} \left[ \sum_{1 \leq i \leq 2n-1} 2n(2n-i)i \Delta_i^2 + 2 \sum_{k < l \leq 2n-1} 2n(2n-l)k \Delta_k \Delta_l \right] \\ &= \frac{1}{2n-1} \left[ \sum_{1 \leq i \leq 2n-1} (2n-i)i \Delta_i^2 + 2 \sum_{k < l \leq 2n-1} (2n-l)k \Delta_k \Delta_l \right]. \end{aligned}$$

As a result, for a general stratification  $\lambda$ , the loss function (46) equals

$$L(\lambda | g, X^{(n)}) = \sum_{1 \leq s \leq S} \frac{1}{n_s - 1} \left[ \sum_{1 \leq i \leq n_s - 1} (n_s - i)i \Delta_{i,s}^2 + 2 \sum_{k < l \leq n_s - 1} (n_s - l)k \Delta_{k,s} \Delta_{l,s} \right]. \quad (97)$$

Note that  $g^{\text{mm}}(x) = Mx$  simultaneously maximizes (97) for every  $\lambda$ . But we know the stratification that solves

$$\min_{\lambda \in \Lambda} L(\lambda | g^{\text{mm}}, X^{(n)})$$

is the “optimal non-bipartite matching” of  $X$  on  $\mathbf{R}$ , i.e.  $\lambda^x$ . ■

For a prespecified  $\theta_0 \in \mathbf{R}$ , consider the problem of testing (33) at level  $\alpha \in (0, 1)$ . We use the test in (36) by setting  $\hat{g}_m = \nu$ .

**Corollary E.1.** *Suppose the treatment assignment scheme satisfies Assumption 2.1 and  $Q$  satisfies Assumption 5.1 and  $h = \nu$  satisfies Assumption B.1 with  $\tau = \frac{1}{2}$ . Then, for the problem of testing (33) at level  $\alpha \in (0, 1)$ ,  $\phi_n^\nu$  satisfies*

$$\lim_{n \rightarrow \infty} E[\phi_n^\nu(W^{(n)})] = \alpha,$$

whenever  $Q$  additionally satisfies the null hypothesis, i.e.,  $\theta(Q) = \theta_0$ .

For other specifications of  $G$  in (94), there does not exist a clean result as Theorem E.1, as illustrated by the following example.

**Example E.7.** Let  $n = 4$  and  $X_1 = (0, 0)'$ ,  $X_2 = (1, 0)'$ ,  $X_3 = (0, 1)'$ ,  $X_4 = (1, 1)'$ . Let  $n = 4$  and define

$$\begin{aligned}\lambda^0 &= \{\{1, 2, 3, 4\}\} \\ \lambda^1 &= \{\{1, 2\}, \{3, 4\}\} \\ \lambda^2 &= \{\{1, 3\}, \{2, 4\}\} \\ \lambda^3 &= \{\{1, 4\}, \{2, 3\}\}.\end{aligned}$$

Let  $G$  be as defined in (95) with  $M = 1$ . Then  $\lambda^0$  solves (94). Indeed, for  $\lambda = \lambda^1$ , the worst case occurs at  $h^{(n)} = (0, \sqrt{2} - 1, \sqrt{2} - 1, \sqrt{2})$ , with the loss equal to 2. For  $\lambda = \lambda^2$  or  $\lambda^3$ , the worst case occurs at  $h^{(n)} = (0, 1, 1, 0)$ , with the loss equal to 2. In contrast, the worst case for  $\lambda = \lambda^0$  occurs at  $h^{(n)} = (0, \sqrt{2} - 1, \sqrt{2} - 1, \sqrt{2})$ , and the loss is  $(10 - 4\sqrt{2})/3 < 2$ . ■

The key reason why (94) is hard to solve when  $p > 1$  is that the choice set  $\Lambda$  is not convex. In principle, we could convexify the problem by considering the  $\text{co}(\Lambda)$ , the convex hull of  $\Lambda$ . That amounts to allowing for mixing over (potentially a large number of) matched-pair designs, which is hard to interpret and is almost never used in practice. Although  $\Lambda$  is not convex, we can still provide computational strategies to solve (94). Note that  $L(\lambda|h^{(n)})$  is convex in  $h^{(n)}$ , which combined with Assumption E.1 implies that the inner maximum in (94) is attained on the vertices of  $G$ , which we denote by  $V$ . Then, the minimax problem is equivalent to

$$\min_{\lambda \in \Lambda} \max_{h^{(n)} \in V} L(\lambda|h^{(n)}). \quad (98)$$

We now apply results from graph theory to reformulate (98) into Mixed Integer Linear Programs (MILPs). We first recall some definitions from the graph theory and connect them to the optimal stratification problem. For more details, see [Bertsimas and Tsitsiklis \(1997\)](#).

An undirected graph  $\Gamma = (N, E)$  consists of a set of nodes  $N$  and a set of edges  $E$ . Each element of  $E$  is an unordered pair  $\{i, j\}$  where  $i \in N$  and  $j \in N$ . Define  $q_e = 1$  if  $e \in E$  and define  $\mathbf{q} = (q_e)_{e \in E}$ . Define  $q_{ij} = q_{i,j}$ . The degree of  $i$  is defined as  $d_i = \sum_j q_{ij}$ . The graph  $\Gamma$  is complete if  $q_{ij} = 1$  for all  $i \neq j$ . A subset  $U$  of  $N$  is a clique in  $\Gamma$  if  $\{i, j\} \in E$  for all  $i, j \in U$ . The set of induced edges by  $U$  is  $E(U) = \{\{i, j\} \in E : i, j \in U, i \neq j\}$ . A clique partition of  $\Gamma$  is  $\Gamma^C = (N, E(U_1, \dots, U_S))$  for  $E(U_1, \dots, U_S) = \cup_{s=1}^S E(U_s)$  where each  $U_s$  is a clique in  $\Gamma^C$  (and  $\Gamma$ ), and  $\{U_s\}_{s=1}^S$  is a partition of  $N$ , i.e.,  $N = \cup_{s=1}^S U_s$  and  $U_s \cap U_t = \emptyset$  for  $s \neq t$ .

In terms of stratification, a unit is a node and an edge  $\{i, j\} \in E$  if units  $i$  and  $j$  are in the same stratum. A stratum is a clique. A stratification  $\lambda = \{\lambda_s\}_{s=1}^S$  of  $N = \{1, \dots, n\}$  induces a clique partition  $\Gamma^\lambda = (N, E(\lambda_1, \dots, \lambda_S))$  of  $\Gamma = (N, E)$  for  $E = \{\{i, j\} : i, j \in N, i \neq j\}$  where the size of each clique  $\lambda_s$  is even, or equivalently the degree of each node in  $\Gamma^\lambda$  is odd.

Define  $c_e = (h_i - h_j)^2$  as the cost of edge  $e = \{i, j\} \in E$ ,  $\mathbf{c} = (c_e)_{e \in E}$  and  $C = \{\mathbf{c} : h^{(n)} \in V\}$ . By (46),

$$L(\lambda|h^{(n)}) = L(\lambda|h, X^{(n)}) = \sum_{1 \leq s \leq S} \frac{1}{n_s - 1} \sum_{i, j \in \lambda_s, i < j} (h_i - h_j)^2.$$

If  $n_s \equiv 2$ , then it equals

$$\sum_{e \in E} c_e q_e.$$

If  $n_s > 2$  for some  $s$ , then we need to introduce other binary variables to indicate  $n_s$ . The minimax problem (94) is equivalent to the following MILP which solves the cost minimization problem over size-bounded stratifications within  $\Lambda$ , i.e.,  $\lambda$  with  $n_s \leq 2K$  for all  $s$ .

$$\begin{aligned} \min_{\mathbf{q}} \quad & z & (99) \\ \text{subject to} \quad & \sum_{e \in E} c_e \left( \sum_{1 \leq k \leq K} \frac{u_{ik}}{2k-1} \right) I\{i \in e\} \leq z, \text{ for all } \mathbf{c} \in C, \\ & \sum_{i \in N} q_{il} = \sum_{1 \leq k \leq K} (2k-1)u_{ik}, \text{ for all } i \in N, \\ & u_{ik} \in \{0, 1\}, \text{ for all } i \in N, 1 \leq k \leq K, \\ & q_{e_1} + q_{e_2} - q_{e_3} \leq 1, \text{ for all } e_1, e_2, e_3 \in E, \\ & q_e \in \{0, 1\}, \text{ for all } e \in E. \end{aligned} \quad (100)$$

We impose an upper bound on the size of each stratum,  $2k$ .  $u_{ik}, k = 1, \dots, K-1$  are binary indicators of whether the stratum of unit  $i$  has size  $2k$ . The first set of constraints express the loss function (46). The second set of constraints say the degree of each node is  $2k-1$ , the stratum size minus one. The third set of constraints restrict  $u_{ik}$  to be binary. The fourth and the most important set of constraints, (100), are called triangle inequalities in the clique partition literature. See Grötschel and Wakabayashi (1990). They ensure that the solution to (99) is indeed a clique partition, i.e., a stratification. However, our problem differs from the standard clique partition problem in two ways: we only allow an even number of units within each clique; and the final weights on each edge in the total cost depends on the degrees of either of its nodes, rather than being a constant.

The program (99) is computationally intensive even when  $k = 2$  and becomes prohibitive quickly as  $n$  increases. Therefore, we consider two relaxations of it. The first relaxation is to optimize over  $\Lambda^p$  instead of  $\Lambda$ . For a matched-pair design  $\lambda = \{\{\pi(2s-1), \pi(2s)\} : 1 \leq s \leq n\}$ ,

$$L(\lambda|h, X^{(n)}) = \sum_{1 \leq s \leq n} (h_{\pi(2s-1)} - h_{\pi(2s)})^2.$$

As a result, we introduce the program as

$$\begin{aligned} \min_{\mathbf{q}} \quad & z & (101) \\ \text{subject to} \quad & \sum_{e \in E} c_e q_e \leq z, \text{ for all } \mathbf{c} \in C, \\ & \sum_{j \in N} q_{ij} = 1, \text{ for all } i \in N, \\ & q_e \in \{0, 1\}, \text{ for all } e \in E. \end{aligned}$$

The solution to (101) is  $\lambda^{\text{mm}} = \{e \in E : q_e = 1\}$ . We define the permutation  $\pi^{\text{mm}}$  such that  $\lambda^{\text{mm}} = \{\{\pi^{\text{mm}}(2s - 1), \pi^{\text{mm}}(2s)\} : 1 \leq s \leq n\}$ . (101) is feasible even when  $n$  is large and requires substantially less computational budget than (99). Moreover, as simulation evidence in Section in Table 4 shows, the solution to (101) is frequently the same with (99) for a small  $K$  and (102).

The second relaxation is the following hierarchical procedure.

**Algorithm E.1.**

1. Solve (101). Denote the solution by  $\mathbf{q}^0$  and denote  $\Lambda^0 = \{e \in E : q_e = 1\}$ .

2. For  $k \geq 0$ , repeat steps (a) and (b) below.

(a) For  $\mathbf{q}^k = (q_{AB}^k)_{A,B \in \Lambda^k, A \neq B}$ , solve

$$\begin{aligned} & \min_{\mathbf{q}^k} z \\ & \text{subject to} \quad \sum_{A,B \in \Lambda^k} q_{AB} c_{AB} + \sum_{A \in \Lambda^k} c_A \leq z, \text{ for all } \mathbf{c} \in C, \\ & \quad \sum_{B \in \Lambda^k} q_{AB}^k \leq 1, \quad \text{for all } A \in \Lambda^k, \\ & \quad q_{AB}^k \in \{0, 1\}, \text{ for all } A, B \in \Lambda^k, \end{aligned} \tag{102}$$

where  $c_A = L(\lambda|g, X_A)$ , for  $X_A = \{X_i : i \in A\}$  and  $c_{AB} = c_{A \cup B} - c_A - c_B$ .

(b) Update

$$\Lambda^{k+1} = \{A \cup B : q_{AB}^k = 1\} \cup \{A : \sum_{B \in \Lambda^k} q_{AB}^k = 0\}$$

until  $\Lambda^{k^*} = \Lambda^{k^*+1}$ . Collect  $\Lambda^{k^*}$  as the solution.

Algorithm E.1 iteratively decides whether to merge pairs of strata or not. The algorithm stops when no pairwise merging of existing strata reduces the worst-case loss.

We now study the properties of minimax matching in a small simulation study. We compare both the actual and worst-case losses under different stratifications. In the following model, we construct a bounded polyhedron  $G$  around  $g^{(n)}$ . We then calculate both the actual losses  $L(\lambda|g^{(n)})$  and worst-case losses  $\max_{h^{(n)} \in G} L(\lambda|h^{(n)})$  across different stratifications. We set  $g(x) = x' \beta$  and

$$G = \{X^{(n)} \beta : \beta \in \mathcal{B}\},$$

where  $\mathcal{B}$  is a polyhedron such that  $\beta \in \mathcal{B}$ .

**Model MM**  $2n = 24; p = 2; X_{i,1} = 0$  for  $1 \leq i \leq 8, X_{i,2} = 1$  for  $9 \leq i \leq 24; X_{i,2} \sim N(0, 1)^2$  i.i.d. across  $i; g(x) = x' \beta, \beta = (1, 1)'; \mathcal{B} = \mathcal{B}_1 \times \mathcal{B}_2, \mathcal{B}_1 = \beta_1 + \gamma_1 \times [0, 1], \mathcal{B}_2 = \beta_2 + \gamma_2 \times [-1, 1]; \gamma \in \{(0.5, 0.5)', (2, 2)', (0, 2)', (2, 0)'\}$ .

We randomly generate  $X^{(n)}$  in 100 replications and summarize

- (a) ratios of the values of the actual loss against those under infeasible optimal stratifications.
- (b) ratios of the values of the worst-case loss against those under size-bounded minimax stratifications with  $k = 2$ .

We consider the following stratifications:



**Oracle** infeasible optimal stratification in (18).

**by1** :  $\lambda_1 = \{i : 1 \leq i \leq 8\}, \lambda_2 = \{i : 9 \leq i \leq 24\}$ .

**by2** two strata separated by the sample median of  $X_{i,2}$ .

**2by2** four strata as the cross product of **by1** and **by2**.

**MP2** matching on  $X_{i,2}$  only, i.e., stratification in (21) with  $\hat{g}_m(x) = x_1$ .

**MPcell** within each value of  $X_{i,1}$ , optimal matched-pair design using  $X_{i,2}$ .

**MMpair** the minimax matched-pair design in (101).

**MMbdd** the size-bounded minimax stratification in (99) with  $k = 2$ .

**MMhier** results from the hierarchical procedure in Algorithm E.1.

In Model MM, **MMpair** and **MMhier** have the same solution with **MMbdd** (which we know weakly dominates **MMpair**) most of the time, while other stratifications which do not incorporate minimax consideration sometimes generate much larger worst-case losses.

		<b>Oracle</b>	<b>by1</b>	<b>by2</b>	<b>2by2</b>	<b>MP2</b>	<b>MPcell</b>	<b>MMpair</b>	<b>MMbdd</b>	<b>MMhier</b>
$\gamma = (0.5, 0.5)$	Actual	25%	<b>1.0000</b>	4.0109	2.3318	1.8158	1.0619	1.0256	1.0000	1.0000
		50%	<b>1.0000</b>	7.2394	3.8858	2.9807	1.2631	1.5291	1.0000	1.0000
		75%	<b>1.0000</b>	13.7890	7.7959	6.5012	1.8567	4.4629	1.0001	1.0001
		Mean	<b>1.0000</b>	13.6242	7.0691	7.6378	1.7480	5.5226	1.0346	1.0346
	Worst-case	25%	1.0000	4.0109	2.3381	1.7832	1.0481	1.0243	1.0000	<b>1.0000</b>
		50%	1.0000	6.9420	3.6908	2.8469	1.1858	1.4011	1.0000	<b>1.0000</b>
		75%	1.0003	11.9445	6.7125	5.9020	1.6146	3.9388	1.0000	<b>1.0000</b>
		Mean	1.0212	10.3183	5.4894	5.6169	1.4884	3.8240	1.0000	<b>1.0000</b>
$\gamma = (2, 2)$	Actual	25%	<b>1.0000</b>	4.4595	2.7007	2.1185	1.0700	1.0994	1.0000	1.0000
		50%	<b>1.0000</b>	9.2109	4.1580	3.5446	1.3348	1.8127	1.0096	1.0096
		75%	<b>1.0000</b>	14.5268	6.8864	6.9304	1.8268	4.0986	1.3038	1.3038
		Mean	<b>1.0000</b>	13.5257	7.3036	6.3795	1.7873	3.8773	1.2997	1.2997
	Worst-case	25%	1.0000	3.8897	2.2736	1.8008	1.0408	1.0315	1.0000	<b>1.0000</b>
		50%	1.0126	6.2500	3.1816	2.6923	1.1604	1.4000	1.0000	<b>1.0000</b>
		75%	1.2516	10.1048	5.0563	4.3542	1.6279	2.8357	1.0000	<b>1.0000</b>
		Mean	1.2390	8.5778	4.9668	3.9661	1.4436	2.2735	1.0000	<b>1.0000</b>
$\gamma = (0, 1)$	Actual	25%	<b>1.0000</b>	4.1720	2.5479	1.8497	1.0397	1.1857	1.0000	1.0000
		50%	<b>1.0000</b>	7.4458	3.8469	3.3647	1.2599	1.7892	1.0135	1.0135
		75%	<b>1.0000</b>	14.1891	7.6734	6.3794	1.7666	3.1199	1.1199	1.1199
		Mean	<b>1.0000</b>	12.4138	6.8864	5.5793	1.8987	2.8784	1.1301	1.1301
	Worst-case	25%	1.0000	4.3021	2.3348	1.8989	1.0012	1.2292	1.0000	<b>1.0000</b>
		50%	1.0077	7.2928	3.4658	3.6051	1.0450	1.5861	1.0000	<b>1.0000</b>
		75%	1.1138	16.6540	6.7290	6.8655	1.2165	3.7622	1.0000	<b>1.0000</b>
		Mean	1.1128	12.0228	5.8405	5.4142	1.2350	2.8276	1.0000	<b>1.0000</b>
$\gamma = (1, 0)$	Actual	25%	<b>1.0000</b>	3.5310	2.1679	2.0152	1.0654	1.0985	1.0000	1.0000
		50%	<b>1.0000</b>	8.5908	4.1682	3.8322	1.2567	1.9700	1.0481	1.0481
		75%	<b>1.0000</b>	17.9252	8.6984	8.2448	1.8296	3.6598	1.5850	1.5850
		Mean	<b>1.0000</b>	14.7115	8.3951	6.6366	1.6191	3.8705	1.7197	1.7197
	Worst-case	25%	1.0000	2.9528	2.4142	1.5418	1.1470	1.0000	1.0000	<b>1.0000</b>
		50%	1.0435	4.6975	3.3215	2.1056	1.5634	1.0211	1.0000	<b>1.0000</b>
		75%	1.6225	9.0650	5.4879	3.9089	2.6225	1.6384	1.0000	<b>1.0000</b>
		Mean	1.6231	7.8535	6.4319	3.6442	2.3219	1.8804	1.0000	<b>1.0000</b>

Table 4: Ratios of values of the actual loss under all stratifications against those under the infeasible optimal stratifications (**Oracle**) and ratios of values of the worst-case loss under all stratifications against those under size-bounded minimax stratifications (**MMbdd**) in Model MM. Benchmarks are displayed in bold face.

## **F AEA RCT Registry**

The following experiments in the AEA RCT Registry use matched-pair designs: AEARCTR-0000086, 0000171, 0000293, 0000443, 0000481, 0000550, 0000578, 0000587, 0000644, 0000688, 0000721, 0000983, 0000986, 0001034, 0001097, 0001218, 0001370, 0001591, 0001607, 0001712, 0001714, 0001778, 0001992, 0001995, 0002010, 0002125, 0002132, 0002282, 0002585, 0002622, 0002664, 0002750, 0002776, 0003056, 0003076, 0003524, 0003581, 0003629, 0003648, 0003779, 0003814, 0003933, 0003994, 0004024, 0004042, 0004022.

## References

- ABADIE, A. and IMBENS, G. W. (2008). Estimation of the Conditional Variance in Paired Experiments. *Annales d'Économie et de Statistique* 175–187.
- AKER, J. C., KSOLL, C. and LYBBERT, T. J. (2012). Can mobile phones improve learning? Evidence from a field experiment in Niger. *American Economic Journal: Applied Economics*, **4** 94–120.
- ALATAS, V., BANERJEE, A., HANNA, R., OLKEN, B. A. and TOBIAS, J. (2012). Targeting the poor: Evidence from a field experiment in Indonesia. *American Economic Review*, **102** 1206–40.
- ANGRIST, J. and LAVY, V. (2009). The effects of high stakes high school achievement awards: Evidence from a randomized trial. *American Economic Review*, **99** 1384–1414.
- ARMITAGE, P., BERRY, G. and MATTHEWS, J. N. S. (2008). *Statistical methods in medical research*. John Wiley & Sons.
- ASHRAF, N., BERRY, J. and SHAPIRO, J. M. (2010). Can higher prices stimulate product use? Evidence from a field experiment in Zambia. *American Economic Review*, **100** 2383–2413.
- ASHRAF, N., KARLAN, D. and YIN, W. (2006). Deposit collectors. *Advances in Economic Analysis & Policy*, **5**.
- ATHEY, S. and IMBENS, G. W. (2017). The econometrics of randomized experiments. In *Handbook of Economic Field Experiments*, vol. 1. Elsevier, 73–140.
- BAI, Y. (2019). Randomization under permutation invariance. Working paper.
- BAI, Y., SHAIKH, A. and ROMANO, J. P. (2019). Inference in experiments with matched pairs. Working paper.
- BAILEY, R. A. (2004). *Association schemes: Designed experiments, algebra and combinatorics*, vol. 84. Cambridge University Press.
- BANERJEE, A., CHASSANG, S., MONTERO, S. and SNOWBERG, E. (2019). A theory of experimenters.
- BANERJEE, A., DUFLO, E., GLENNERSTER, R. and KINNAN, C. (2015). The miracle of microfinance? Evidence from a randomized evaluation. *American Economic Journal: Applied Economics*, **7** 22–53.
- BARRIOS, T. (2013). Optimal stratification in randomized experiments. Working paper.
- BELLEÇ, P. C., DALALYAN, A. S., GRAPPIN, E., PARIS, Q. and OTHERS (2018). On the prediction loss of the lasso in the partially labeled setting. *Electronic Journal of Statistics*, **12** 3443–3472.
- BELLONI, A., CHEN, D., CHERNOZHUKOV, V. and HANSEN, C. (2012). Sparse models and methods for optimal instruments with an application to eminent domain. *Econometrica*, **80** 2369–2429.
- BELLONI, A., CHERNOZHUKOV, V. and HANSEN, C. (2014). Inference on treatment effects after selection among high-dimensional controls. *The Review of Economic Studies*, **81** 608–650.

- BERRY, J., KARLAN, D. and PRADHAN, M. (2018). The impact of financial education for youth in Ghana. *World Development*, **102** 71–89.
- BERTRAND, M. and DUFLO, E. (2017). Field experiments on discrimination. In *Handbook of Economic Field Experiments*, vol. 1. Elsevier, 309–393.
- BERTSIMAS, D. and TSITSIKLIS, J. N. (1997). *Introduction to linear optimization*, vol. 6.
- BHARGAVA, S. and MANOLI, D. (2015). Psychological frictions and the incomplete take-up of social benefits: Evidence from an IRS field experiment. *American Economic Review*, **105** 3489–3529.
- BIAU, G. (2012). Analysis of a random forests model. *Journal of Machine Learning Research*, **13** 1063–1095.
- BOGACHEV, V. I. (2007). *Measure theory*. Springer, Berlin–New York.
- BOLD, T., KIMENYI, M., MWABU, G., NG’ANG’A, A. and SANDEFUR, J. (2018). Experimental evidence on scaling up education reforms in Kenya. *Journal of Public Economics*, **168** 1–20.
- BRUHN, M., LEÃO, L. D. S., LEGOVINI, A., MARCHETTI, R. and ZIA, B. (2016). The impact of high school financial education: Evidence from a large-scale evaluation in Brazil. *American Economic Journal: Applied Economics*, **8** 256–295.
- BRUHN, M. and MCKENZIE, D. (2009). In pursuit of balance: Randomization in practice in development field experiments. *American Economic Journal: Applied Economics*, **1** 200–232. Publisher: American Economic Association.
- BUGNI, F. A., CANAY, I. A. and SHAIKH, A. M. (2018). Inference under covariate-adaptive randomization. *Journal of the American Statistical Association*, **113** 1784–1796.
- BUGNI, F. A., CANAY, I. A. and SHAIKH, A. M. (2019). Inference under covariate-adaptive randomization with multiple treatments. *Quantitative Economics*, **10** 1747–1785.
- BURSZTYN, L., FERMAN, B., FIORIN, S., KANZ, M. and RAO, G. (2018). Status goods: Experimental evidence from platinum credit cards. *The Quarterly Journal of Economics*, **133** 1561–1595.
- BURSZTYN, L., FIORIN, S., GOTTLIEB, D. and KANZ, M. (2019). Moral incentives in credit card debt repayment: Evidence from a field experiment. *Journal of Political Economy*. Forthcoming.
- BÜHLMANN, P. and VAN DE GEER, S. (2011). *Statistics for high-dimensional data: Methods, theory and applications*. Springer Science & Business Media.
- CALLEN, M., GULZAR, S., HASANAIN, A., KHAN, M. Y. and REZAEI, A. (2018). Data and policy decisions: Experimental evidence from Pakistan. Working paper.
- CALLEN, M., ISAQZADEH, M., LONG, J. D. and SPRENGER, C. (2014). Violence and risk preference: Experimental evidence from Afghanistan. *American Economic Review*, **104** 123–48.
- CARNEIRO, P., LEE, S. and WILHELM, D. (2019). Optimal data collection for randomized control trials. *The Econometrics Journal*. Forthcoming.

- CASABURI, L. and MACCHIAVELLO, R. (2019). Demand and supply of infrequent payments as a commitment device: Evidence from Kenya. *American Economic Review*, **109** 523–55.
- CHAMBAZ, A., VAN DER LAAN, M. J. and ZHENG, W. (2015). Targeted covariate-adjusted response-adaptive LASSO-based randomized controlled trials. *Modern Adaptive Randomized Clinical Trials: Statistical, Operational, and Regulatory Aspects* 345–368.
- CHATTERJEE, S. (2013). Assumptionless consistency of the Lasso. *arXiv preprint arXiv:1303.5817*.
- CHEN, X. (2007). Large sample sieve estimation of semi-nonparametric models. In *Handbook of Econometrics* (J. J. Heckman and E. E. Leamer, eds.), vol. 6. Elsevier, 5549–5632.
- CHEN, X. and WHITE, H. (1999). Improved rates and asymptotic normality for nonparametric neural network estimators. *IEEE Transactions on Information Theory*, **45** 682–691.
- CHEN, Y. and YANG, D. Y. (2019). The impact of media censorship: 1984 or Brave New World? *American Economic Review*, **109** 2294–2332.
- CHERNOZHUKOV, V., FERNÁNDEZ-VAL, I. and LUO, Y. (2018). The sorted effects method: Discovering heterogeneous effects beyond their averages. *Econometrica*, **86** 1911–1938.
- CHERNOZHUKOV, V., NEWEY, W. and SANTOS, A. (2015). Constrained conditional moment restriction models.
- CHETVERIKOV, D., SANTOS, A. and SHAIKH, A. M. (2018). The econometrics of shape restrictions. *Annual Review of Economics*, **10** 31–63.
- CHONG, A., COHEN, I., FIELD, E., NAKASONE, E. and TORERO, M. (2016). Iron deficiency and schooling attainment in Peru. *American Economic Journal: Applied Economics*, **8** 222–55.
- COX, D. and REID, N. (2000). *The theory of the design of experiments*. Chapman & Hall/CRC Monographs on Statistics & Applied Probability, CRC Press.
- CRÉPON, B., DEVOTO, F., DUFLO, E. and PARIENTÉ, W. (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** 123–50.
- DE CHAISEMARTIN, C. and RAMIREZ-CUELLAR, J. (2019). At what level should one cluster standard errors in paired experiments? *arXiv preprint arXiv:1906.00288*.
- DELLAVIGNA, S. and POPE, D. (2018). What motivates effort? Evidence and expert forecasts. *The Review of Economic Studies*, **85** 1029–1069.
- DENIL, M., MATHESON, D. and DE FREITAS, N. (2014). Narrowing the gap: Random forests in theory and in practice. In *International Conference on Machine Learning*. 665–673.
- DIZON-ROSS, R. (2019). Parents’ beliefs about their children’s academic ability: Implications for educational investments. *American Economic Review*, **109** 2728–2765.

- DUFLO, E. and BANERJEE, A. (2017). *Handbook of field experiments*. Elsevier Science.
- DUFLO, E., DUPAS, P. and KREMER, M. (2015). Education, HIV, and early fertility: Experimental evidence from Kenya. *American Economic Review*, **105** 2757–97.
- DUFLO, E., GLENNERSTER, R. and KREMER, M. (2007). Using randomization in development economics research: A toolkit. In *Handbook of Development Economics*, vol. 4. Elsevier, 3895–3962.
- DUPAS, P., KARLAN, D., ROBINSON, J. and UBFAL, D. (2018). Banking the unbanked? Evidence from three countries. *American Economic Journal: Applied Economics*, **10** 257–97.
- DUPAS, P. and ROBINSON, J. (2013). Savings constraints and microenterprise development: Evidence from a field experiment in Kenya. *American Economic Journal: Applied Economics*, **5** 163–92.
- FARRELL, M. H., LIANG, T. and MISRA, S. (2018). Deep neural networks for estimation and inference: Application to causal effects and other semiparametric estimands. *arXiv preprint arXiv:1809.09953*.
- FOGARTY, C. B. (2018a). On mitigating the analytical limitations of finely stratified experiments. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **80** 1035–1056.
- FOGARTY, C. B. (2018b). Regression-assisted inference for the average treatment effect in paired experiments. *Biometrika*, **105** 994–1000.
- FREEDMAN, D. A. (2008). On regression adjustments to experimental data. *Advances in Applied Mathematics*, **40** 180–193.
- FRYER, J., ROLAND G, DEVI, T. and HOLDEN, R. T. (2017). Vertical versus horizontal incentives in education: Evidence from randomized trials. Working paper.
- FRYER, R. (2017). Management and student achievement: Evidence from a randomized field experiment. Working paper.
- FRYER, R. (2018). The ”pupil” factory: Specialization and the production of human capital in schools. *American Economic Review*, **108** 616–656.
- GELMAN, A. and HILL, J. (2006). *Data analysis using regression and multilevel/hierarchical models*. Cambridge University Press.
- GLENNERSTER, R. and TAKAVARASHA, K. (2013). *Running randomized evaluations: A practical guide*. Princeton University Press.
- GLEWWE, P., PARK, A. and ZHAO, M. (2016). A better vision for development: Eyeglasses and academic performance in rural primary schools in China. *Journal of Development Economics*, **122** 170–182.
- GROH, M. and MCKENZIE, D. (2016). Macroinsurance for microenterprises: A randomized experiment in post-revolution Egypt. *Journal of Development Economics*, **118** 13–25.
- GRÖTSCHEL, M. and WAKABAYASHI, Y. (1990). Facets of the clique partitioning polytope. *Mathematical Programming*, **47** 367–387.

- GYÖRFI, L., KOHLER, M., KRZYŻAK, A. and WALK, H. (2006). *A distribution-free theory of nonparametric regression*. Springer Science & Business Media.
- HAHN, J., HIRANO, K. and KARLAN, D. (2011). Adaptive experimental design using the propensity score. *Journal of Business & Economic Statistics*, **29** 96–108.
- HEARD, K., O'TOOLE, E., NAIMPALLY, R. and BRESSLER, L. (2017). *Real world challenges to randomization and their solutions*. Boston, MA: Abdul Latif Jameel Poverty Action Lab.
- HOOVER, P. M. (1989). Maximality of randomized optimal designs. *The Annals of Statistics*, **17** 1315–1324.
- HORTON, J. J., RAND, D. G. and ZECKHAUSER, R. J. (2011). The online laboratory: Conducting experiments in a real labor market. *Experimental economics*, **14** 399–425.
- HSU, H. and LACHENBRUCH, P. A. (2007). Paired t-test. Wiley Online Library, 1–3.
- IMAI, K., KING, G., NALL, C. and OTHERS (2009). The essential role of pair matching in cluster-randomized experiments, with application to the Mexican universal health insurance evaluation. *Statistical Science*, **24** 29–53.
- IMBENS, G. W. (2011). Experimental design for unit and cluster randomized trials.
- IMBENS, G. W. and RUBIN, D. B. (2015). *Causal inference in statistics, social, and biomedical sciences*. Cambridge University Press.
- JOHANSSON, P., SCHULTZBERG, M. A. and RUBIN, D. (2019). On optimal re-randomization designs. Working paper.
- KALLUS, N. (2018). Optimal a priori balance in the design of controlled experiments. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **80** 85–112.
- KARLAN, D. and APPEL, J. (2016). *Failing in the field: What we can learn when field research goes wrong*. Princeton University Press.
- KARLAN, D. and WOOD, D. H. (2017). The effect of effectiveness: Donor response to aid effectiveness in a direct mail fundraising experiment. *Journal of Behavioral and Experimental Economics*, **66** 1–8.
- KARLAN, D. S. and ZINMAN, J. (2008). Credit elasticities in less-developed economies: Implications for microfinance. *American Economic Review*, **98** 1040–68.
- KASY, M. (2016). Why experimenters might not always want to randomize, and what they could do instead. *Political Analysis*, **24** 324–338.
- KHAN, A. Q., KHWAJA, A. I. and OLKEN, B. A. (2019). Making moves matter: Experimental evidence on incentivizing bureaucrats through performance-based postings. *American Economic Review*, **109** 237–70.
- KITAGAWA, T. and TETENOV, A. (2018). Who should be treated? Empirical welfare maximization methods for treatment choice. *Econometrica*, **86** 591–616.

- KUZIEMKO, I., NORTON, M. I., SAEZ, E. and STANTCHEVA, S. (2015). How elastic are preferences for redistribution? Evidence from randomized survey experiments. *American Economic Review*, **105** 1478–1508.
- LI, K.-C. (1983). Minimality for randomized designs: Some general results. *The Annals of Statistics*, **11** 225–239.
- LI, Q. and RACINE, J. S. (2007). *Nonparametric econometrics: Theory and practice*. Princeton University Press.
- LI, X., DING, P. and RUBIN, D. B. (2018). Asymptotic theory of rerandomization in treatment–control experiments. *Proceedings of the National Academy of Sciences*, **115** 9157–9162.
- LIN, W. (2013). Agnostic notes on regression adjustments to experimental data: Reexamining Freedman’s critique. *The Annals of Applied Statistics*, **7** 295–318.
- LIN, Y., ZHU, M. and SU, Z. (2015). The pursuit of balance: An overview of covariate-adaptive randomization techniques in clinical trials. *Contemporary Clinical Trials*, **45** 21–25.
- LIST, J. A. and RASUL, I. (2011). Field experiments in labor economics. vol. 4 of *Handbook of Labor Economics*. Elsevier, 103 – 228.
- LITMAN, L., ROBINSON, J. and ABBERBOCK, T. (2017). TurkPrime.com: A versatile crowdsourcing data acquisition platform for the behavioral sciences. *Behavior Research Methods*, **49** 433–442.
- MANSKI, C. F. (2004). Statistical treatment rules for heterogeneous populations. *Econometrica*, **72** 1221–1246.
- MASON, W. and SURI, S. (2012). Conducting behavioral research on Amazon’s Mechanical Turk. *Behavior Research Methods*, **44** 1–23.
- MBAKOP, E. and TABORD-MEEHAN, M. (2018). Model selection for treatment choice: Penalized welfare maximization. Working paper.
- MORGAN, K. L. and RUBIN, D. B. (2015). Rerandomization to balance tiers of covariates. *Journal of the American Statistical Association*, **110** 1412–1421.
- MORGAN, K. L., RUBIN, D. B. and OTHERS (2012). Rerandomization to improve covariate balance in experiments. *The Annals of Statistics*, **40** 1263–1282.
- MOSES, L. E. (2006). Matched pairs t-tests. In *Encyclopedia of Statistical Sciences*. American Cancer Society.
- MUNKRES, J. R. (1997). *Analysis on manifolds*. Westview Press.
- MURALIDHARAN, K., SINGH, A. and GANIMIAN, A. J. (2019). Disrupting education? Experimental evidence on technology-aided instruction in India. *American Economic Review*, **109** 1426–60.
- NEWBY, W. K. (1997). Convergence rates and asymptotic normality for series estimators. *Journal of Econometrics*, **79** 147–168.



- PANAGOPOULOS, C. and GREEN, D. P. (2008). Field experiments testing the impact of radio advertisements on electoral competition. *American Journal of Political Science*, **52** 156–168.
- PAOLACCI, G. and CHANDLER, J. (2014). Inside the Turk: Understanding Mechanical Turk as a participant pool. *Current Directions in Psychological Science*, **23** 184–188.
- PETERS, J., LANGBEIN, J. and ROBERTS, G. (2016). Policy evaluation, randomized controlled trials, and external validity—A systematic review. *Economics Letters*, **147** 51–54.
- PUKELSHEIM, F. (2006). *Optimal design of experiments*. Classics in Applied Mathematics, Society for Industrial and Applied Mathematics.
- RIACH, P. A. and RICH, J. (2002). Field experiments of discrimination in the market place. *The Economic Journal*, **112** F480–F518.
- ROSENBERGER, W. F. and LACHIN, J. M. (2015). *Randomization in clinical trials: Theory and Practice*. John Wiley & Sons.
- ROTHER, C. (2020). Flexible Covariate Adjustments in Randomized Experiments. Working paper.
- RUDIN, W. (1976). *Principles of mathematical analysis*, vol. 3. McGraw-hill New York.
- SCHULTZBERG, M. A. and JOHANSSON, P. (2019). Optimal designs and asymptotic inference. Working paper.
- SCORNET, E., BIAU, G., VERT, J.-P. and OTHERS (2015). Consistency of random forests. *The Annals of Statistics*, **43** 1716–1741.
- SONDHEIMER, R. M. and GREEN, D. P. (2010). Using experiments to estimate the effects of education on voter turnout. *American Journal of Political Science*, **54** 174–189.
- SPIVAK, M. (1965). *Calculus on manifolds*.
- STEINWART, I. and CHRISTMANN, A. (2008). *Support vector machines*. Springer Science & Business Media.
- TABORD-MEEHAN, M. (2020). Stratification trees for adaptive randomization in randomized controlled trials. Working paper.
- WAGER, S. and WALTHER, G. (2015). Adaptive concentration of regression trees, with application to random forests. *arXiv preprint arXiv:1503.06388*.
- WHITE, H. (1990). Connectionist nonparametric regression: Multilayer feedforward networks can learn arbitrary mappings. *Neural Networks*, **3** 535–549.
- WHITE, H. (2013). An introduction to the use of randomised control trials to evaluate development interventions. *Journal of Development Effectiveness*, **5** 30–49. Publisher: Taylor & Francis.
- WU, C.-F. (1981). On the robustness and efficiency of some randomized designs. *The Annals of Statistics*, **9** 1168–1177.