

# On the Efficiency of Finely Stratified Experiments \*

Yuehao Bai

Department of Economics  
University of Southern California

[yuehao.bai@usc.edu](mailto:yuehao.bai@usc.edu)

Jizhou Liu

Booth School of Business  
University of Chicago

[jliu32@chicagobooth.edu](mailto:jliu32@chicagobooth.edu)

Azeem M. Shaikh

Department of Economics  
University of Chicago

[amshaikh@uchicago.edu](mailto:amshaikh@uchicago.edu)

Max Tabord-Meehan

Department of Economics  
University of Chicago

[maxtm@uchicago.edu](mailto:maxtm@uchicago.edu)

August 23, 2024

## Abstract

This paper studies the use of finely stratified designs for the efficient estimation of a large class of treatment effect parameters that arise in the analysis of experiments. By a “finely stratified” design, we mean experiments in which units are divided into groups of a fixed size and a proportion within each group is assigned to a binary treatment uniformly at random. The class of parameters considered are those that can be expressed as the solution to a set of moment conditions constructed using a known function of the observed data. They include, among other things, average treatment effects, quantile treatment effects, and local average treatment effects as well as the counterparts to these quantities in experiments in which the unit is itself a cluster. In this setting, we establish two results. First, we show that under a finely stratified design, the naïve method of moments estimator achieves the same asymptotic variance as what could typically be attained under alternative treatment assignment schemes only through *ex post* covariate adjustment. Second, we argue that in fact the naïve method of moments estimator under a finely stratified design is asymptotically efficient by deriving a lower bound on the asymptotic variance of “regular” estimators of the parameter of interest in the form of a convolution theorem. This result accommodates a large class of possible treatment assignment schemes that are used routinely throughout the sciences, such as stratified block randomization and matched pairs. In this sense, “finely stratified” experiments are attractive because they lead to efficient estimators of treatment effect parameters “by design” rather than through *ex post* covariate adjustment and thereby remain “hands above the table.”

KEYWORDS: Convolution Theorem, Efficiency, Experiment, Experimental design, Finely stratified experiment, Matched pairs, Randomized controlled trial

JEL classification codes: C12, C14

---

\*We thank the seminar participants at UC Riverside, UC Irvine, Cornell University, and Syracuse University for helpful comments. The fourth author acknowledges support from NSF grant SES-2149408.

# 1 Introduction

This paper studies the use of finely stratified designs for the efficient estimation of a large class of treatment effect parameters that arise in the analysis of experiments. By a “finely stratified” design, we mean experiments in which units are divided into groups of a fixed size and a proportion within each group is assigned to a binary treatment uniformly at random. For example, when the fixed size equals two and the marginal probability of treatment assignment is specified to be one half, such a design is simply a matched pairs design. The class of parameters considered are those that can be expressed as the solution to a set of moment conditions constructed using a known function of the observed data. This class of parameters includes many treatment effect parameters of interest: average treatment effects (ATEs), quantile treatment effects, and local average treatment effects as well as the counterparts to these quantities in experiments in which the unit is itself a cluster.

In the setting described above, we establish two results. First, we study the asymptotic properties of a naïve method of moments estimator under a finely stratified design. Here, by a naïve method of moments estimator, we mean an estimator constructed using a direct sample analog of the moment conditions. For example, in the case of the ATE, such an estimator is given by a difference-in-means. We show that under a finely stratified design, the naïve method of moments estimator achieves the same asymptotic variance as what could typically be attained under alternative treatment assignment schemes only through *ex post* covariate adjustment. Such adjustment strategies frequently involve the nonparametric estimation of conditional expectations or similar quantities; see, for example, [Zhang et al. \(2008\)](#), [Tsiatis et al. \(2008\)](#), [Jiang et al. \(2022a\)](#), [Jiang et al. \(2022b\)](#) and [Rafi \(2023\)](#).<sup>1</sup> Second, we argue that in fact the naïve method of moments estimator under a finely stratified design is asymptotically efficient by deriving a lower bound on the asymptotic variance of “regular” estimators of the parameter of interest in the form of a convolution theorem. This result accommodates a large class of possible treatment assignment schemes that are used routinely throughout the sciences, such as stratified block randomization and matched pairs.<sup>2</sup> We emphasize that, since treatment assignment schemes used in experiments are generally not i.i.d., classical results in, for example, [van der Vaart \(1998\)](#), do not directly apply. Together, these two results show that finely stratified experiments lead to efficient estimators that prioritize transparency in that they preclude the researcher from “data snooping” associated with *ex post* nonparametric covariate adjustment. Importantly, concerns with this type of data snooping are not completely eliminated by typical pre-analysis plans because such adjustments involve choices, such as the choice of nonparametric estimator or tuning parameters, that are often not pre-registered prior to the experiment.<sup>3</sup> The estimators are therefore attractive because they effectively perform nonparametric regression adjustment “by design” ([Cytrynbaum, 2023b](#)) and thereby remain “hands above the table” ([Freedman, 2008](#); [Lin, 2013](#)).

Our paper builds upon two strands of literature. The first strand of literature concerns the analysis of

---

<sup>1</sup>For related results in the context of observational data see [Newey \(1994\)](#), [Hahn \(1998\)](#), [Heckman et al. \(1998\)](#), [Imbens \(2004\)](#), [Frolich \(2007\)](#), [Firpo \(2007\)](#), [Imbens et al. \(2007\)](#), [Farrell \(2015\)](#), [Abadie and Imbens \(2016\)](#), and [Chernozhukov et al. \(2017\)](#).

<sup>2</sup>For a discussion of such treatment assignment schemes focused on clinical trials, see [Rosenberger and Lachin \(2015\)](#); for reviews focused on development economics, see [Duflo et al. \(2007\)](#) and [Bruhn and McKenzie \(2009\)](#).

<sup>3</sup>We emphasize that parametric covariate adjustment would, in general, not lead to efficiency. Furthermore, it too would require choices of covariates and functional form that are also often not pre-registered prior to the experiment.

finely stratified experiments. Within this literature, our analysis is most closely related to [Bai et al. \(2022\)](#), who derive the asymptotic behavior of the difference-in-means estimator of the ATE when treatment is assigned according to a matched pairs design, and [Cytrynbaum \(2023b\)](#), who develops related results for an experimental design referred to as “local randomization” that permits the proportion of units assigned to treatment to vary with the baseline covariates. Beyond settings that study estimation of the ATE, [Bai et al. \(2024c\)](#) develops results for the analysis of different cluster-level average treatment effects and [Jiang et al. \(2021\)](#) develop results analogous to those in [Bai et al. \(2022\)](#) for suitable estimators of the quantile treatment effect. To our knowledge, our paper is the first to analyze the properties of finely stratified experiments in a general framework which accommodates all of the above parameters as well as any parameter that can be characterized as the solutions to a set of moment conditions involving a known function of the observed data. We emphasize that none of the above papers formally establish the asymptotic efficiency of finely stratified experiments. The second strand of literature concerns bounds on the efficiency with which treatment effect parameters can be estimated in experiments. Two important recent papers in this literature studying efficiency bounds in the special case of estimating the ATE are [Armstrong \(2022\)](#) and [Rafi \(2023\)](#). Even in this special case, their results differ from ours in important and empirically relevant ways; [Remark 4.4](#) provides an in-depth discussion of the connection between these results and ours. See also [Bai \(2022\)](#) for some finite-sample optimality properties of matched pairs designs for estimation of the ATE.

The remainder of this paper is organized as follows. In [Section 2](#), we describe our setup and notation. We emphasize in particular the way in which our framework can accommodate various treatment effect parameters of interest. [Section 3](#) derives the asymptotic behavior of the naïve method of moments estimator of our parameter of interest when treatment is assigned using a finely stratified design. In [Section 4](#), we develop our lower bound on the asymptotic variance of “regular” estimators of these parameters and show that it is achieved by the naïve method of moments estimator in a finely stratified design. In [Section 5](#), we illustrate the practical relevance of our theoretical results by comparing the mean-squared errors of the naïve method of moments estimators for the ATE and local average treatment effect when treatment status is assigned according to a matched pairs design versus that of an estimator using *ex post* covariate adjustment when treatment status is assigned in an i.i.d. fashion. Finally, we conclude in [Section 6](#) with some recommendations for empirical practice guided by both these simulations and our theoretical results. Proofs of all results can be found in the Appendix.

## 2 Setup and Motivation

Let  $A_i \in \{0, 1\}$  denote the treatment status of the  $i$ th unit, and let  $X_i \in \mathbf{R}^{d_x}$  denote their observed, baseline covariates. For  $a \in \{0, 1\}$ , let  $R_i(a) \in \mathbf{R}^{d_r}$  denote a *vector* of potential responses. As we illustrate below, considering a vector of responses allows us to accommodate certain parameters of interest. Let  $R_i \in \mathbf{R}^{d_r}$  denote the vector of observed responses obtained from  $R_i(a)$  once treatment is assigned. As usual, the observed responses and potential responses are related to treatment status by the relationship

$$R_i = R_i(1)A_i + R_i(0)(1 - A_i) . \tag{1}$$

We assume throughout that our sample consists of  $n$  units. For any random vector indexed by  $i$ , for example  $A_i$ , we define  $A^{(n)} = (A_1, \dots, A_n)$ . Let  $P_n$  denote the distribution of the observed data  $(R^{(n)}, A^{(n)}, X^{(n)})$ , and  $Q_n$  the distribution of  $(R^{(n)}(1), R^{(n)}(0), X^{(n)})$ . We assume  $Q_n = Q^n$ , where  $Q$  is the marginal distribution of  $(R_i(1), R_i(0), X_i)$ . Given  $Q_n$ ,  $P_n$  is then determined by (1) and the mechanism for determining treatment assignment. We assume that treatment assignment is performed such that a standard unconfoundedness assumptions holds and such that the probability of assignment given  $X$  is some known constant for every  $1 \leq i \leq n$ , as is often the case in most experiments:

**Assumption 2.1.** Treatment status is assigned so that

$$(R^{(n)}(1), R^{(n)}(0)) \perp\!\!\!\perp A^{(n)} | X^{(n)}, \quad (2)$$

and such that  $P\{A_i = 1 | X_i = x\} = \eta$ , for some  $\eta \in (0, 1)$  for all  $1 \leq i \leq n$ .

Assumption 2.1 restricts the probability of assignment to be the fixed fraction  $\eta$  across the entire experimental sample, but this restriction can be weakened so that  $\eta$  is replaced by  $\eta(X_i)$  for many of our subsequent results: see, in particular, Remarks 3.5, 4.3, and 4.4. Given Assumption 2.1, it can be shown that  $(X_i, A_i, R_i)$  are identically distributed for  $1 \leq i \leq n$ , and their marginal distribution does not change with  $n$  (see Lemma B.5 in the Appendix). As a consequence, we denote the marginal distribution of  $(X_i, A_i, R_i)$  by  $P$ . We consider parameters  $\theta_0 \in \Theta \subset \mathbf{R}^{d_\theta}$  which can be defined as the solution to a set of moment equalities. In particular, let  $m : \mathbf{R}^{d_x} \times \{0, 1\} \times \mathbf{R}^{d_r} \rightarrow \mathbf{R}^{d_\theta}$  be a known measurable function, and we consider parameters  $\theta_0$  that uniquely solve the moment equality

$$E_P[m(X_i, A_i, R_i, \theta_0)] = 0. \quad (3)$$

We emphasize that  $m(\cdot)$  is not a function of any unknown nuisance parameters, but may depend on the known value of  $\eta$  in Assumption 2.1. We present five examples of well-known parameters that can be described as (functions of) solutions to a set of moment conditions as in (3).

**Example 2.1** (Average Treatment Effect). Let  $Y_i(a) = R_i(a)$  denote a scalar potential outcome for the  $i$ th unit under treatment  $a \in \{0, 1\}$ , and let  $Y_i = R_i$  denote the observed outcome. Let  $\theta_0 = E_Q[Y_i(1) - Y_i(0)]$  denote the average treatment effect (ATE). Under Assumption 2.1,  $\theta_0$  solves the moment condition in (3) with

$$m(X_i, A_i, R_i, \theta) = \frac{Y_i A_i}{\eta} - \frac{Y_i(1 - A_i)}{1 - \eta} - \theta. \quad (4)$$

For a list of papers which consider estimators based on (4), see Hirano and Imbens (2001) and Hirano et al. (2003). ■

**Example 2.2** (Quantile Treatment Effect). Let  $Y_i(a) = R_i(a)$  denote a scalar potential outcome for the  $i$ th unit under treatment  $a \in \{0, 1\}$ , and let  $Y_i = R_i$  denote the observed outcome. Let  $\tau \in (0, 1)$  and  $\theta_0 = (\theta_0(1), \theta_0(0))' = (q_{Y(1)}(\tau), q_{Y(0)}(\tau))'$ , where

$$q_{Y(a)}(\tau) = \inf\{\lambda \in \mathbf{R} : Q\{Y_i(a) \leq \lambda\} \geq \tau\}.$$

In other words,  $\theta_0$  is defined to be the vector of  $\tau$ th quantiles of the marginal distributions of  $Y_i(1)$  and  $Y_i(0)$ . If we assume  $q_{Y(a)}(\tau)$  is unique for  $a \in \{0, 1\}$  in the sense that  $Q\{Y(a) \leq q_{Y(a)}(\tau) + \epsilon\} > Q\{Y(a) \leq q_{Y(a)}(\tau)\}$  for all  $\epsilon > 0$ , then it follows from Assumption 2.1 and Lemma 1 in Firpo (2007) that  $\theta_0$  solves the moment condition in (3) with

$$m(X_i, A_i, R_i, \theta) = \left( \frac{A_i(\tau - I\{Y_i \leq \theta(1)\})}{(1 - A_i)(\tau - I\{Y_i \leq \theta(0)\})} \frac{\eta}{1 - \eta} \right),$$

for  $\theta = (\theta^{(1)}, \theta^{(0)})'$ . Note that the quantile treatment effect  $q_{Y(1)}(\tau) - q_{Y(0)}(\tau)$  can then be defined as  $h(\theta_0)$  where  $h : \mathbf{R}^2 \rightarrow \mathbf{R}$  is given by  $h(s, t) = s - t$ . ■

**Example 2.3** (Local Average Treatment Effect). Let  $(\tilde{Y}_i(a), D_i(a)) = R_i(a)$  denote the vector of potential outcomes and treatment take-up under treatment  $a \in \{0, 1\}$ , and let  $(Y_i, D_i) = R_i$  denote the vector of observed outcomes and treatment take-up. Note here that  $\tilde{Y}_i(a)$  corresponds to the potential outcome under assignment  $a \in \{0, 1\}$  and not to the potential outcome for a given take-up  $D_i = d$ . Suppose  $E_Q[D_i(1) - D_i(0)] \neq 0$  and let

$$\theta_0 = \frac{E_Q[\tilde{Y}_i(1) - \tilde{Y}_i(0)]}{E_Q[D_i(1) - D_i(0)]}.$$

It then follows from Assumption 2.1 that  $\theta_0$  solves the moment condition in (3) with

$$m(X_i, A_i, R_i, \theta) = \frac{Y_i A_i}{\eta} - \frac{Y_i(1 - A_i)}{1 - \eta} - \theta \left( \frac{D_i A_i}{\eta} - \frac{D_i(1 - A_i)}{1 - \eta} \right). \quad (5)$$

If we further assume instrument monotonicity (i.e.,  $P\{D_i(1) \geq D_i(0)\} = 1$ ) and instrument exclusion, then  $\theta_0$  could be re-interpreted as the local average treatment effect (LATE) in the sense of Imbens and Angrist (1994). ■

**Example 2.4** (Weighted Average Treatment Effect). Let  $Y_i(a) = R_i(a)$  denote a scalar potential outcome for the  $i$ th unit under treatment  $a \in \{0, 1\}$ , and let  $Y_i = R_i$  denote the observed outcome. Let

$$\theta_0 = E_Q \left[ \frac{\omega(X_i)}{E_Q[\omega(X_i)]} (Y_i(1) - Y_i(0)) \right],$$

for some known function  $\omega : \mathbf{R}^{d_x} \rightarrow \mathbf{R}$ . It then follows from Assumption 2.1 that  $\theta_0$  solves the moment condition in (3) with

$$m(X_i, A_i, R_i, \theta) = \omega(X_i) \left( \frac{Y_i A_i}{\eta} - \frac{Y_i(1 - A_i)}{1 - \eta} \right) - \omega(X_i) \theta.$$

Note that  $\theta_0$  defined in this way can accommodate the (cluster) size-weighted and equally-weighted average treatment effects considered in Bugni et al. (2022) and Bai et al. (2024c) in the context of cluster-level randomized controlled trials. ■

**Example 2.5** (Log-Odds Ratio). Let  $Y_i(a) = R_i(a) \in \{0, 1\}$  denote a binary potential outcome for the  $i$ th unit under treatment  $a \in \{0, 1\}$ , and let  $Y_i = R_i$  denote the observed outcome. Suppose  $0 < P\{Y_i(a) =$

$0\} < 1$  for  $a \in \{0, 1\}$ , and let  $\theta_0 = (\theta_0(1), \theta_0(2))'$ , where

$$\theta_0(1) = \text{logit}(E_Q[Y_i(0)]) ,$$

$$\theta_0(2) = \text{logit}(E_Q[Y_i(1)]) - \text{logit}(E_Q[Y_i(0)]) ,$$

with  $\text{logit}(z) = \log(\frac{z}{1-z})$ , so that  $\theta_0(2)$  denotes the log-odds ratio of treatment 1 relative to treatment 0. It follows from Assumption 2.1 that  $\theta_0$  solves the moment condition in (3) with

$$m(X_i, A_i, R_i, \theta) = \begin{pmatrix} 1 - A_i \\ A_i \end{pmatrix} (Y_i - \text{expit}(\theta(1) + \theta(2)A_i)) ,$$

where  $\text{expit}(z) = \frac{\exp(z)}{1+\exp(z)}$ . The log-odds ratio can then be defined as  $h(\theta_0)$  where  $h : \mathbf{R}^2 \rightarrow \mathbf{R}$  is given by  $h(s, t) = t$ . This parameter appears in, for example, Zhang et al. (2008). ■

Additional examples could be obtained by considering combinations of Examples 2.1–2.5. For instance, combining the moment functions from Examples 2.3 and 2.4 would result in a weighted LATE parameter. Beyond these examples, certain treatment effect contrasts could also be related to the structural parameters in, for instance, a model of supply in demand: see, for example, the model estimated in Casaburi and Reed (2022).

Throughout the rest of the paper we consider the asymptotic properties of the method of moments estimator  $\hat{\theta}_n$  for  $\theta_0$  which is constructed as a solution to the sample analogue of (3):

$$\frac{1}{n} \sum_{1 \leq i \leq n} m(X_i, A_i, R_i, \hat{\theta}_n) = 0 . \tag{6}$$

Note that  $\hat{\theta}_n$  as defined in (6) is closely related to standard estimators of the parameter  $\theta_0$  in specific examples. For instance, in Example 2.1,

$$\hat{\theta}_n = \frac{1}{\eta} \sum_{1 \leq i \leq n} Y_i A_i - \frac{1}{1-\eta} \sum_{1 \leq i \leq n} Y_i (1 - A_i) ,$$

so that  $\hat{\theta}_n$  is a Horvitz-Thompson analogue of the standard difference-in-means estimator for the ATE. In Example 2.3,

$$\hat{\theta}_n = \frac{\frac{1}{\eta} \sum_{1 \leq i \leq n} Y_i A_i - \frac{1}{1-\eta} \sum_{1 \leq i \leq n} Y_i (1 - A_i)}{\frac{1}{\eta} \sum_{1 \leq i \leq n} D_i A_i - \frac{1}{1-\eta} \sum_{1 \leq i \leq n} D_i (1 - A_i)} ,$$

so that  $\hat{\theta}_n$  is a Horvitz-Thompson analogue of the standard Wald estimator for the local average treatment effect.

To illustrate the key contribution of our paper, note that if  $A^{(n)}$  were assigned i.i.d., independently of  $X^{(n)}$ , then it can be shown under mild conditions on  $m(\cdot)$  (see, for instance, Theorem 5.1 in van der Vaart, 1998) that the naïve method of moments estimator satisfies

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{d} N(0, \mathbb{V}) ,$$

where

$$\mathbb{V} = M^{-1} E_P[m(X_i, A_i, R_i, \theta_0)m(X_i, A_i, R_i, \theta_0)'](M^{-1})' , \quad (7)$$

with  $M = \frac{\partial}{\partial \theta} E_P[m(X, A, R, \theta)] \Big|_{\theta=\theta_0}$ . In Section 3, we show that if we assign  $A^{(n)}$  using a finely stratified design (i.e., a treatment assignment scheme which uses the covariates  $X^{(n)}$  to block units into groups of fixed size: see Assumption 3.1 below for a formal definition), then

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{d} N(0, \mathbb{V}_*) ,$$

where  $\mathbb{V} \geq \mathbb{V}_*$  (see Theorem 3.1). Under i.i.d. assignment, the naïve method of moment estimator  $\hat{\theta}_n$  cannot generally attain  $\mathbb{V}_*$ , but an estimator that attains  $\mathbb{V}_*$  could instead be constructed by appropriately “augmenting” the moment function, and then considering an estimator which solves the augmented moment equation. For instance, if we consider the ATE in Example 2.1, then it is straightforward to show that the following augmented moment function identifies  $\theta_0$ :

$$m^*(X_i, A_i, R_i, \theta) = \left( \frac{A_i(Y_i - \mu_1(X_i))}{\eta} - \frac{(1 - A_i)(Y_i - \mu_0(X_i))}{1 - \eta} + \mu_1(X_i) - \mu_0(X_i) \right) - \theta , \quad (8)$$

where  $\mu_a(X_i) = E_Q[Y_i(a)|X_i]$ . This choice of  $m^*(\cdot)$  produces the well known doubly-robust moment condition for estimating the ATE (Robins et al., 1995; Hahn, 1998). It can then be shown that an appropriately constructed two-step estimator, where  $\mu_1(\cdot)$  and  $\mu_0(\cdot)$  are non-parametrically estimated in a first step, attains  $\mathbb{V}_*$  (Tsiatis et al., 2008; Farrell, 2015; Chernozhukov et al., 2017; Rafi, 2023). Intuitively, the estimator obtained from the augmented moment function  $m^*(\cdot)$  performs nonparametric regression adjustment by exploiting the information contained in  $X^{(n)}$  that may not have been captured in the original moment function  $m(\cdot)$ . Similar nonparametric regression adjustments based on augmented moment equations have been developed for other parameters of interest (Zhang et al., 2008; Belloni et al., 2017; Jiang et al., 2022a,b). In this sense, we show that “fine stratification” can perform nonparametric regression adjustment “by design” for the large class of parameters that can be expressed in terms of moment conditions of the form given in (3), thus generalizing similar observations made in Bai et al. (2022), Bai (2022), and Cytrynbaum (2023b) in the special case of estimating the ATE.

Earlier work on efficient treatment effect estimation has noted that the variance  $\mathbb{V}_*$  is in fact the efficiency bound for estimating  $\theta_0$  under i.i.d. assignment. A natural follow-up question is whether or not  $\mathbb{V}_*$  continues to be the efficiency bound for estimating  $\theta_0$  under a finely stratified design, or more generally for complex experimental designs which induce dependence in the treatment assignments across individuals in the experiment. In Section 4, we show that  $\mathbb{V}_*$  continues to be the efficiency bound for estimating  $\theta_0$  for a large class of treatment assignment schemes with a fixed marginal probability of treatment assignment, which includes finely stratified designs as a special case. We can thus conclude that, from the perspective of asymptotic efficiency, finely stratified designs are optimal experimental designs for a broad range of treatment effect estimation problems.

### 3 The Asymptotic Variance of Finely Stratified Experiments

In this section, we derive the limiting distribution of the method of moments estimator  $\hat{\theta}_n$  when treatment is assigned by fine stratification over the baseline covariates  $X^{(n)}$ . Such assignment mechanisms use the covariates  $X^{(n)}$  to group units with “similar” covariate values into blocks of fixed size, and then assign treatment completely at random within each block. In order to describe this assignment mechanism formally, we require some further notation to define the blocks of units. Let  $\ell$  and  $k$  be arbitrary positive integers with  $\ell \leq k$  and set  $\eta = \ell/k$ . For simplicity, assume that  $n$  is divisible by  $k$ . We then represent blocks of units using a partition of  $\{1, \dots, n\}$  given by

$$\left\{ \lambda_j = \lambda_j(X^{(n)}) \subseteq \{1, \dots, n\}, 1 \leq j \leq n/k \right\},$$

with  $|\lambda_j| = k$ . Because of its possible dependence on  $X^{(n)}$ ,  $\{\lambda_j : 1 \leq j \leq n/k\}$  encompasses a variety of different ways of blocking the  $n$  units according to the observed, baseline covariates. Given such a partition, we assume that treatment status is assigned as described in the following assumption:

**Assumption 3.1.** Treatment status is assigned so that  $(R^{(n)}(1), R^{(n)}(0)) \perp\!\!\!\perp A^{(n)} | X^{(n)}$  and, conditional on  $X^{(n)}$ ,

$$\{(A_i : i \in \lambda_j) : 1 \leq j \leq n/k\}$$

are i.i.d. and each uniformly distributed over all permutations of  $\underbrace{(0, 0, \dots, 0)}_{k-\ell}, \underbrace{(1, 1, \dots, 1)}_{\ell}$ .

The assignment mechanism described in Assumptions 3.1 generalizes the definition of a matched pairs design. In particular, we recover a matched pairs design if we set  $(\ell, k) = (1, 2)$ , with  $\eta = 1/2$ . Indeed, suppose  $n$  is even and consider pairing the experimental units into  $n/2$  pairs, represented by the sets

$$\{\pi(2j-1), \pi(2j)\} \text{ for } j = 1, \dots, n/2,$$

where  $\pi = \pi_n(X^{(n)})$  is a permutation of  $n$  elements. Because of its possible dependence on  $X^{(n)}$ ,  $\pi$  encompasses a broad variety of ways of pairing the  $n$  units according to the observed, baseline covariates  $X^{(n)}$ . Given such a  $\pi$ , we assume that treatment status is assigned so that Assumption 3.1 holds and, conditional on  $X^{(n)}$ ,  $(A_{\pi(2j-1)}, A_{\pi(2j)}), j = 1, \dots, n/2$  are i.i.d. and each uniformly distributed over the values in  $\{(0, 1), (1, 0)\}$ . For some examples of such an assignment mechanism being used in practice, see, for instance, Angrist and Lavy (2009), Banerjee et al. (2015), and Bruhn et al. (2016).

**Remark 3.1.** Note that Assumption 3.1 generalizes matched pairs designs along two dimensions: first, it allows for treatment fractions other than  $\eta = 1/2$ . Second, it allows for choices of  $\ell$  and  $k$  which are not relatively prime. For instance, if we set  $(\ell, k) = (2, 4)$ , then  $\eta = 1/2$  as in matched pairs, but now the assignment mechanism blocks units into groups of size 4 and assigns two units to treatment, two units to control. Although Theorem 3.1 below establishes that allowing for this level of flexibility has no effect on the asymptotic properties of our estimator, in our experience we have found that designs which employ these treatment “replicates” in each block can simplify the construction of variance estimators in practice. See Appendix A for further discussion. ■



Our analysis will require some discipline on the way in which the blocks are formed. In particular, we will require that the units in each block be “close” in terms of their baseline covariates in the sense described by the following assumption:

**Assumption 3.2.** The blocks used in determining treatment status satisfy

$$\frac{1}{n} \sum_{1 \leq j \leq n/k} \max_{i, i' \in \lambda_j} \|X_i - X_{i'}\|^2 \xrightarrow{P} 0.$$

Bai et al. (2022) and Cytrynbaum (2023b) discuss blocking algorithms that satisfy Assumption 3.2. When  $X_i \in \mathbf{R}$  and  $E[X_i^2] < \infty$ , a simple algorithm that satisfies Assumption 3.2 is to simply order units from smallest to largest and then block adjacent units into blocks of size  $k$ . In the case of matched pairs, if  $\dim(X_i) > 1$  and  $E[\|X_i\|^d] < \infty$  for  $d \geq \dim(X_i) + 1$ , then Assumption 3.2 is satisfied by the `nbpmatching` algorithm in `R` that minimizes the sum of squared distances of  $X$  within pairs. See Appendix A of Bai et al. (2024c) for details.

The next two sets of assumptions allow us to derive the large-sample properties of  $\hat{\theta}_n$ . We impose Assumption 3.3 to establish the consistency of  $\hat{\theta}_n$ , and we further impose Assumption 3.4 to establish its limiting distribution. In what follows, when writing expectations and variances, we suppress the subscripts  $P$  and  $Q$  whenever doing so does not lead to confusion.

**Assumption 3.3.** Let  $m(\cdot) = (m_s(\cdot) : 1 \leq s \leq d_\theta)'$ . Then the moment functions are such that

- (a) For every  $\epsilon > 0$ ,  $\inf_{\theta \in \Theta : \|\theta - \theta_0\| > \epsilon} \|E[m(X_i, A_i, R_i, \theta)]\| > 0$ .
- (b) For  $1 \leq s \leq d_\theta$ ,  $\{m_s(x, a, r, \theta) : \theta \in \Theta\}$  with  $a \in \{0, 1\}$  fixed is a VC-class of functions in  $(x, r)$ .
- (c) For  $1 \leq s \leq d_\theta$ ,  $\{m_s(x, a, r, \theta) : \theta \in \Theta\}$  is pointwise measurable in the sense that there exists a countable set  $\Theta^*$  such that for each  $\theta \in \Theta$ , there exists a sequence  $\{\theta_m\} \subset \Theta^*$  such that  $m_s(x, a, r, \theta_m) \rightarrow m_s(x, a, r, \theta)$  as  $m \rightarrow \infty$  for all  $x, a, r$ .
- (d)  $E \left[ \sup_{\theta \in \Theta^*} \|m(X, a, R(a), \theta)\| \right] < \infty$  for  $a \in \{0, 1\}$ .
- (e) For some  $K < \infty$ ,

$$\sup_{\theta \in \Theta^*} \|E[m(X, a, R(a), \theta) | X = x] - E[m(X, a, R(a), \theta) | X = x']\| \leq K \|x - x'\|$$

for all  $x, x' \in \mathbf{R}^{d_x}$ .

- (f)  $E \left[ \sup_{\theta \in \Theta^*} \|m(X, a, R(a), \theta)\| \middle| X = x \right]$  is Lipschitz for  $a \in \{0, 1\}$ .

Assumption 3.3(a) is a standard assumption to ensure the solution to (3) is “well separated.” It appears as a condition, for instance, in Theorem 5.9 in van der Vaart (1998). Assumption 3.3(b) can be readily verified in Examples 2.1–2.5 because the moment conditions are either constructed as linear functions in  $\theta$  (multiplied or composed with fixed functions), or dependent on  $\theta$  through indicator functions. Assumption

3.3(c) is a standard condition to guarantee the measurability of the supremum of a suitable class of functions. In particular, it allows us to define expectations of suprema without invoking outer expectations. See Example 2.3.4 in [van der Vaart and Wellner \(1996\)](#) for details. Assumption 3.3(d) guarantees the existence of an envelope function needed to establish a uniform law of large numbers. Assumptions 3.3(e)–(f) mirror common assumptions used when studying matched pairs designs to ensure units that are close in terms of the baseline covariates are also close in terms of their moments.

**Assumption 3.4.** Let  $m(\cdot) = (m_s(\cdot) : 1 \leq s \leq d_\theta)'$ . The moment functions are such that

- (a)  $E[m(X_i, A_i, R_i, \theta)]$  is differentiable at  $\theta_0$  with a nonsingular derivative  $M$ .
- (b) For  $\Theta^*$  in Assumption 3.3(c),  $E\left[\sup_{\theta \in \Theta^*} \|m(X, a, R(a), \theta)\|^2\right] < \infty$  for  $a \in \{0, 1\}$ .
- (c) For  $1 \leq s \leq d_\theta$ ,  $E[(m_s(X, a, R(a), \theta) - m_s(X, a, R(a), \theta_0))^2] \rightarrow 0$  as  $\theta \rightarrow \theta_0$  for  $a \in \{0, 1\}$ .
- (d) For  $1 \leq s \leq d_\theta$ ,  $\{E[m_s(X, a, R(a), \theta)|X = x] : \theta \in \Theta\}$  is a VC-class of functions for  $a \in \{0, 1\}$ .
- (e) For  $\Theta^*$  in Assumption 3.3(c) and some  $K < \infty$ , for  $1 \leq s \leq d_\theta$ ,

$$\begin{aligned} \sup_{\theta \in \Theta^*} |E[m_s^2(X, a, R(a), \theta)|X = x] - E[m_s^2(X, a, R(a), \theta)|X = x']| &\leq K\|x - x'\| \\ \sup_{\theta \in \Theta^*} |E[m_s(X, a, R(a), \theta)m_s(X, a, R(a), \theta_0)|X = x] \\ &- E[m_s(X, a, R(a), \theta)m_s(X, a, R(a), \theta_0)|X = x']| \leq K\|x - x'\|, \end{aligned}$$

for all  $x, x' \in \mathbf{R}^{d_x}$ ,  $a \in \{0, 1\}$

- (f) For  $\Theta^*$  in Assumption 3.3(c),  $E\left[\sup_{\theta \in \Theta^*} \|m(X, a, R(a), \theta)\|^2 \middle| X = x\right]$  is Lipschitz for  $a \in \{0, 1\}$ .

Assumption 3.4(a) is a standard assumption used when deriving the properties of  $Z$ -estimators. See, for instance, Theorem 3.1 in [Newey and McFadden \(1994\)](#) and Theorem 5.21 in [van der Vaart \(1998\)](#). Because differentiability is imposed on their expectations instead of the moment functions themselves, the moment functions are allowed to be nonsmooth as in Example 2.2. Assumption 3.4(b) guarantees the existence of an envelope function needed to establish a uniform law of large numbers. Assumption 3.4(c) implies the moment functions are mean-square continuous in  $\theta$ . Assumptions 3.4(e)–(f) again mirror common assumptions used when studying matched pairs to ensure units that are close in terms of the baseline covariates are also close in terms of their moments. Assumption 3.4(d) is again readily verified in Examples 2.1, 2.3–2.5 because  $\theta$  enters separably in these examples. To verify the assumption for Example 2.2, note that for any random variables  $Y(a), X$ , the subgraphs  $\{(x, t) : t < P\{Y(a) \leq \theta(a)|X = x\}\}$  are linearly ordered in  $\theta(a)$  because the conditional distribution function is increasing in  $\theta(a)$ . Therefore, the class of subgraphs is VC with index 2 (see, for instance, the last sentence of the proof of Lemma 2.6.16 in [van der Vaart and Wellner, 1996](#)).

**Remark 3.2.** We note that some of the assumptions imposed in Assumptions 3.3 and 3.4 are seemingly more stringent than the low-level conditions considered in previous papers which study inference for certain specific parameters of interest under matched pairs designs ([Bai et al., 2022](#); [Jiang et al., 2021](#); [Cytrynbaum,](#)

2023b; Bai et al., 2024a). We suspect that, with more delicate arguments, some of these assumptions could be weakened for specific parameters of interest: for example, following an approximation argument in Cytrynbaum (2023b) in the special case of the ATE, we expect Assumptions 3.3(e) and 3.4(e) could be dropped whenever the moment functions  $m_s(\cdot)$  are linear in the parameter  $\theta$ . However, in order to accommodate more general, possibly nonlinear moment functions, we do not pursue this further in the paper. ■

The following theorem establishes the asymptotic variance of the “naïve” method of moments estimator when the treatment assignment mechanism is finely stratified in the sense of satisfying Assumptions 3.1–3.2.

**Theorem 3.1.** *Suppose the treatment assignment mechanism satisfies Assumptions 3.1–3.2 and the moment functions satisfy Assumptions 3.3–3.4. Let  $\hat{\theta}_n$  be defined as in (6). Then,*

$$\sqrt{n}(\hat{\theta}_n - \theta_0) = \frac{1}{\sqrt{n}} \sum_{1 \leq i \leq n} \psi^*(X_i, A_i, R_i, \theta_0) + o_P(1) . \quad (9)$$

where

$$\begin{aligned} \psi^*(X_i, A_i, R_i, \theta_0) &= -M^{-1} \left( I\{A_i = 1\} (m(X_i, 1, R_i, \theta_0) - E[m(X_i, 1, R_i(1), \theta_0)|X_i]) \right. \\ &\quad + I\{A_i = 0\} (m(X_i, 0, R_i, \theta_0) - E[m(X_i, 0, R_i(0), \theta_0)|X_i]) \\ &\quad \left. + \eta E[m(X_i, 1, R_i(1), \theta_0)|X_i] + (1 - \eta) E[m(X_i, 0, R_i(0), \theta_0)|X_i] \right) , \end{aligned}$$

and  $M = \frac{\partial}{\partial \theta} E[m(X, A, R, \theta)] \Big|_{\theta=\theta_0}$ . Further, we have that

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{d} N(0, \mathbb{V}_*) , \quad (10)$$

where

$$\mathbb{V}_* = \text{Var}[\psi^*(X_i, A_i, R_i, \theta_0)] . \quad (11)$$

In order to make Theorem 3.1 useful for inference about  $\theta_0$ , we describe in Appendix A an estimator  $\hat{\mathbb{V}}_n$  of  $\mathbb{V}_*$  and sketch a proof of its consistency.

**Remark 3.3.** By comparing the variance expression in (7) to the variance expression for  $\mathbb{V}_*$ , we obtain

$$\begin{aligned} \mathbb{V} - \mathbb{V}_* &= \eta(1 - \eta) M^{-1} E[(E[m(X_i, 1, R_i(1)|X_i] - E[m(X_i, 0, R_i(0)|X_i]) \\ &\quad \times (E[m(X_i, 1, R_i(1)|X_i] - E[m(X_i, 0, R_i(0)|X_i])') (M^{-1})' , \end{aligned}$$

which is positive semidefinite. From this, we conclude that the asymptotic variance of the “naïve” method of moments estimator  $\hat{\theta}_n$  is lower in a finely stratified design compared to i.i.d. assignment. In Section 4, we will further show that  $\mathbb{V}_*$  is the lowest possible asymptotic variance among “regular” estimators for  $\theta_0$  in a large class of treatment assignment schemes, including both i.i.d. assignment and finely stratified designs. When

$d_\theta = 1$ , we may express  $\mathbb{V} - \mathbb{V}_*$  in terms of the “nonparametric  $R^2$ .” In particular,  $\mathbb{V} - \mathbb{V}_*$  is proportional to  $E[R^2(g_i, X_i) \text{Var}[g_i]]$ , where

$$R^2(g_i, X_i) = \frac{\text{Var}[E[g_i|X_i]]}{\text{Var}[g_i]}, \quad (12)$$

and  $g_i = m(X_i, 1, R_i(1), \theta_0) - m(X_i, 0, R_i(0), \theta_0)$ . The quantity in (12) measures how much of the variation in  $g_i$  can be explained nonparametrically by  $X_i$ . See, for instance, Chernozhukov et al. (2024). ■

**Remark 3.4.** Note it follows from (3) that

$$\eta E_Q[m(X_i, 1, R_i(1), \theta_0)] + (1 - \eta) E_Q[m(X_i, 0, R_i(0), \theta_0)] = E_P[m(X_i, A_i, R_i, \theta_0)] = 0,$$

so that  $E[\psi^*(X_i, A_i, R_i, \theta_0)] = 0$ . It is further straightforward to show using Assumption 2.1 that

$$\begin{aligned} \mathbb{V}_* &= \text{Var}[\psi^*(X_i, A_i, R_i, \theta_0)] \\ &= M^{-1} \left( E[\eta \text{Var}[m(X_i, 1, R_i(1), \theta_0)|X_i] + (1 - \eta) \text{Var}[m(X_i, 0, R_i(0), \theta_0)|X_i]] \right. \\ &\quad \left. + \text{Var}[\eta E[m(X_i, 1, R_i(1), \theta_0)|X_i] + (1 - \eta) E[m(X_i, 0, R_i(0), \theta_0)|X_i]] \right) (M^{-1})' \end{aligned} \quad (13)$$

For instance, in the special case of the ATE (Example 2.1) we obtain that

$$\begin{aligned} \text{Var}[\psi^*(X_i, A_i, R_i, \theta_0)] &= E \left[ \frac{\text{Var}[Y_i(1)|X_i]}{\eta} + \frac{\text{Var}[Y_i(0)|X_i]}{1 - \eta} \right. \\ &\quad \left. + (E[Y_i(1) - Y_i(0)|X_i] - E[Y_i(1) - Y_i(0)])^2 \right], \end{aligned} \quad (14)$$

which matches the asymptotic variance derived in Bai et al. (2022) for matched pairs. Theorem 3.1 however accommodates a much larger class of parameters including those introduced in Examples 2.2–2.5. ■

**Remark 3.5.** Although Theorem 3.1 is focused on the case where  $\eta(X_i) = \eta$  is a constant, straightforward modifications of the treatment assignment mechanism described in Assumptions 3.1–3.2 can be applied in more general settings. For instance, suppose  $\eta(X_i)$  takes on a finite set of values  $\{\eta_1, \dots, \eta_S\}$ , we could then simply implement a finely stratified experiment over each set  $\{i : \eta(X_i) = \eta_s\}$  for  $1 \leq s \leq S$ . In other words, separately within each stratum defined by the units for which  $\eta(X_i) = \eta_s$ , employ the assignment mechanism described in Assumptions 3.1–3.2 with  $\ell/k = \eta_s$ . For more general functions  $\eta(X_i)$ , we conjecture that we could employ the local randomization procedure proposed in Cytrynbaum (2023b). ■

## 4 Efficiency Bound

An inspection of the asymptotic variance in (11) reveals that  $\mathbb{V}_*$  in fact coincides with the classical efficiency bound for estimating  $\theta_0$  with i.i.d. assignment. For example, the variance of the difference-in-means estimator in (14) equals the efficiency bound derived in Hahn (1998) for estimating the ATE with a known marginal treatment probability  $\eta$ . Therefore, another way to interpret our result in Theorem 3.1 is that the standard i.i.d. efficiency bound can be attained by a naïve method of moments estimator under a finely stratified design. On the other hand, because treatment status is not independent in a finely stratified design, a natural follow-

up question is whether or not the efficiency bound for estimating  $\theta_0$  changes relative to what can be obtained under i.i.d. assignment once we allow for more general assignment mechanisms. In this section, we show that  $\mathbb{V}_*$  continues to be the efficiency bound for the class of parameters introduced in Section 2 under a more general class of treatment assignment mechanisms. The main restriction on treatment assignment is given by Assumption 2.1, which requires the marginal treatment probability to be known and equal to  $\eta$ . As mentioned earlier and explained in Remark 4.3 below, it is possible to relax this requirement so that  $\eta$  can be replaced by a known function  $\eta(X_i)$ . For a discussion of how our efficiency bound compares with other results in the literature, see Remark 4.4.

We impose the following high-level assumption on the assignment mechanism:

**Assumption 4.1.** The treatment assignment mechanism is such that for any integrable Lipschitz functions  $\gamma_0, \gamma_1 : \mathbf{R}^{d_x} \rightarrow \mathbf{R}$ ,

$$\frac{1}{n} \sum_{1 \leq i \leq n} \sum_{a \in \{0,1\}} I\{A_i = a\} \gamma_a(X_i) \xrightarrow{P} \eta E[\gamma_1(X_i)] + (1 - \eta) E[\gamma_0(X_i)] .$$

In other words, Assumption 4.1 requires that the assignment mechanism admits a law of large numbers for “well-behaved” functions of the covariate values. Examples 4.1–4.2 illustrate that the assumption holds for a large class of treatment assignment mechanisms used in practice.

**Example 4.1** (Covariate-adaptive randomization). Let  $S : \mathbf{R}^{d_x} \rightarrow \mathcal{S} = \{1, \dots, |\mathcal{S}|\}$  be a function that maps the covariates into a set of discrete “strata.” Assume that treatment status is assigned so that  $(R^{(n)}(1), R^{(n)}(0), X^{(n)}) \perp\!\!\!\perp A^{(n)} | \mathcal{S}^{(n)}$ , and that for  $s \in \mathcal{S}$ ,

$$\frac{\sum_{1 \leq i \leq n} I\{S_i = s, A_i = 1\}}{\sum_{1 \leq i \leq n} I\{S_i = s\}} \xrightarrow{P} \eta .$$

This high-level assumption accommodates stratified assignment mechanisms commonly used in empirical practice (see, for instance, Duflo et al., 2015; Dizon-Ross, 2019). It follows from Lemma C.4 in Bugni et al. (2019) that for any integrable functions  $\gamma_0, \gamma_1$ ,

$$\frac{1}{n} \sum_{1 \leq i \leq n} \sum_{a \in \{0,1\}} I\{A_i = a\} \gamma_a(X_i) \xrightarrow{P} \sum_{s \in \mathcal{S}} P\{S_i = s\} (\eta E[\gamma_1(X_i) | S_i = s] + (1 - \eta) E[\gamma_0(X_i) | S_i = s]) .$$

Therefore, Assumption 4.1 is satisfied. ■

**Example 4.2** (Matched pairs). Suppose  $n$  is even and we assign treatment using a finely stratified design with  $(\ell, k) = (1, 2)$ . As discussed at the beginning of Section 3, such a design is also known as a matched pairs design. Assume that the pairing algorithm  $\pi_n(X^{(n)})$  results in pairs that are “close” in the sense of Assumption 3.2. It then follows from the proof of Lemma S.1.5 of Bai et al. (2022) that

$$\frac{1}{n} \sum_{1 \leq i \leq n} \sum_{a \in \{0,1\}} I\{A_i = a\} \gamma_a(X_i) \xrightarrow{P} \frac{1}{2} E[\gamma_0(X_i)] + \frac{1}{2} E[\gamma_1(X_i)] .$$

Therefore, Assumption 4.1 is satisfied. ■

Next, we impose the following high-level assumption on the distributions  $Q$  and  $P$ :

**Assumption 4.2.** The distributions  $Q$  and  $P$  are such that

- (a)  $\text{Var}[m(X_i, a, R_i(a), \theta_0) | X_i = x]$  is a Lipschitz function.
- (b)  $\theta_0$  is uniquely determined by (3) and  $M = \frac{\partial}{\partial \theta} E[m(X_i, A_i, R_i, \theta)] \Big|_{\theta=\theta_0}$  is invertible.

Assumption 4.2(a) is a smoothness condition that is required in settings where  $X_i$  is continuous to ensure that the function  $\psi^*(\cdot)$  we derive in Theorem 4.1 below is in fact the efficient influence function. Note that, if we strengthen Assumption 4.1 to hold for all integrable functions  $\gamma_0, \gamma_1$  then Assumption 4.2(a) could be dropped, with no change to the resulting bound derived below. However, our argument for attaining the bound using finely stratified designs in Section 3 (see, in particular, Remark 3.2) requires similar smoothness conditions, and so we maintain it here as well. Assumption 4.2(b) is a standard assumption used when deriving the properties of  $Z$ -estimators and repeats Assumption 3.4(a).

We now present an efficiency bound for the parameter  $\theta_0$  introduced in Section 2. Formally, we characterize the bound via a convolution theorem that applies to all “regular” estimators of the parameter  $\theta_0$ , where “regular” here should be understood in the standard sense necessary to rule out, for instance, super-efficient estimators (see, for instance, Example 8.1 in van der Vaart, 1998). In stating our theorem we leave the precise definition of “regular” and related assumptions to Appendix B.1. In the paragraph following the statement of the theorem we provide some more details on the nature of our result.

**Theorem 4.1.** *Suppose Assumptions 2.1 and 4.1–4.2 hold, and maintain the additional regularity conditions (20), (21) and Assumption B.1 described in Appendix B.1. Let  $\tilde{\theta}_n$  be any “regular” estimator of the parameter  $\theta_0$  in the sense of (24) in Appendix B.1. Then,*

$$\sqrt{n}(\tilde{\theta}_n - \theta_0) \xrightarrow{d} L ,$$

where

$$L = N(0, \mathbb{V}_*) * B ,$$

for  $\mathbb{V}_*$  in (11) and some fixed probability measure  $B$  which is specific to the estimator  $\tilde{\theta}_n$ .

Given Theorem 4.1 we call  $\mathbb{V}_* = \text{Var}[\psi^*(X_i, A_i, R_i, \theta_0)]$  the efficiency bound for  $\theta_0$ , since our result shows that this is the lowest asymptotic variance attainable by any regular estimator under our assumptions. We note that our assumptions on the assignment mechanism preclude us from applying results based on “standard” arguments (see, for instance, van der Vaart, 1998). Specifically, if we define a tangent set as the collection of score functions of “smooth” one-dimensional parametric sub-models in an appropriate sense, then we are not able to guarantee that the resulting tangent set is linear (or even a convex cone) while *simultaneously* verifying that the likelihood ratio process is locally asymptotically normal for arbitrary assignment mechanisms which satisfy Assumption 4.1. Instead, we proceed by justifying an application of Corollary 3.1 in Armstrong (2022) combined with the convolution Theorem 3.11.2 in van der Vaart and Wellner (1996) to each  $d_\theta$ -dimensional parametric submodel separately, and then arguing that the supremum over all such submodels is attained by  $\text{Var}[\psi^*]$  under Assumption 4.2.

**Remark 4.1.** Following similar arguments as those in Remark 3.4, we can deduce that our efficiency bound recovers well-known bounds for common parameters (like those presented in Examples 2.1–2.3) in the setting of i.i.d. assignment. For example, we have noted in the case of the ATE (Example 2.1) that (14) matches the efficiency bound under i.i.d. assignment derived in Hahn (1998). See Rafi (2023) and Armstrong (2022) for related results in the context of stratified and adaptive experiments. Straightforward calculation also implies that, for the quantile treatment effect (Example 2.2), the efficiency bound is given by

$$E \left[ \frac{1}{\eta} \frac{F_1(\theta_0(1)|X_i)(1 - F_1(\theta_0(1)|X_i))}{f_1(\theta_0(1))^2} + \frac{1}{1 - \eta} \frac{F_0(\theta_0(0)|X_i)(1 - F_0(\theta_0(0)|X_i))}{f_0(\theta_0(0))^2} + \left( \frac{F_1(\theta_0(1)|X_i) - \tau}{f_1(\theta_0(1))} - \frac{F_0(\theta_0(0)|X_i) - \tau}{f_0(\theta_0(0))} \right)^2 \right],$$

which matches the efficiency bound under i.i.d. assignment derived in Firpo (2007) when the propensity score is set to  $\eta$ . ■

**Remark 4.2.** The efficiency bound in Theorem 4.1 is attained by finely stratified experiments as in Theorem 3.1 if no additional covariates are available for estimation beyond the set of covariates  $X_i$  used in the design. In practice, researchers may consider adjusting for additional baseline covariates in order to improve efficiency. Suppose additional covariates  $W^{(n)}$  are available, Assumption 3.1 is modified such that

$$(R^{(n)}(1), R^{(n)}(0), W^{(n)}) \perp\!\!\!\perp A^{(n)} | X^{(n)},$$

and Assumption 4.2(a) is modified such that  $\text{Var}[m(X_i, a, R_i(a), \theta_0) | X_i = x, W_i = w]$  is Lipschitz. When  $d_\theta = 1$ , it can be shown that the efficiency bound, allowing for additional covariate adjustment based on  $X_i$  and  $W_i$ , is

$$\mathbb{V}_* - \eta(1 - \eta)M^{-1}E[\text{Var}[E[g_i | X_i, W_i] | X_i]](M^{-1})', \quad (15)$$

where  $g_i = m(X_i, 1, R_i(1), \theta_0) - m(X_i, 0, R_i(0), \theta_0)$ . Then, as in Remark 3.3, the gain in efficiency obtained by adjusting for  $X_i$  and  $W_i$  in equation (15) is proportional to

$$E[R^2(g_i, X_i, W_i) \text{Var}[g_i | X_i]],$$

where

$$R^2(g_i, X_i, W_i) = \frac{\text{Var}[E[g_i | X_i, W_i] | X_i]}{\text{Var}[g_i | X_i]}$$

is the nonparametric  $R^2$  from regressing  $g_i$  on  $X_i$  and  $W_i$ , when matching on the variables  $X_i$ . As a result, the scope for improving efficiency by adjusting for additional covariates is limited if  $R^2(g_i, X_i, W_i)$  is small. In the case of estimating the ATE,

$$g_i = \frac{Y_i(1)}{\eta} + \frac{Y_i(0)}{1 - \eta},$$

so the scope for improvement depends on how much additional variation in the weighted potential outcomes can be explained by  $W_i$  beyond  $X_i$ . If researchers select matching variables for which they believe  $R^2(g_i, X_i, W_i)$  to be small, then the additional gain from adjusting for the remaining covariates will neces-

sarily be limited. ■

**Remark 4.3.** Although we focus on the case where  $\eta_i(X_i) = P\{A_i = 1|X_i\} = \eta$  is a constant, the proof of Theorem 4.1 holds when  $\eta_i(x) = \eta(x)$  for  $1 \leq i \leq n$ , where  $\eta(x)$  is an arbitrary known and fixed function. In these settings, Lemma B.4 shows that the efficiency bound equals

$$\begin{aligned} \mathbb{V}_* &= \text{Var}[\psi^*(X_i, A_i, R_i, \theta_0)] \\ &= M^{-1} \left( E[\eta(X_i) \text{Var}[m(X_i, 1, R_i(1), \theta_0)|X_i] + (1 - \eta(X_i)) \text{Var}[m(X_i, 0, R_i(0), \theta_0)|X_i]] \right. \\ &\quad \left. + \text{Var}[\eta(X_i)E[m(X_i, 1, R_i(1), \theta_0)|X_i] + (1 - \eta(X_i))E[m(X_i, 0, R_i(0), \theta_0)|X_i]] \right) (M^{-1})', \end{aligned} \quad (16)$$

so that the only difference from (13) is that  $\eta$  is replaced by  $\eta(X_i)$ . Consider Example 2.1 and note the moment condition for the ATE is now given by

$$m(X_i, A_i, R_i, \theta) = \frac{Y_i A_i}{\eta(X_i)} - \frac{Y_i(1 - A_i)}{1 - \eta(X_i)} - \theta. \quad (17)$$

Straightforward calculation implies that in this example, the efficiency bound in (16) becomes

$$E \left[ \frac{\text{Var}[Y_i(1)|X_i]}{\eta(X_i)} + \frac{\text{Var}[Y_i(0)|X_i]}{1 - \eta(X_i)} + (E[Y_i(1) - Y_i(0)|X_i] - E[Y_i(1) - Y_i(0)])^2 \right], \quad (18)$$

which again matches the efficiency bound under i.i.d. assignment in Hahn (1998). If we additionally impose that  $\eta(X_i) = \eta(S(X_i))$  for  $S$  taking on finitely many values as in Example 4.1, then the bound could be achieved by employing the modified design described in Remark 3.5. ■

**Remark 4.4.** Here, we comment on how Theorem 4.1 relates to prior efficiency bounds in experiments with general assignment mechanisms. For the case of estimating the ATE, Armstrong (2022) derives an efficiency bound over a very large class of assignment mechanisms, including even response-adaptive designs, and shows that the bound is attained when units are assigned to treatment (control) with conditional probability proportional to the conditional variance of the potential outcome under treatment (control). This type of assignment is sometimes referred to as the Neyman allocation. On the other hand, our results show that his bound may be quite loose whenever the assignment proportions are restricted to be anything not equal to the Neyman allocation, which is, of course, unknown. For example, his bound is not informative about what can be achieved if the assignment proportions were set to one half regardless of whether or not the conditional outcome variances across treatment and control are equal. Such settings frequently arise in practice due to logistical constraints or the absence of pilot data with which to estimate conditional variances of potential outcomes under treatment and control. Furthermore, as argued in Cai and Rafi (2022), even if pilot data is available, these quantities may be estimated so poorly that exogenously constraining the assignment proportions to one half leads to more efficient estimates of the ATE in practice. Motivated by such concerns, Rafi (2023) derives an efficiency bound for the ATE over the class of “coarsely-stratified” assignment mechanisms studied in Bugni et al. (2019), where the stratum-level assignment proportions are restricted *a priori* by the experimenter. This framework, however, rules out finely stratified designs. Finally, we once again emphasize that our analysis, unlike these other papers, applies to a general class of treatment effect parameters, including the ATE as a special case. ■



## 5 Simulations

In this section, we illustrate the results in Sections 3 and 4 with a simulation study. Specifically, we set  $\eta = 1/2$ , and compare the mean-squared errors (MSE) obtained from the “naive” estimator  $\hat{\theta}_n$  and various adjusted estimators, for i.i.d. treatment assignment versus matched pairs assignment (see Remark 3.1 and Example 4.2). In Section 5.1, we present the model specifications and estimators for estimating the ATE as in Example 2.1. In Section 5.2, we present the model specifications and estimators for estimating the LATE as in Example 2.3. Section 5.3 reports the simulation results.

### 5.1 Average Treatment Effect

In this section, we present model specifications and estimators for estimating the ATE as in Example 2.1. Recall that in this case the moment function we consider is given by

$$m(X_i, R_i, A_i, \theta) = \frac{Y_i A_i}{\eta} - \frac{Y_i(1 - A_i)}{1 - \eta} - \theta ,$$

with  $R_i = Y_i$ . For  $a \in \{0, 1\}$  and  $1 \leq i \leq n$ , the potential outcomes are generated according to the equation:

$$Y_i(a) = \mu_a(X_i) + \sigma_a(X_i)\epsilon_i .$$

In each of the specifications,  $((X_i, \epsilon_i) : 1 \leq i \leq n)$  are i.i.d; for  $1 \leq i \leq n$ ,  $X_i$  and  $\epsilon_i$  are independent.

**Model 1:**  $\mu_0(X_i) = X_i + (X_i^2 - 1)/3$ ,  $\mu_1(X_i) = 0.2 + \mu_0(X_i)$ ,  $\epsilon_i \sim N(0, 1)$ ,  $X_i \sim N(0, 1)$  and  $\sigma_a(X_i) = 2$ .

**Model 2:** As in Model 1, but  $\mu_a(X_i) = 0.2I\{a = 1\} + \gamma_a(\sin(X_i) + X_i) + (X_i^2 - 1)/3$  where  $\gamma_1 = 1$  and  $\gamma_0 = -1$ , and  $\sigma_a(X_i) = (1 + a)X_i^2$ .

**Model 3:** As in Model 2, but  $\mu_1(X_i) = 0.2 + 3(X_i^2 - 1)$  and  $\mu_0(X_i) = 0$ .

We consider the following three estimators for the ATE:

**Unadjusted Estimator:**

$$\hat{\theta}_n^{\text{unadj}} = \frac{1}{n/2} \sum_{1 \leq i \leq n} (Y_i A_i - Y_i(1 - A_i)) .$$

**Adjusted Estimator 1:**

$$\hat{\theta}_n^{\text{adj},1} = \frac{1}{n} \sum_{1 \leq i \leq n} (2A_i(Y_i - \hat{\mu}_1^Y(X_i)) - 2(1 - A_i)(Y_i - \hat{\mu}_0^Y(X_i)) + \hat{\mu}_1^Y(X_i) - \hat{\mu}_0^Y(X_i)) ,$$

where  $\hat{\mu}_a^Y(X_i)$  is constructed by running a least squares regression of  $Y_i$  on  $(1, X_i, X_i^2)$  using the sample from  $A_i = a$ .

**Adjusted Estimator 2:**

$$\hat{\theta}_n^{\text{adj},2} = \frac{1}{n} \sum_{1 \leq i \leq n} (2A_i(Y_i - \hat{\mu}_1^Y(X_i)) - 2(1 - A_i)(Y_i - \hat{\mu}_0^Y(X_i)) + \hat{\mu}_1^Y(X_i) - \hat{\mu}_0^Y(X_i)) ,$$

where  $\hat{\mu}_a^Y(X_i)$  is constructed by running a least squares regression of  $Y_i$  on  $(1, X_i, X_i^2, X_i 1\{X_i > t\})$  where  $t$  is the sample median using the sample from  $A_i = a$ .

The first estimator  $\hat{\theta}_n^{\text{unadj}}$  is the method of moments estimator given by the solution to (6). The second and third estimators  $\hat{\theta}_n^{\text{adj},1}$  and  $\hat{\theta}_n^{\text{adj},2}$  are covariate-adjusted estimators which can be obtained as two-step method of moments estimators from solving the ‘‘augmented’’ moment equation (8) described in the discussion at the end of Section 2.  $\hat{\theta}_n^{\text{adj},1}$  and  $\hat{\theta}_n^{\text{adj},2}$  differ in the choice of basis functions used in the construction of the estimators  $\hat{\mu}_a(x)$ . Note that by the double-robustness property of the augmented estimating equation (8), it can be shown that the adjusted estimators  $\hat{\theta}_n^{\text{adj},1}$ ,  $\hat{\theta}_n^{\text{adj},2}$  are consistent and asymptotically normal regardless of the choice of estimators  $\hat{\mu}_a(x)$ , but consistency of  $\hat{\mu}_a(x)$  to  $\mu_a(x)$  would ensure that  $\hat{\theta}_n^{\text{adj},1}$ ,  $\hat{\theta}_n^{\text{adj},2}$  are efficient under i.i.d. assignment (Robins et al., 1995; Tsiatis et al., 2008; Chernozhukov et al., 2017).

## 5.2 Local Average Treatment Effect

In this section, we present the model specifications and estimators for estimating the LATE as in Example 2.3. Recall that in this case the moment condition we consider is given by

$$m(X_i, A_i, R_i, \theta) = \frac{Y_i A_i}{\eta} - \frac{Y_i(1 - A_i)}{1 - \eta} - \theta \left( \frac{D_i A_i}{\eta} - \frac{D_i(1 - A_i)}{1 - \eta} \right) ,$$

with  $R_i = (Y_i, D_i)$ . The outcome is determined by the relationship  $Y_i = D_i Y_i(1) + (1 - D_i) Y_i(0)$ , where  $Y_i(d) = \mu_d(X_i) + \sigma_a(X_i) \epsilon_i$  follows the same outcome model as in the ATE setup of Section 5.1. In addition, we have  $D_i = A_i D_i(1) + (1 - A_i) D_i(0)$ , where

$$D_i(0) = I \{ \alpha_0 + \alpha(X_i) > \varepsilon_{1,i} \} ,$$

$$D_i(1) = \begin{cases} I \{ \alpha_1 + \alpha(X_i) > \varepsilon_{2,i} \} & \text{if } D_i(0) = 0 \\ 1 & \text{otherwise} \end{cases} .$$

For each outcome model, we set  $\alpha_0 = 0.5$ ,  $\alpha_1 = 1$ ,  $\alpha(X_i) = X_i + (X_i^2 - 1)/3$  and  $\varepsilon_{1,i}, \varepsilon_{2,i} \sim N(0, 4)$ .

We consider the following three estimators for the LATE:

**Unadjusted Estimator:**

$$\hat{\theta}_n^{\text{unadj}} = \frac{\sum_{1 \leq i \leq n} (Y_i A_i - Y_i(1 - A_i))}{\sum_{1 \leq i \leq n} (D_i A_i - D_i(1 - A_i))} .$$

**Adjusted Estimator 1:**

$$\hat{\theta}_n^{\text{adj},1} = \frac{\sum_{1 \leq i \leq n} (2A_i(Y_i - \hat{\mu}_1^Y(X_i)) - 2(1 - A_i)(Y_i - \hat{\mu}_0^Y(X_i)) + \hat{\mu}_1^Y(X_i) - \hat{\mu}_0^Y(X_i))}{\sum_{1 \leq i \leq n} (2A_i(D_i - \hat{\mu}_1^D(X_i)) - 2(1 - A_i)(D_i - \hat{\mu}_0^D(X_i)) + \hat{\mu}_1^D(X_i) - \hat{\mu}_0^D(X_i))} ,$$

where  $\hat{\mu}_a^Y(X_i)$  is estimated by running a least squares regression of  $Y_i$  on  $(1, X_i, X_i^2)$  using the sample from  $A_i = a$ , and  $\hat{\mu}_a^D(X_i)$  is estimated using logistic regressions using the same set of regressors using the sample from  $A_i = a$ .

**Adjusted Estimator 2:** As in Adjusted Estimator 1, but  $\hat{\mu}_a^Y(X_i)$  and  $\hat{\mu}_a^D(X_i)$  are estimated respectively by running a least squares and logistic regressions of  $Y_i$  on  $(1, X_i, X_i^2, X_i 1\{X_i > t\})$  where  $t$  is the sample median.

Similarly to Section 5.1,  $\hat{\theta}_n^{\text{unadj}}$  solves (6) for the moment condition given in (5). The second and third estimators are covariate adjusted estimators which can be obtained as two-step method of moments estimators from solving an “augmented” version of the moment condition (5) (see, for instance, Chernozhukov et al., 2018; Jiang et al., 2022a).

### 5.3 Simulation Results

Table 5.3 displays the ratio of the MSE for each design/estimator pair relative to the MSE of the unadjusted estimator under i.i.d. assignment, computed across 2000 Monte Carlo replications. As expected given our theoretical results, we find that the empirical MSEs of the naive unadjusted estimator under a matched pairs design closely match the empirical MSEs of the covariate adjusted estimators under i.i.d. assignment.

## 6 Recommendations for Empirical Practice

We conclude with some recommendations for empirical practice based on our theoretical results. Overall, our findings highlight the general benefit of fine stratification for designing efficient experiments: finely stratified experiments “automatically” perform fully-efficient regression adjustment for a large class of interesting parameters. This generalizes similar observations made by Bai et al. (2022), Bai (2022) and Cytrynbaum (2023b) for the special case of estimating the ATE.

One caveat to this result, however, is that it crucially hinges on the assumption that units within each block are sufficiently “close” (Assumption 3.2), and such a condition becomes difficult to satisfy as the dimension of  $X_i$  increases. For this reason, we recommend that practitioners construct their blocks using a small subset of the baseline covariates that they believe have the highest explanatory power in terms of the nonparametric  $R^2$  in (12); the baseline level of the experimental outcomes, for example, is typically believed to be one such covariate (see, in particular, Bruhn and McKenzie, 2009). The experimental data can then be analyzed efficiently using an unadjusted method-of-moments estimator.

If one wishes to perform regression adjustment with additional covariates beyond those used for blocking, then this can be done *ex post*. As discussed in Remark 4.2, the scope for improvement from covariate adjustment is limited by the nonparametric  $R^2$  from the regression of the moment functions on the additional covariates, conditional on the ones used for matching; if one has already matched on the covariates with the highest explanatory power, then the potential gain in efficiency from adjusting for these additional covariates

Table 1: MSE ratios relative to unadjusted estimator under i.i.d. assignment

		I.I.D. assignment			Matched pairs	
		Model	Unadjusted	Adjusted 1	Adjusted 2	Unadjusted
$n = 200$	ATE	1	1.0000	0.4580	0.4637	0.4530
		2	1.0000	0.9836	1.0090	1.0291
		3	1.0000	0.7473	0.7615	0.7415
	LATE	1	1.0000	0.4391	0.4175	0.4115
		2	1.0000	0.8967	0.9724	0.8813
		3	1.0000	0.5191	0.5002	0.4936
$n = 400$	ATE	1	1.0000	0.4616	0.4640	0.4471
		2	1.0000	0.9778	1.0470	1.0042
		3	1.0000	0.7535	0.7364	0.7293
	LATE	1	1.0000	0.4267	0.4583	0.4327
		2	1.0000	0.8754	0.9570	0.9375
		3	1.0000	0.5671	0.5349	0.5240
$n = 1000$	ATE	1	1.0000	0.4518	0.4568	0.4453
		2	1.0000	0.9874	0.9730	0.9186
		3	1.0000	0.7374	0.7099	0.7018
	LATE	1	1.0000	0.4437	0.4301	0.4259
		2	1.0000	0.8415	0.8735	0.8618
		3	1.0000	0.5408	0.4879	0.4884
$n = 2000$	ATE	1	1.0000	0.4677	0.4835	0.4867
		2	1.0000	0.9780	0.9163	0.9105
		3	1.0000	0.7272	0.7552	0.7611
	LATE	1	1.0000	0.4702	0.4473	0.4871
		2	1.0000	0.8692	0.8545	0.8202
		3	1.0000	0.4973	0.5044	0.5079

Note: For each model, the MSE of the unadjusted estimator under i.i.d. assignment are normalized to one and the other columns contain the ratios of the MSEs against that of the unadjusted estimator under i.i.d. assignment. MSEs are calculated across 2000 replications.

may be limited. We further caution that care must be taken to ensure that the adjustment is performed in such a way that it guarantees a gain in efficiency: see [Bai et al. \(2024b\)](#) and [Cytrynbaum \(2023a\)](#) for related discussions. Recent work has developed such methods of covariate adjustment for specific parameters of interest (see, for instance, [Bai et al., 2024a,b,c](#); [Cytrynbaum, 2023a](#)), but we leave the development of a method of covariate adjustment which applies at the level of generality considered in this paper to future work.

## Supplemental Appendix (For Online Publication)

### A Construction of the Variance Estimator

For notational convenience we focus on the leading case in which  $d_\theta = 1$ . Note that a similar construction is valid when  $d_\theta > 1$  simply by replacing all quantities with their matrices counterparts. First note that in certain examples (including Examples 2.1 and 2.3–2.5), the analog principle suggests that a natural estimator for  $M$  is given by

$$\widehat{M}_n = \frac{1}{n} \sum_{1 \leq i \leq n} \left. \frac{\partial}{\partial \theta} m(X_i, A_i, R_i, \theta) \right|_{\theta = \hat{\theta}_n}.$$

Under suitable conditions, it follows directly from following arguments in each of the papers mentioned above that  $\widehat{M}_n \xrightarrow{P} M$ .<sup>4</sup> Therefore, it suffices to construct a consistent estimator for the “meat” in (13). By the law of total variance, this middle component equals  $\Sigma_1 + \Sigma_2$ , where

$$\begin{aligned} \Sigma_1 &= \eta \text{Var}[m(X_i, 1, R_i(1), \theta_0)] + (1 - \eta) \text{Var}[m(X_i, 0, R_i(0), \theta_0)] \\ \Sigma_2 &= -\eta(1 - \eta) E\left[ \left( E[m(X_i, 1, R_i(1), \theta_0) | X_i] - E[m(X_i, 1, R_i(1), \theta_0)] \right. \right. \\ &\quad \left. \left. - \left( E[m(X_i, 0, R_i(0), \theta_0) | X_i] - E[m(X_i, 0, R_i(0), \theta_0)] \right) \right)^2 \right] \\ &= -\eta(1 - \eta) \left( E[E[m(X_i, 1, R_i(1), \theta_0) | X_i]^2] + E[E[m(X_i, 0, R_i(0), \theta_0) | X_i]^2] \right. \\ &\quad \left. - 2E[E[m(X_i, 1, R_i(1), \theta_0) | X_i] E[m(X_i, 0, R_i(0), \theta_0) | X_i]] \right. \\ &\quad \left. - (E[m(X_i, 1, R_i(1), \theta_0)] - E[m(X_i, 0, R_i(0), \theta_0)])^2 \right) \end{aligned}$$

For  $a \in \{0, 1\}$ , define

$$\hat{\mu}_n(a) = \frac{1}{\eta_a n} \sum_{1 \leq i \leq n} I\{A_i = a\} m(X_i, A_i, R_i, \hat{\theta}_n),$$

where  $\eta_1 = \eta$  and  $\eta_0 = 1 - \eta$ . The analog principle again suggests that a natural estimator for  $\Sigma_1$  is

$$\hat{\Sigma}_{1,n} = \frac{1}{n} \sum_{1 \leq i \leq n} I\{A_i = 1\} (m(X_i, A_i, R_i, \hat{\theta}_n) - \hat{\mu}_n(1))^2 + \frac{1}{n} \sum_{1 \leq i \leq n} I\{A_i = 0\} (m(X_i, A_i, R_i, \hat{\theta}_n) - \hat{\mu}_n(0))^2.$$

To estimate  $\Sigma_2$ , we first define

$$\hat{\varsigma}_n(0, 1) = \frac{k}{n} \sum_{1 \leq j \leq n/k} \frac{1}{\ell(k - \ell)} \sum_{i, i' \in \lambda_j: A_i = 1, A_{i'} = 0} m(X_i, A_i, R_i, \hat{\theta}_n) m(X_{i'}, A_{i'}, R_{i'}, \hat{\theta}_n).$$

Next, define

$$\hat{\varsigma}_n(1, 1) = \begin{cases} \frac{k}{n} \sum_{1 \leq j \leq n/k} \frac{1}{\binom{\ell}{2}} \sum_{i < i' \in \lambda_j: A_i = A_{i'} = 1} m(X_i, A_i, R_i, \hat{\theta}_n) m(X_{i'}, A_{i'}, R_{i'}, \hat{\theta}_n) & \text{if } \ell > 1 \\ \frac{2k}{n} \sum_{1 \leq j \leq \frac{n}{2k}} \sum_{i \in \lambda_{2j}, i' \in \lambda_{2j-1}: A_i = A_{i'} = 1} m(X_i, A_i, R_i, \hat{\theta}_n) m(X_{i'}, A_{i'}, R_{i'}, \hat{\theta}_n) & \text{if } \ell = 1. \end{cases}$$

<sup>4</sup>In examples including Example 2.2 where  $m$  is nonsmooth in  $\theta$ ,  $M$  may consist of components which require nonparametric estimators, and in such cases bootstrap procedures may be preferable. See, for instance, Jiang et al. (2021).

Similarly, define

$$\hat{\zeta}_n(0, 0) = \begin{cases} \frac{k}{n} \sum_{1 \leq j \leq n/k} \frac{1}{\binom{k-\ell}{2}} \sum_{i < i' \in \lambda_j: A_i = A_{i'} = 0} m(X_i, A_i, R_i, \hat{\theta}_n) m(X_{i'}, A_{i'}, R_{i'}, \hat{\theta}_n) & \text{if } k - \ell > 1 \\ \frac{2k}{n} \sum_{1 \leq j \leq \frac{n}{2k}} \sum_{i \in \lambda_{2j}, i' \in \lambda_{2j-1}: A_i = A_{i'} = 0} m(X_i, A_i, R_i, \hat{\theta}_n) m(X_{i'}, A_{i'}, R_{i'}, \hat{\theta}_n) & \text{if } k - \ell = 1. \end{cases}$$

Finally, define

$$\hat{\Sigma}_{2,n} = -\eta(1 - \eta)(\hat{\zeta}_n(1, 1) + \hat{\zeta}_n(0, 0) - 2\hat{\zeta}_n(0, 1) - (\hat{\mu}_n(1) - \hat{\mu}_n(0))^2).$$

The estimator  $\hat{\zeta}_n(1, 1)$  is constructed in one of two ways depending on the number of treated units in each block. If more than one unit in each block is treated, then we take the averages of all pairwise products of the treated units in each block, and average them across all blocks. We call this a “within block” estimator. If instead only one unit in each block is treated, then we take the product of two treated units in *adjacent* blocks. We call this a “between block” estimator, and note that similar constructions have been used previously in [Abadie and Imbens \(2008\)](#), [Bai et al. \(2022\)](#), [Bai et al. \(2024c\)](#), and [Cytrynbaum \(2023b\)](#). The estimator  $\hat{\zeta}_n(0, 0)$  is constructed similarly. [Bai et al. \(2024d\)](#) compare the finite-sample properties of the “within block” and “between block” variance estimators via simulation. Their findings are that experimental designs which allow for a “within block” variance estimator have better small-sample inferential performance, at the cost of slightly increasing the mean-squared error of the estimator  $\hat{\theta}_n$ , relative to experimental designs which require the use of the “between block” variance estimator. Under suitable assumptions, it follows from similar arguments to those in [Bai \(2022\)](#) and [Bai et al. \(2024d\)](#) that  $\hat{\Sigma}_{1,n} \xrightarrow{P} \Sigma_1$  and  $\hat{\Sigma}_{2,n} \xrightarrow{P} \Sigma_2$ . A natural estimator for  $\mathbb{V}_*$  is therefore given by

$$\hat{\mathbb{V}}_n = \widehat{M}_n^{-2} \left( \hat{\Sigma}_{1,n} + \hat{\Sigma}_{2,n} \right).$$

Thus, provided  $M$  is invertible, we have that  $\hat{\mathbb{V}}_n \xrightarrow{P} \mathbb{V}_*$ .

## B Proofs of Main Results

### B.1 Proofs for Section 4

Recall that  $P_n$  denotes the distribution of the observed data  $(X^{(n)}, A^{(n)}, R^{(n)})$ , and  $Q$  denotes the marginal distribution of the vector  $(R_i(1), R_i(0), X_i)$ . Note that any treatment assignment mechanism satisfying [Assumption 2.1](#) can be represented as a function of  $X^{(n)}$  and some additional randomization device  $U_n \in \mathbf{R}$ . Let  $p_n^{U_n}$  denote the density function for  $U_n$  with respect to a dominating measure  $\mu^U$ . In what follows, we consider a family  $\{Q_t : t \in \mathbb{R}^{d_\theta}\}$  of marginal distributions indexed by  $t$ , and let  $q_t^X$  denote the density function for  $X_i$  with respect to a dominating measure  $\mu^X$ ,  $q_t^{R(a)|X}(r|x)$  denote the conditional density of  $R_i(a)$  given  $X_i$  with respect to a dominating measure  $\mu^R$ . Further let  $P_{t,n}$  denote the distribution of  $Z^{(n)}$  and note it is jointly determined by  $Q_t$  and the distribution of  $U_n$ . We require that  $Q_0 = Q$  and  $P_{0,n} = P_n$

and define  $q^X = q_0^X$  and  $q^{R(a)|X} = q_0^{R(a)|X}$ . As a consequence, the density function of  $P_{t,n}$  is given by

$$\ell_n = p_n^U(U_n) \prod_{1 \leq i \leq n} q_t^X(X_i) \prod_{1 \leq i \leq n} \prod_{a \in \{0,1\}} q_t^{R(a)|X}(R_i|X_i)^{I\{A_i=a\}}. \quad (19)$$

Because the density  $p_n^U$  does not depend on  $t$ , and in general we will only concern ourselves with the ratio of likelihoods at different values of  $t$  (so that  $p_n^U$  in the ratio will cancel), in what follows we suppress the dependence on  $n$  and simply identify the distribution  $P_{t,n}$  with its corresponding marginal distribution  $P_t$ . We consider a parametric submodel  $\{P_t : t \in \mathbf{R}^{d_\theta}\}$ , where  $P_0 = P$ , such that the following holds for  $g = (g^X, g^{R(1)|X}, g^{R(0)|X})$ , each component of which is a  $d_\theta$ -dimensional function:

(a) As  $t \rightarrow 0$ ,

$$\int \frac{1}{\|t\|^2} \left( q_t^X(x)^{1/2} - q^X(x)^{1/2} - \frac{1}{2} q^X(x)^{1/2} t' g^X(x) \right)^2 d\mu^X(x) \rightarrow 0. \quad (20)$$

(b) For  $a \in \{0, 1\}$ ,  $E_Q[g^{R(a)|X}(R(a)|X)g^{R(a)|X}(R(a)|X)'|X = x]$  is Lipschitz and for  $Q$ -almost every  $x$ , as  $t \rightarrow 0$ ,

$$\frac{1}{\|t\|^2} \iint \left( q_t^{R(a)|X}(r|x)^{1/2} - q^{R(a)|X}(r|x)^{1/2} - \frac{1}{2} q^{R(a)|X}(r|x)^{1/2} t' g^{R(a)|X}(r|x) \right)^2 d\mu^R(r) q^X(x) d\mu^X(x) \rightarrow 0. \quad (21)$$

In what follows, we will further index a parametric submodel by its associated function  $g$ , denoted by  $P_{t,g}$ , to emphasize the role of  $g$ . Similarly we denote the density of  $Q_{t,g}$  by  $q_{t,g}$ .

Define the information of  $X$  as  $I^X = E_Q[g^X(X)g^X(X)']$ . Define the conditional information of  $R(a)$  given  $X = x$  as

$$I^{R(a)|X}(x) = E_Q[g^{R(a)|X}(R(a)|X)g^{R(a)|X}(R(a)|X)'|X = x].$$

**Lemma B.1.** *For a parametric submodel  $\{P_{t,g} : t \in \mathbf{R}^{d_\theta}\}$  with  $P_{0,g} = P$  that satisfies (20)–(21),*

(a)  $I^X < \infty$ .

(b)  $E_Q[g^X(X)] = 0$ .

(c)  $E_Q[g^{R(a)|X}(R(a)|X)g^{R(a)|X}(R(a)|X)'] < \infty$  and hence  $I^{R(a)|X}(X) < \infty$  with probability one under  $Q$ .

(d)  $E_Q[g^{R(a)|X}(R(a)|X)|X] = 0$  with probability one under  $Q$ .

PROOF. (a) and (b) follow from Lemma 14.2.1 in [Lehmann and Romano \(2022\)](#). (c) follows from the same lemma. In order to show (d), fix  $t_n \rightarrow 0$ . Note (21) and Markov's inequality imply that along a subsequence  $t_{n_k}$ ,

$$\frac{1}{\|t_{n_k}\|^2} \int \left( q_{t_{n_k}}^{R(a)|X}(r|x)^{1/2} - q^{R(a)|X}(r|x)^{1/2} - \frac{1}{2} q^{R(a)|X}(r|x)^{1/2} t_{n_k}' g^{R(a)|X}(r|x) \right)^2 d\mu^R(r) \rightarrow 0$$

for  $Q$ -almost every  $x$ . Along that subsequence, another application of Lemma 14.2.1 in [Lehmann and Romano \(2022\)](#) implies (d). ■

For  $t \in \mathbf{R}^{d_\theta}$ , the log-likelihood ratio between  $P_{t/\sqrt{n},g}$  and  $P_0 = P$  is

$$L_{n,t}(g) = \frac{1}{n} \sum_{1 \leq i \leq n} \log \frac{q_{t/\sqrt{n},g}^X(X_i)}{q^X(X_i)} + \frac{1}{n} \sum_{1 \leq i \leq n} \sum_{a \in \{0,1\}} I\{A_i = a\} \log \frac{q_{t/\sqrt{n},g}^{R(a)|X}(R_i|X_i)}{q^{R(a)|X}(R_i|X_i)}.$$

The following lemma establishes an expansion of the log-likelihood ratio and local asymptotic normality of  $\{P_{t/\sqrt{n},g}\}$ .

**Lemma B.2.** *Suppose the treatment assignment mechanism satisfies Assumption 2.1 and  $g$  satisfies (20)–(21). Then,*

$$L_{n,t}(g) = \frac{1}{\sqrt{n}} \sum_{1 \leq i \leq n} t' s_g(X_i, A_i, R_i) - \frac{1}{2} t' I^X t - \frac{1}{2n} \sum_{1 \leq i \leq n} \sum_{a \in \{0,1\}} I\{A_i = a\} t' I^{R(a)|X}(X_i) t + o_P(1),$$

where

$$s_g(x, a, r) = g^X(x) + I\{a = 1\} g^{R(1)|X}(r|x) + I\{a = 0\} g^{R(0)|X}(r|x) \quad (22)$$

and  $I = I^X + \eta E_Q[I^{R(1)|X}(X_i)] + (1 - \eta) E_Q[I^{R(0)|X}(X_i)]$ . If in addition the assignment mechanism satisfies Assumption 4.1, then, under  $P_0$ ,

$$L_{n,t}(g) \xrightarrow{d} N\left(-\frac{1}{2} t' I t, t' I t\right),$$

PROOF. The first result follows from Theorem 3.1 of Armstrong (2022). The second result follows from Corollary 3.1 of Armstrong (2022) given Assumption 4.1 and the assumption that  $I^{R(a)|X}(x)$  is Lipschitz. ■

We emphasize that Lemma B.5 implies

$$\sum_{1 \leq i \leq n} s_g(X_i, A_i, R_i)$$

is the sum of  $n$  identically distributed, despite possibly dependent, random variables. Therefore, in what follows, quantities like  $E_P[s_g]$  are well defined.

Let the following condition collect the properties of the functions  $g$  that are of interest to us:

**Condition B.1.** The function  $g$  satisfies that  $E_P[g^X(X)] = 0$ ,  $E_P[g^X(X)g^X(X)'] < \infty$ ,  $E_P[g^{R(a)|X}(R|X)g^{R(a)|X}(R|X)'] < \infty$ ,  $E_P[g^{R(a)|X}(R|X)|X] = 0$  with probability one, and  $E_P[g^{R(a)|X}(R|X)g^{R(a)|X}(R|X)'|X = x]$  is Lipschitz for  $a \in \{0, 1\}$ . In addition,  $I$  is nonsingular.

We note that for any  $g$  that satisfies Condition B.1, there exists a parametric submodel  $\{P_{t,g} : t \in \mathbf{R}^{d_\theta}\}$  such that (20)–(21) hold. Such a construction follows from the construction on p.69 in Tsiatis (2006) and can be done separately for  $g_1(x)$  and  $g_2^a(r|x)$  for each  $x$  separately so that they satisfy (20)–(21).

Let  $\theta(P) \in \mathbf{R}^{d_\theta}$  be a parameter of interest. Further suppose that for each  $g$  satisfying Condition B.1, there exists a  $d_\theta \times 1$  vector of functions  $\psi^* \in L^2(P)$  such that for all  $t \in \mathbf{R}^{d_\theta}$ , as  $n \rightarrow \infty$ ,

$$\sqrt{n}(\theta(P_{t/\sqrt{n},g}) - \theta(P)) \rightarrow E_P[\psi^* s_g' t]. \quad (23)$$



We provide explicit conditions which guarantee this is possible when  $\theta(P)$  is defined by (3), in Lemma B.4 below.

We recall an estimator  $\tilde{\theta}_n$  for  $\theta(P)$  is regular if for all  $g$  and  $t \in \mathbf{R}^{d_\theta}$ ,

$$\sqrt{n}(\tilde{\theta}_n - \theta(P_{t/\sqrt{n},g})) \xrightarrow{P_{t/\sqrt{n},g}} L \quad (24)$$

for a fixed probability measure  $L$ .

The following lemma establishes a convolution theorem for regular estimators:

**Lemma B.3.** *Suppose  $\theta$  satisfies (23). Let  $\tilde{\theta}_n$  be a regular estimator for  $\theta$ . Further suppose that  $\psi^* = s_g$  for some function  $g$  satisfying Condition B.1. Then,*

$$L = N(0, E_P[\psi^* \psi^{*'}]) * B ,$$

where  $B$  is a fixed probability measure.

PROOF. In what follows, for each  $g$  satisfying Condition B.1, we consider the linear subspace given by

$$\mathcal{M}_g = \{t' s_g : t \in \mathbf{R}^{d_\theta}\} .$$

Note that  $t' s_g$  appears in the expansion of the log-likelihood ratio between  $P_{t/\sqrt{n},g}$  and  $P$ . We first derive the Riesz representer along the parametric subspace  $\mathcal{M}_g$ . In particular, for each  $b \in \mathbf{R}^{d_\theta}$ , we solve for  $w(b) \in \mathbf{R}^{d_\theta}$  via the property that,

$$b' E_P[\psi^* s'_g t] = E_P[w(b)' s_g s'_g t]$$

needs to hold for all  $t \in \mathbf{R}^{d_\theta}$  and get

$$w(b) = E[s_g s'_g]^{-1} E[s_g \psi^{*'}] b .$$

Therefore, the Riesz representer is

$$E_P[\psi^* s'_g] E[s_g s'_g]^{-1} s_g .$$

It then follows from the local asymptotic normality established in Lemma B.2 and Theorem 3.11.2 in van der Vaart and Wellner (1996) that

$$L = N(0, V_g) * B_g ,$$

where

$$V_g = E_P[\psi^* s'_g] E_P[s_g s'_g]^{-1} E[s_g \psi^{*'}]$$

and  $B_g$  is a fixed probability measure. Furthermore, by a standard projection argument, in particular the fact that the second moment of  $\psi^* - E_P[\psi^* s'_g] E_P[s_g s'_g]^{-1} s_g$  is positive semi-definite, it can be shown that  $V_g$  is maximized in the matrix sense when  $s_g = \psi^*$ . Note this maximum is attained by our assumption that  $\psi^* = s_g$  for some  $g$  satisfying Condition B.1. The conclusion then follows. ■

To apply Lemma B.3 to the setting in Section 4, we establish the form of  $\psi^*$  in (23) for the parameter  $\theta_0 = \theta(P)$  defined by (3). Define  $\eta(X_i) = P\{A_i = 1|X_i\}$ . Note that

$$0 = E_P[m(X_i, A_i, R_i, \theta(P))] = E_Q[m(X, 1, R(1), \theta(P))\eta(X)] + E_Q[m(X, 0, R(0), \theta(P))(1 - \eta(X))] . \quad (25)$$

**Lemma B.4.** *Suppose the treatment assignment mechanism satisfies Assumptions 2.1 and 4.1. Fix a function  $g$  that satisfies Condition B.1. Suppose (20)–(21) holds. Fix  $t \in \mathbf{R}^{d_\theta}$  and consider a one-dimensional submodel  $\{P_{t/\sqrt{n}, g}\}$  such that*

$$\begin{aligned} E_{Q_{t/\sqrt{n}}} [m(X, a, R(a), \theta(P))^2] &= O(1) \\ E_{Q^X} [E_{Q_{t/\sqrt{n}}^{R(a)|X}} [m(X, a, R(a), \theta(P))^2|X]] &= O(1) \\ E_{Q_{t/\sqrt{n}}^X} [E_{Q^{R(a)|X}} [m(X, a, R(a), \theta(P))^2|X]] &= O(1) \end{aligned} \quad (26)$$

as  $n \rightarrow \infty$  and  $\theta(P_{t/\sqrt{n}, g})$  is uniquely determined by (25). Then,  $\theta(P_{t/\sqrt{n}, g})$  defined by (25) satisfies

$$\begin{aligned} &\sqrt{n}(\theta(P_{t/\sqrt{n}, g}) - \theta(P)) \\ &\rightarrow M^{-1}E_P[m(X_i, A_i, R_i, \theta(P))(g^X(X_i) + I\{A_i = 1\}g^{R(1)|X}(R_i|X_i) + I\{A_i = 0\}g^{R(0)|X}(R_i|X_i))]'t \\ &= E_P[\psi^*(X_i, A_i, R_i, \theta(P))(g^X(X_i) + I\{A_i = 1\}g^{R(1)|X}(R_i|X_i) + I\{A_i = 0\}g^{R(0)|X}(R_i|X_i))]'t , \end{aligned}$$

where

$$\begin{aligned} &\psi^*(X_i, A_i, R_i, \theta(P)) \\ &= M^{-1} \left( \eta(X_i)E_Q[m(X_i, 1, R_i(1), \theta(P))|X_i] + (1 - \eta(X_i))E_Q[m(X_i, 0, R_i(0), \theta(P))|X_i] \right. \\ &\quad + I\{A_i = 1\}(m(X_i, 1, R_i, \theta(P)) - E_Q[m(X_i, 1, R_i(1), \theta(P))|X_i]) \\ &\quad \left. + I\{A_i = 0\}(m(X_i, 0, R_i, \theta(P)) - E_Q[m(X_i, 0, R_i(0), \theta(P))|X_i]) \right) . \end{aligned}$$

PROOF. In what follows, we only use the property that the quadratic mean derivative of  $P_{t/\sqrt{n}, g}$  is given by  $s'_g t$ . Therefore, for ease of notation we consider a generic one-dimensional submodel  $\{P_\nu : \nu \in [-\epsilon, \epsilon]\}$  that satisfies (20)–(21) for some  $g = (g^X, g^{R(1)|X}, g^{R(0)|X})$ , each component of which is a one-dimensional function. (25) implies

$$\begin{aligned} 0 = \int m(x, 1, r, \theta(P_\nu))q_\nu^{R(1)|X}(r|x)d\mu^R(r)\eta(x)q_\nu^X(x)d\mu^X(x) \\ + \int m(x, 0, r, \theta(P_\nu))q_\nu^{R(0)|X}(r|x)d\mu^R(r)(1 - \eta(x))q_\nu^X(x)d\mu^X(x) \end{aligned}$$

Note that

$$\begin{aligned} &\int m(x, 1, r, \theta(P))q_\nu^{R(1)|X}(r|x)d\mu^R(r)\eta(x)q_\nu^X(x)d\mu^X(x) \\ &- \int m(x, 1, r, \theta(P))q_\nu^{R(1)|X}(r|x)d\mu^R(r)\eta(x)q_\nu^X(x)d\mu^X(x) = \gamma_1(\nu) + \gamma_2(\nu) + \gamma_3(\nu) + \gamma_4(\nu) , \end{aligned}$$

where

$$\begin{aligned}
\gamma_1(\nu) &= \int m(x, 1, r, \theta(P)) (q_\nu^{R(1)|X}(r|x)^{1/2} - q^{R(1)|X}(r|x)^{1/2}) q_\nu^{R(1)|X}(r|x)^{1/2} d\mu^R(r) \eta(x) q_\nu^X(x) d\mu^X(x) \\
\gamma_2(\nu) &= \int m(x, 1, r, \theta(P)) (q_\nu^{R(1)|X}(r|x)^{1/2} - q^{R(1)|X}(r|x)^{1/2}) q^{R(1)|X}(r|x)^{1/2} d\mu^R(r) \eta(x) q_\nu^X(x) d\mu^X(x) \\
\gamma_3(\nu) &= \int m(x, 1, r, \theta(P)) q^{R(1)|X}(r|x) d\mu^R(r) \eta(x) (q_\nu^X(x)^{1/2} - q^X(x)^{1/2}) q_\nu^X(x)^{1/2} d\mu^X(x) \\
\gamma_4(\nu) &= \int m(x, 1, r, \theta(P)) q^{R(1)|X}(r|x) d\mu^R(r) \eta(x) (q_\nu^X(x)^{1/2} - q^X(x)^{1/2}) q^X(x)^{1/2} d\mu^X(x) .
\end{aligned}$$

It follows from the Cauchy-Schwarz inequality that

$$\begin{aligned}
& \frac{1}{\nu} \gamma_4(\nu) - \int m(x, 1, r, \theta(P)) q^{R(1)|X}(r|x) d\mu^R(r) \eta(x) \frac{1}{2} g^X(x) q^X(x)^{1/2} \times q^X(x)^{1/2} d\mu^X(x) \\
& \leq \int \left( m(x, 1, r, \theta(P))^2 q^{R(1)|X}(r|x) d\mu^R(r) \eta(x)^2 q^X(x) d\mu^X(x) \right)^{1/2} \\
& \quad \times \left( \int q^{R(1)|X}(r|x) d\mu^R(r) \left( \frac{1}{\nu} (q_\nu^X(x)^{1/2} - q^X(x)^{1/2}) - \frac{1}{2} g^X(x) q^X(x)^{1/2} \right)^2 d\mu^X(x) \right)^{1/2} \rightarrow 0
\end{aligned}$$

by the assumption that  $E_P[m(X, a, R(a), \theta(P))^2] < \infty$ , the facts that  $0 \leq \eta(x) \leq 1$ ,  $\int q^{R(1)|X}(r|x) d\mu^R(r) = 1$ , and (20). Similar arguments implies as  $\nu \rightarrow 0$ ,

$$\frac{1}{\nu} \gamma_1(\nu) - \int m(x, 1, r, \theta(P)) \frac{1}{2} g^{R(1)|X}(r|x) q^{R(1)|X}(r|x) d\mu^R(r) \eta(x) q^X(x) d\mu^X(x) \rightarrow 0$$

because  $E_{P_\nu}[m(X, a, R(a), \theta(P))^2] = O(1)$  as  $\nu \rightarrow 0$ . The limits of  $\gamma_2(\nu)$  and  $\gamma_3(\nu)$  can be derived following similar arguments using the last two conditions in (26). Combining all previous results yields

$$\begin{aligned}
& \left. \frac{\partial}{\partial \nu} E_{P_\nu}[m(X, A, R, \theta(P))] \right|_{\nu=0} \\
& = E_Q[m(X, 1, R(1), \theta(P))(g^X(X) + g^{R(1)|X}(R|X))\eta(X)] \\
& \quad + E_Q[m(X, 0, R(0), \theta(P))(g^X(X) + g^{R(0)|X}(R|X))(1 - \eta(X))] \\
& = E_P[m(X, A, R, \theta(P))(g^X(X) + I\{A = 1\}g^{R(1)|X}(R) + I\{A = 0\}g^{R(0)|X}(R))] .
\end{aligned}$$

On the other hand, by definition

$$M_{\theta(P)} = \left. \frac{\partial}{\partial \theta} E_P[m(X, A, R, \theta)] \right|_{\theta=\theta(P)} .$$

The formula for the derivative therefore follows from the implicit function theorem (in particular, because we have assumed the existence of  $\theta(P_\nu)$  along the path, it follows from the last part of the proof of Theorem 3.2.1 in Krantz and Parks (2013)). The second equality follows from Lemma B.5 together with Condition B.1. ■

Finally, to prove Theorem 4.1 we require the following additional regularity condition:

**Assumption B.1.** For every function  $g$  satisfying Condition B.1 and every  $t \in \mathbf{R}^{d_\theta}$  there exists a submodel

$P_{t/\sqrt{n},g}$  for which (26) holds as  $n \rightarrow \infty$ , and  $\theta(P_{t/\sqrt{n},g})$  is uniquely determined by (25).

This assumption guarantees that every element satisfying Condition B.1 has a corresponding path for which we can apply Lemma B.4. A similar assumption appears in Chen and Santos (2018) (see their Assumption 4.1(iv)). Note that a simple sufficient condition for the first part of Assumption B.1 is that  $m(x, a, r, \theta(P))$  is a bounded function in  $(x, r)$  on the support of  $(X, R(a))$ . The second part of Assumption B.1 can be verified easily in specific examples (see, for instance, Examples 2.1–2.5 in the main text). Alternatively, Assumption B.1 could be avoided by assuming that we can differentiate under the integral in the final step of the proof of Lemma B.4, from which we would immediately obtain the expression for the pathwise derivative. See, for instance, Newey (1994) and Chen et al. (2008).

PROOF OF THEOREM 4.1. First note  $\theta$  satisfies (23) because of Lemma B.4 and Assumption B.1. The result then follows from Lemma B.3 upon noting that  $\psi^* = s_g$  for some  $g$  that satisfies Condition B.1 because of Assumption 4.2. ■

## B.2 Proof of Theorem 3.1

First note (10) follows from (9) and the same proof as that of Lemma B.3 in Bai (2022). To establish (9), we follow the proof of Theorem 5.21 in van der Vaart (1998). We start by noting that because Assumptions 3.1–3.2, 3.3(e), and 3.4(b) hold, it follows from the same proof as that of Lemma B.3 in Bai (2022) that

$$\begin{aligned} \frac{1}{\sqrt{n}} \sum_{1 \leq i \leq n} m(X_i, A_i, R_i, \theta_0) &= \frac{1}{\sqrt{n}} \sum_{1 \leq i \leq n} \sum_{a \in \{0,1\}} I\{A_i = a\} (m(X_i, a, R_i(a), \theta_0) - E_Q[m(X_i, a, R_i(a), \theta_0)|X_i]) \\ &\quad + \frac{\eta}{\sqrt{n}} \sum_{1 \leq i \leq n} (E_Q[m(X_i, 1, R_i(1), \theta_0)|X_i] - E_Q[m(X_i, 1, R_i(1), \theta_0)]) \\ &\quad + \frac{(1-\eta)}{\sqrt{n}} \sum_{1 \leq i \leq n} (E_Q[m(X_i, 0, R_i(0), \theta_0)|X_i] - E_Q[m(X_i, 0, R_i(0), \theta_0)]) + o_P(1) . \end{aligned}$$

where we note  $\eta E_Q[m(X_i, 1, R_i(1), \theta_0)] + (1-\eta) E_Q[m(X_i, 0, R_i(0), \theta_0)] = E_P[m(X_i, A_i, R_i, \theta_0)] = 0$  by (3). Therefore, by the proof of Theorem 5.21 in van der Vaart (1998), it suffices to show

$$\mathbb{L}_n(\hat{\theta}_n) \xrightarrow{P} 0 , \tag{27}$$

where

$$\begin{aligned} \mathbb{L}_n(\theta) &= \frac{1}{\sqrt{n}} \sum_{1 \leq i \leq n} (m(X_i, A_i, R_i, \theta) - E_P[m(X_i, A_i, R_i, \theta)]) \\ &\quad - \frac{1}{\sqrt{n}} \sum_{1 \leq i \leq n} (m(X_i, A_i, R_i, \theta_0) - E_P[m(X_i, A_i, R_i, \theta_0)]) . \end{aligned}$$

To accomplish this, we study  $m_s$  for  $1 \leq s \leq d_\theta$  separately. It follows from Assumption 3.3(c), (d), and the proof of Proposition 8.11 in Kosorok (2008) that

$$\begin{aligned} & \sup_{\theta \in \Theta: \|\theta - \theta_0\| < \delta} \left| \frac{1}{\sqrt{n}} \sum_{1 \leq i \leq n} (m_s(X_i, A_i, R_i, \theta) - E_P[m_s(X_i, A_i, R_i, \theta)]) \right. \\ & \quad \left. - \frac{1}{\sqrt{n}} \sum_{1 \leq i \leq n} (m_s(X_i, A_i, R_i, \theta_0) - E_P[m_s(X_i, A_i, R_i, \theta_0)]) \right| \\ &= \sup_{\theta \in \Theta^*: \|\theta - \theta_0\| < \delta} \left| \frac{1}{\sqrt{n}} \sum_{1 \leq i \leq n} (m_s(X_i, A_i, R_i, \theta) - E_P[m_s(X_i, A_i, R_i, \theta)]) \right. \\ & \quad \left. - \frac{1}{\sqrt{n}} \sum_{1 \leq i \leq n} (m_s(X_i, A_i, R_i, \theta_0) - E_P[m_s(X_i, A_i, R_i, \theta_0)]) \right|, \end{aligned}$$

and thus since  $\hat{\theta}_n \xrightarrow{P} \theta_0$  by Lemma B.6, to show (27) it suffices to argue that for every  $\epsilon > 0$ ,

$$\lim_{\delta \downarrow 0} \limsup_{n \rightarrow \infty} P \left\{ \sup_{\theta \in \Theta^*: \|\theta - \theta_0\| < \delta} |\mathbb{L}_n^{(s)}(\theta)| > \epsilon \right\} = 0, \quad (28)$$

where

$$\begin{aligned} \mathbb{L}_n^{(s)}(\theta) &= \frac{1}{\sqrt{n}} \sum_{1 \leq i \leq n} (m_s(X_i, A_i, R_i, \theta) - E_P[m_s(X_i, A_i, R_i, \theta)]) \\ & \quad - \frac{1}{\sqrt{n}} \sum_{1 \leq i \leq n} (m_s(X_i, A_i, R_i, \theta_0) - E_P[m_s(X_i, A_i, R_i, \theta_0)]). \end{aligned}$$

Consider the following decomposition:

$$|\mathbb{L}_n^{(s)}(\theta)| \leq \sum_{a \in \{0,1\}} (\mathbb{L}_{1,a,n}^{(s)}(\theta) + \mathbb{L}_{2,a,n}^{(s)}(\theta) + \mathbb{L}_{3,a,n}^{(s)}(\theta)),$$

where

$$\begin{aligned} \mathbb{L}_{1,a,n}^{(s)}(\theta) &= \left| \frac{1}{\sqrt{n}} \sum_{1 \leq i \leq n} I\{A_i = a\} (m_s(X_i, a, R_i(a), \theta) - E_Q[m_s(X_i, a, R_i(a), \theta)|X_i]) \right. \\ & \quad \left. - \frac{1}{\sqrt{n}} \sum_{1 \leq i \leq n} I\{A_i = a\} (m_s(X_i, a, R_i(a), \theta_0) - E_Q[m_s(X_i, a, R_i(a), \theta_0)|X_i]) \right| \\ \mathbb{L}_{2,a,n}^{(s)}(\theta) &= \left| \frac{1}{\sqrt{n}} \sum_{1 \leq i \leq n} (I\{A_i = a\} - \eta_a) \right. \\ & \quad \left. \times (E_Q[m_s(X_i, a, R_i(a), \theta)|X_i] - E_Q[m_s(X_i, a, R_i(a), \theta_0)|X_i]) \right| \\ \mathbb{L}_{3,a,n}^{(s)}(\theta) &= \left| \frac{\eta_a}{\sqrt{n}} \sum_{1 \leq i \leq n} (E_Q[m_s(X_i, a, R_i(a), \theta)|X_i] - E_Q[m_s(X, a, R(a), \theta)]) \right. \\ & \quad \left. - (E_Q[m_s(X_i, a, R_i(a), \theta_0)|X_i] - E_Q[m_s(X, a, R(a), \theta_0)]) \right|, \end{aligned}$$

where  $\eta_1 = \eta$  and  $\eta_0 = 1 - \eta$ . Then to establish (28), it suffices to establish that

$$\lim_{\delta \downarrow 0} \limsup_{n \rightarrow \infty} P \left\{ \sup_{\theta \in \Theta^* : \|\theta - \theta_0\| < \delta} \mathbb{L}_{\ell, a, n}^{(s)}(\theta) > \epsilon \right\} = 0. \quad (29)$$

for  $\ell \in \{1, 2, 3\}$  and  $a \in \{0, 1\}$ .

**Step 1.** First we consider  $\mathbb{L}_{3, a, n}^{(s)}$ . It follows from Assumption 3.4(d) and Theorems 2.5.2 and 2.6.7 in van der Vaart and Wellner (1996) that the class of functions

$$\{E_Q[m_s(X, a, R(a), \theta)|X = x] : \theta \in \Theta^*\},$$

is Donsker, and thus we obtain by Theorem 3.34 in Dudley (2014) that

$$\lim_{\delta \downarrow 0} \limsup_{n \rightarrow \infty} P \left\{ \sup_{\theta \in \Theta^* : \rho_Q(\theta, \theta_0) < \delta} \mathbb{L}_{3, a, n}^{(s)}(\theta) > \epsilon \right\} = 0,$$

where  $\rho_Q(\theta, \theta_0) = E_Q[(E_Q[m_s(X, a, R(a), \theta)|X] - E_Q[m_s(X, a, R(a), \theta_0)|X])^2]$ . We then obtain (29) for  $\ell = 3$  since, by Assumption 3.4(c) as  $\theta \rightarrow \theta_0$ ,

$$\begin{aligned} \rho_Q(\theta, \theta_0) &= E_Q[(E_Q[m_s(X, a, R(a), \theta)|X] - E_Q[m_s(X, a, R(a), \theta_0)|X])^2] \\ &\leq E_Q[(m_s(X, a, R(a), \theta) - m_s(X, a, R(a), \theta_0))^2] \rightarrow 0. \end{aligned} \quad (30)$$

**Step 2.** Next, we study  $\mathbb{L}_{2, a, n}^{(s)}$ . Define

$$f(X, \theta) = E[m_s(X, a, R(a), \theta)|X] - E[m_s(X, a, R(a), \theta_0)|X].$$

Note

$$\mathbb{L}_{2, a, n}^{(s)}(\theta) = C \left| \frac{1}{\sqrt{n/k}} \sum_{1 \leq j \leq n/k} \alpha_j(\theta) \right|,$$

for some constant  $C > 0$ , where  $\alpha_j(\theta) \in \{\frac{1}{\ell} \sum_{i \in I} f(X_i, \theta) - \frac{1}{k-\ell} \sum_{i \in \lambda_j \setminus I} f(X_i, \theta) : I \subset \lambda_j, |I| = \ell\}$ ,  $E[\alpha_j(\theta)|X^{(n)}] = 0$ , and  $\alpha_j(\theta), 1 \leq j \leq n/k$  are independent conditional on  $X^{(n)}$ .

Define

$$\begin{aligned} h(x_1, \dots, x_k, \theta) &= \sup_{I \subset \{1, \dots, k\}, |I| = \ell} \left( \frac{1}{\ell} \sum_{i \in I} f(x_i, \theta) - \frac{1}{k-\ell} \sum_{i \in \{1, \dots, k\} \setminus I} f(x_i, \theta) \right) \\ &\quad - \inf_{I \subset \{1, \dots, k\}, |I| = \ell} \left( \frac{1}{\ell} \sum_{i \in I} f(x_i, \theta) - \frac{1}{k-\ell} \sum_{i \in \{1, \dots, k\} \setminus I} f(x_i, \theta) \right) \end{aligned}$$

and the classes of functions

$$\mathcal{H}_\delta = \{h(x_1, \dots, x_k, \theta) : \theta \in \Theta^*, \|\theta - \theta_0\| < \delta\}$$

$$\mathcal{H}_\infty = \{h(x_1, \dots, x_k, \theta) : \theta \in \Theta^*\}.$$

Let  $P_n^\dagger$  denote a measure that puts mass  $\frac{k}{n}$  on each of  $(X_i : i \in \lambda_j), 1 \leq j \leq n/k$ . It follows from a generalized version of Hoeffding's inequality (see, for instance, Theorem 2.2.6 in [Vershynin, 2018](#)) that conditional on  $X^{(n)}$ ,

$$\left\{ \frac{1}{\sqrt{n/k}} \sum_{1 \leq j \leq n/k} \alpha_j(\theta) : \theta \in \Theta^*, \|\theta - \theta_0\| < \delta \right\}$$

is sub-Gaussian for the seminorm

$$\|h\|_{P_n^\dagger} = \left( \int h^2 dP_n^\dagger \right)^{1/2}.$$

Let  $N(\epsilon, \mathcal{H}_\delta, L_2(P_n^\dagger))$  denote the covering number of  $\mathcal{H}_\delta$  with respect to  $\|\cdot\|_{P_n^\dagger}$ . Let  $\delta_n \downarrow 0$  be an arbitrary decreasing sequence. It follows from the maximal inequality in Corollary 2.2.8 in [van der Vaart and Wellner \(1996\)](#) (note  $\alpha_j(\theta_0) = 0$ ) that

$$E \left[ \sup_{\theta \in \Theta^* : \|\theta - \theta_0\| < \delta_n} \mathbb{L}_{2,a,n}^{(s)}(\theta) \middle| X^{(n)} \right] \lesssim \int_0^\infty \sqrt{\log N(\epsilon, \mathcal{H}_{\delta_n}, L_2(P_n^\dagger))} d\epsilon. \quad (31)$$

The upper limit of the integral is in fact  $c_n$ , where

$$\begin{aligned} c_n^2 &= \sup_{\theta \in \Theta^* : \|\theta - \theta_0\| < \delta_n} \frac{4k}{n} \sum_{1 \leq j \leq n/k} \sup_{I \in \lambda_j, |I|=\ell} \left( \frac{1}{\ell} \sum_{i \in I} f(X_i, \theta) - \frac{1}{k-\ell} \sum_{i \in \lambda_j \setminus I} f(X_i, \theta) \right)^2 \\ &\lesssim \frac{1}{n} \sum_{1 \leq j \leq n/k} \max_{i, i' \in \lambda_j} \|X_i - X_{i'}\|^2 \xrightarrow{P} 0 \end{aligned} \quad (32)$$

by Assumptions [3.2](#) and [3.3\(e\)](#) and the inequality  $(a+b)^2 \leq 2(a^2 + b^2)$ . Moreover,

$$H(x_1, \dots, x_k) = \sum_{1 \leq i \leq k} E \left[ \sup_{\theta \in \Theta^*} |m_s(X, a, R(a), \theta)| + |m_s(X, a, R(a), \theta_0)| \middle| X = x_i \right]$$

is an envelope function for  $\mathcal{H}_\infty$  (and thus  $\mathcal{H}_\delta$  for all  $\delta$ ) and  $E[H^2] < \infty$  by Assumption [3.4\(b\)](#). A change of variable in [\(31\)](#) implies

$$\begin{aligned} E \left[ \sup_{\theta \in \Theta^* : \|\theta - \theta_0\| < \delta_n} \mathbb{L}_{2,a,n}^{(s)}(\theta) \middle| X^{(n)} \right] &\lesssim \int_0^{\frac{c_n}{\|H\|_{P_n^\dagger}}} \sqrt{\log N(\epsilon \|H\|_{P_n^\dagger}, \mathcal{H}_{\delta_n}, L_2(P_n^\dagger))} d\epsilon \|H\|_{P_n^\dagger} \\ &\leq \int_0^{\frac{c_n}{\|H\|_{P_n^\dagger}}} \sup_{\nu} \sqrt{\log N(\epsilon \|H\|_{\nu}, \mathcal{H}_\infty, L_2(\nu))} d\epsilon \|H\|_{P_n^\dagger}, \end{aligned}$$

where the supremum for  $\nu$  is over probability measures with discrete support such that  $\|H\|_{\nu} > 0$ . Also note that if  $\|H\|_{P_n^\dagger} = 0$  then the conditional expectation on the left-hand side is trivially zero, so we can without loss of generality assume  $\|H\|_{P_n^\dagger} > 0$ . The Cauchy-Schwarz inequality implies

$$\begin{aligned} E \left[ \int_0^{\frac{c_n}{\|H\|_{P_n^\dagger}}} \sup_{\nu} \sqrt{\log N(\epsilon \|H\|_{\nu}, \mathcal{H}_\infty, L_2(\nu))} d\epsilon \|H\|_{P_n^\dagger} \right] \\ \leq E \left[ \left( \int_0^{\frac{c_n}{\|H\|_{P_n^\dagger}}} \sup_{\nu} \sqrt{\log N(\epsilon \|H\|_{\nu}, \mathcal{H}_\infty, L_2(\nu))} d\epsilon \right)^2 \right]^{1/2} E[\|H\|_{P_n^\dagger}^2]^{1/2}, \end{aligned}$$

where the supremum is over all measures  $\nu$  with discrete support such that  $\|H\|_{\nu} > 0$ . It follows from

Assumption 3.4(b), the inequality  $(a + b)^2 \leq 2(a^2 + b^2)$ , and the conditional Jensen's inequality that

$$E[\|H\|_{P_n^\dagger}^2] \lesssim E\left[\frac{1}{n} \sum_{1 \leq i \leq n} E\left[\sup_{\theta \in \Theta^*} |m_s(X_i, a, R_i(a), \theta)| \middle| X_i\right]^2\right] \leq E\left[\sup_{\theta \in \Theta^*} |m_s(X, a, R(a), \theta)|^2\right] < \infty .$$

On the other hand,

$$\|H\|_{P_n^\dagger}^2 \geq \frac{1}{n} \sum_{1 \leq i \leq n} E\left[\sup_{\theta \in \Theta^*} |m_s(X_i, a, R_i(a), \theta)| \middle| X_i\right]^2 \xrightarrow{P} E\left[E\left[\sup_{\theta \in \Theta^*} |m_s(X, a, R(a), \theta)| \middle| X\right]^2\right], \quad (33)$$

the right-hand side of which can be assumed to be strictly positive, because otherwise  $\sup_{\theta \in \Theta^* : \|\theta - \theta_0\| < \delta_n} \mathbb{L}_{2,a,n}^{(s)}(\theta) = 0$ . Therefore, it follows from (32) and (33) that

$$\frac{c_n}{\|H\|_{P_n^\dagger}} \xrightarrow{P} 0 . \quad (34)$$

From Assumption 3.4(d), Lemma B.7, Theorem 2.6.7 in van der Vaart and Wellner (1996), and Lemma 9.13 in Kosorok (2008), we know

$$\int_0^1 \sup_{\nu} \sqrt{\log N(\epsilon \|H\|_{\nu}, \mathcal{H}_{\infty}, L_2(\nu))} d\epsilon < \infty .$$

Therefore,

$$E\left[\left(\int_0^{\frac{c_n}{\|H\|_{P_n^\dagger}}} \sup_{\nu} \sqrt{\log N(\epsilon \|H\|_{\nu}, \mathcal{H}_{\infty}, L_2(\nu))} d\epsilon\right)^2\right] \rightarrow 0$$

by Lemma B.8 combined with (34) and the continuous mapping theorem. Therefore, it follows from Markov's inequality that

$$\sup_{\theta \in \Theta^* : \|\theta - \theta_0\| < \delta_n} \mathbb{L}_{2,a,n}^{(s)}(\theta) \xrightarrow{P} 0 ,$$

as  $n \rightarrow \infty$ , from which (29) follows (see, for instance, Section 2.1.2 in van der Vaart and Wellner, 1996).

**Step 3.** Finally, we study  $\mathbb{L}_{1,a,n}^{(s)}(\theta)$ . Define

$$\mathbb{B}_n(\theta) = \frac{1}{\sqrt{\eta_a n}} \sum_{1 \leq i \leq n} I\{A_i = a\} (m_s(X_i, a, R_i(a), \theta) - E_Q[m_s(X_i, a, R_i(a), \theta) | X_i]) ,$$

Let  $\delta_n \downarrow 0$  be an arbitrary decreasing sequence. To establish our result we will show

$$\sup_{\theta \in \Theta^* : \|\theta - \theta_0\| < \delta_n} |\mathbb{B}_n(\theta) - \mathbb{B}_n(\theta_0)| \xrightarrow{P} 0 , \quad (35)$$

as  $n \rightarrow \infty$ . As in the proof of Lemma B.6, we define

$$\tilde{P}_n = \frac{1}{\eta_a n} \sum_{1 \leq i \leq n : A_i = a} \delta_{(X_i, R_i(a))} .$$

Define the classes of functions

$$\mathcal{F}_{\theta_0, \infty} = \{m_s(x, a, r(a), \theta) : \theta \in \Theta^*\} .$$



Pick an envelope function for  $\mathcal{F}_{\theta_0, \infty}$  as

$$F = \sup_{\theta \in \Theta^*} |m_s(X, a, R(a), \theta)|.$$

and define

$$\zeta_n^2 = \sup_{\theta \in \Theta^* : \|\theta - \theta_0\| < \delta_n} \int (m_s(X, a, R(a), \theta) - m_s(X, a, R(a), \theta_0))^2 d\tilde{P}_n.$$

**Step 3(a).** Our next goal is to show for every  $\xi > 0$ ,

$$P\{\zeta_n^2 > \xi | X^{(n)}, A^{(n)}\} \xrightarrow{P} 0. \quad (36)$$

To do so, first note by triangle inequality that  $\zeta_n^2 \leq \mathbb{C}_{1,n} + \mathbb{C}_{2,n} + \mathbb{C}_{3,n}$ , where

$$\begin{aligned} \mathbb{C}_{1,n} &= \sup_{\theta \in \Theta^*} \left| \int (m_s(X, a, R(a), \theta) - m_s(X, a, R(a), \theta_0))^2 d\tilde{P}_n \right. \\ &\quad \left. - E \left[ \int (m_s(X, a, R(a), \theta) - m_s(X, a, R(a), \theta_0))^2 d\tilde{P}_n \middle| X^{(n)}, A^{(n)} \right] \right| \\ \mathbb{C}_{2,n} &= \sup_{\theta \in \Theta^*} \left| E \left[ \int (m_s(X, a, R(a), \theta) - m_s(X, a, R(a), \theta_0))^2 d\tilde{P}_n \middle| X^{(n)}, A^{(n)} \right] \right. \\ &\quad \left. - E[(m_s(X, a, R(a), \theta) - m_s(X, a, R(a), \theta_0))^2] \right| \\ \mathbb{C}_{3,n} &= \sup_{\theta \in \Theta^* : \|\theta - \theta_0\| < \delta_n} E[(m_s(X, a, R(a), \theta) - m_s(X, a, R(a), \theta_0))^2]. \end{aligned}$$

Assumption 3.4(c) implies

$$\mathbb{C}_{3,n} = \sup_{\theta \in \Theta^* : \|\theta - \theta_0\| < \delta_n} E[(m_s(X, a, R(a), \theta) - m_s(X, a, R(a), \theta_0))^2] \rightarrow 0. \quad (37)$$

Next, Assumption 3.4(b), (f) and similar arguments to those used to show (56) and (57) are  $o_P(1)$  imply

$$\begin{aligned} \mathbb{C}_{2,n} &= \sup_{\theta \in \Theta^*} \left| E \left[ \int (m_s(X, a, R(a), \theta) - m_s(X, a, R(a), \theta_0))^2 d\tilde{P}_n \middle| X^{(n)}, A^{(n)} \right] \right. \\ &\quad \left. - E[(m_s(X, a, R(a), \theta) - m_s(X, a, R(a), \theta_0))^2] \right| \xrightarrow{P} 0. \quad (38) \end{aligned}$$

Further define

$$\mathcal{G} = \{(m_s(x, a, r(a), \theta) - m_s(x, a, r(a), \theta_0))^2 : \theta \in \Theta^*\}.$$

We then study  $\mathbb{C}_{1,n}$ . We will establish for every  $\xi > 0$ ,

$$P \left\{ \sup_{f \in \mathcal{G}} \left| \int f d\tilde{P}_n - E \left[ \int f d\tilde{P}_n \middle| X^{(n)}, A^{(n)} \right] \right| > \xi \middle| X^{(n)}, A^{(n)} \right\} \xrightarrow{P} 0 \quad (39)$$

as  $n \rightarrow \infty$ . It follows the symmetrization Lemma 6.2 in [Ledoux and Talagrand \(1991\)](#) applied conditional

on  $X^{(n)}, A^{(n)}$  for the distribution

$$\bigotimes_{1 \leq i \leq n: A_i=1} P\{X_i, R_i(1)|X_i\}$$

that

$$\begin{aligned} E \left[ \sup_{f \in \mathcal{G}} \left| \int f d\tilde{P}_n - E \left[ \int f d\tilde{P}_n \middle| X^{(n)}, A^{(n)} \right] \right| \middle| X^{(n)}, A^{(n)} \right] \\ \leq 2E_P \left[ E_\tau \left[ \sup_{f \in \mathcal{G}} \left| \frac{1}{\eta_a n} \sum_{1 \leq i \leq n} \tau_i f(X_i, R_i(a)) \right| \right] \middle| X^{(n)}, A^{(n)} \right], \quad (40) \end{aligned}$$

where  $E_\tau[\cdot]$  should be understood as the expectation with respect to  $(\tau_i, 1 \leq i \leq n)$ , holding all else fixed. Note Assumption 3.3(b) and Theorem 2.6.7 in [van der Vaart and Wellner \(1996\)](#) imply  $\mathcal{F}_{\theta_0, \infty}$  is totally bounded in  $L_2(\tilde{P}_n)$ . Accordingly, for  $\epsilon > 0$ , let  $N(\epsilon, \mathcal{F}_{\theta_0, \infty}, L_2(\tilde{P}_n))$  denote the covering number of  $\mathcal{F}_{\theta_0, \infty}$  with respect to  $L_2(\tilde{P}_n)$ . Let  $f_1, f_2$  be any pair of functions in  $\mathcal{G}$ , where we denote

$$f_j = (m_s(x, a, r(a), \theta_j) - m_s(x, a, r(a), \theta_0))^2, \quad j = 1, 2,$$

then the Cauchy-Schwarz inequality implies

$$\int |f_1 - f_2| d\tilde{P}_n \leq \int |m_s(x, a, r(a), \theta_1) - m_s(x, a, r(a), \theta_2)| 2F d\tilde{P}_n \leq \|m_s(\cdot, \theta_1) - m_s(\cdot, \theta_2)\|_{\tilde{P}_n} 2\|F\|_{\tilde{P}_n}$$

where  $\|\cdot\|_{\tilde{P}_n}$  denotes the  $L_2(\tilde{P}_n)$ -norm. Therefore

$$N(2\epsilon\|F\|_{\tilde{P}_n}^2, \mathcal{G}, L_1(\tilde{P}_n)) \leq N(\epsilon\|F\|_{\tilde{P}_n}, \mathcal{F}_{\theta_0, \infty}, L_2(\tilde{P}_n)). \quad (41)$$

For every  $\epsilon > 0$ , the right-hand side is uniformly bounded across  $n$  by Assumption 3.3(b) and Theorem 2.6.7 in [van der Vaart and Wellner \(1996\)](#). Note it follows from Assumptions 3.1–3.2, 3.4(b), (f), and similar arguments to those in the first part of the proof of Lemma S.1.5 in [Bai et al. \(2022\)](#) that

$$E[\|F\|_{\tilde{P}_n}^2 | X^{(n)}, A^{(n)}] = \frac{1}{\eta_a n} \sum_{1 \leq i \leq n} I\{A_i = a\} E[F^2 | X_i] \xrightarrow{P} E[F^2]. \quad (42)$$

We can assume without loss of generality  $E[F^2] > 0$  because otherwise  $m_s(x, a, r(a), \theta) \equiv 0$ . Therefore,

$$P \left\{ E[\|F\|_{\tilde{P}_n}^2 | X^{(n)}, A^{(n)}] \leq \frac{1}{2} E[F^2] \right\} \rightarrow 0. \quad (43)$$

On the other hand, Assumptions 3.1–3.2, 3.4(b), (f), and similar arguments to those in the last part of the proof of Lemma S.1.5 in [Bai et al. \(2022\)](#) that

$$P \left\{ \left| \|F\|_{\tilde{P}_n}^2 - E[\|F\|_{\tilde{P}_n}^2 | X^{(n)}, A^{(n)}] \right| > \frac{1}{4} E[F^2] \middle| X^{(n)}, A^{(n)} \right\} \xrightarrow{P} 0. \quad (44)$$

(39) now follows from (43)–(44) and similar arguments to those used in the last step of the proof of Lemma B.6.

To conclude (36) holds, note  $\mathbb{C}_{3,n}$  is a sequence of constants and  $\mathbb{C}_{2,n}$  is a function of  $X^{(n)}, A^{(n)}$ , and hence

$$\begin{aligned}
& P\{P\{\zeta_n^2 > \xi | X^{(n)}, A^{(n)}\} > \epsilon\} \\
& \leq P\left\{\mathbb{C}_{2,n} > \frac{\xi}{3}\right\} + P\left\{\mathbb{C}_{3,n} > \frac{\xi}{3}\right\} \\
& \quad + P\left\{P\{\mathbb{C}_{1,n} + \mathbb{C}_{2,n} + \mathbb{C}_{3,n} > \xi | X^{(n)}, A^{(n)}\} > \epsilon, \mathbb{C}_{2,n} \leq \frac{\xi}{3}, \mathbb{C}_{3,n} \leq \frac{\xi}{3}\right\} \\
& \leq P\left\{\mathbb{C}_{2,n} > \frac{\xi}{3}\right\} + P\left\{\mathbb{C}_{3,n} > \frac{\xi}{3}\right\} + P\left\{P\left\{\mathbb{C}_{1,n} > \frac{\xi}{3} \middle| X^{(n)}, A^{(n)}\right\} > \epsilon\right\} \\
& \xrightarrow{P} 0,
\end{aligned}$$

where the convergence follows from (37), (38), and (39).

**Step 3(b).** Next, we show for every  $\xi > 0$ ,

$$P\left\{\frac{\zeta_n^2}{\|F\|_{\tilde{P}_n}^2} > \xi \middle| X^{(n)}, A^{(n)}\right\} \xrightarrow{P} 0. \quad (45)$$

For every  $\epsilon > 0$ ,

$$\begin{aligned}
& P\left\{P\left\{\frac{\zeta_n^2}{\|F\|_{\tilde{P}_n}^2} > \xi \middle| X^{(n)}, A^{(n)}\right\} > \epsilon\right\} \\
& \leq P\left\{E[\|F\|_{\tilde{P}_n}^2 | X^{(n)}, A^{(n)}] \leq \frac{1}{2}E[F^2]\right\} \\
& \quad + P\left\{P\left\{\frac{\zeta_n^2}{\|F\|_{\tilde{P}_n}^2} > \xi \middle| X^{(n)}, A^{(n)}\right\} > \epsilon, E[\|F\|_{\tilde{P}_n}^2 | X^{(n)}, A^{(n)}] > \frac{1}{2}E[F^2]\right\} \\
& \leq P\left\{E[\|F\|_{\tilde{P}_n}^2 | X^{(n)}, A^{(n)}] \leq \frac{1}{2}E[F^2]\right\} \\
& + P\left\{P\left\{\left\{\zeta_n^2 > \frac{1}{4}\xi E[F^2]\right\} \cup \left\{\left|\|F\|_{\tilde{P}_n}^2 - E[\|F\|_{\tilde{P}_n}^2 | X^{(n)}, A^{(n)}]\right| > \frac{1}{4}E[F^2]\right\} \middle| X^{(n)}, A^{(n)}\right\} > \epsilon, \right. \\
& \quad \left. E[\|F\|_{\tilde{P}_n}^2 | X^{(n)}, A^{(n)}] > \frac{1}{2}E[F^2]\right\} \\
& \leq P\left\{E[\|F\|_{\tilde{P}_n}^2 | X^{(n)}, A^{(n)}] \leq \frac{1}{2}E[F^2]\right\} \\
& \quad + P\left\{P\left\{\zeta_n^2 > \frac{1}{4}\xi E[F^2] \middle| X^{(n)}, A^{(n)}\right\} > \frac{\epsilon}{2}\right\} \\
& \quad + P\left\{P\left\{\left|\|F\|_{\tilde{P}_n}^2 - E[\|F\|_{\tilde{P}_n}^2 | X^{(n)}, A^{(n)}]\right| > \frac{1}{4}E[F^2] \middle| X^{(n)}, A^{(n)}\right\} > \frac{\epsilon}{2}\right\} \xrightarrow{P} 0,
\end{aligned}$$

where we use the fact that  $E[\|F\|_{\tilde{P}_n}^2 | X^{(n)}, A^{(n)}]$  is a function of  $X^{(n)}, A^{(n)}$  and the convergence follows from (36) and (43)–(44).

**Step 3(c).** Fix  $\epsilon > 0$ . Following almost verbatim the first part of the proof of Theorem 2.5.2 in [van der Vaart and Wellner \(1996\)](#), with  $\tilde{P}_n$  replacing the empirical measure, we obtain

$$P\left\{\sup_{\theta \in \Theta^*: \|\theta - \theta_0\| < \delta_n} |\mathbb{B}_n(\theta) - \mathbb{B}_n(\theta_0)| > \epsilon \middle| X^{(n)}, A^{(n)}\right\}$$

$$\leq \frac{1}{\epsilon} E \left[ \left( \int_0^{\frac{\zeta_n}{\|F\|_{\tilde{P}_n}}} \sup_{\nu} \sqrt{\log N(\epsilon \|F\|_{\nu}, \mathcal{F}_{\theta_0, \infty}, L_2(\nu))} d\epsilon \right)^2 \middle| X^{(n)}, A^{(n)} \right]^{1/2} E[\|F\|_{\tilde{P}_n}^2 \|X^{(n)}, A^{(n)}\|^{1/2}] , \quad (46)$$

where the supremum for  $\nu$  is over probability measures with discrete supports. Also note that if  $\|F\|_{\tilde{P}_n} = 0$  then the conditional expectation on the left-hand side is trivially zero, so we can without loss of generality assume  $\|F\|_{\tilde{P}_n} > 0$ . Assumption 3.3(b) implies

$$E \left[ \left( \int_0^{\frac{\zeta_n}{\|F\|_{\tilde{P}_n}}} \sup_{\nu} \sqrt{\log N(\epsilon \|F\|_{\nu}, \mathcal{F}_{\theta_0, \infty}, L_2(\nu))} d\epsilon \right)^2 \middle| X^{(n)}, A^{(n)} \right] \leq \left( \int_0^{\infty} \sup_{\nu} \sqrt{\log N(\epsilon \|F\|_{\nu}, \mathcal{F}_{\theta_0, \infty}, L_2(\nu))} d\epsilon \right)^2 < \infty . \quad (47)$$

We now argue

$$E \left[ \left( \int_0^{\frac{\zeta_n}{\|F\|_{\tilde{P}_n}}} \sup_{\nu} \sqrt{\log N(\epsilon \|F\|_{\nu}, \mathcal{F}_{\theta_0, \infty}, L_2(\nu))} d\epsilon \right)^2 \middle| X^{(n)}, A^{(n)} \right] \xrightarrow{P} 0 . \quad (48)$$

Note the last inequality in (47) and the dominated convergence theorem implies that for every  $\epsilon > 0$ , there exists a  $\xi > 0$  such that

$$\left( \int_0^{\xi} \sup_{\nu} \sqrt{\log N(\epsilon \|F\|_{\nu}, \mathcal{F}_{\theta_0, \infty}, L_2(\nu))} d\epsilon \right)^2 < \epsilon . \quad (49)$$

Then consider the following decomposition:

$$\begin{aligned} & E \left[ \left( \int_0^{\frac{\zeta_n}{\|F\|_{\tilde{P}_n}}} \sup_{\nu} \sqrt{\log N(\epsilon \|F\|_{\nu}, \mathcal{F}_{\theta_0, \infty}, L_2(\nu))} d\epsilon \right)^2 \middle| X^{(n)}, A^{(n)} \right] \\ &= E \left[ \left( \int_0^{\frac{\zeta_n}{\|F\|_{\tilde{P}_n}}} \sup_{\nu} \sqrt{\log N(\epsilon \|F\|_{\nu}, \mathcal{F}_{\theta_0, \infty}, L_2(\nu))} d\epsilon \right)^2 I \left\{ \frac{\zeta_n}{\|F\|_{\tilde{P}_n}} \leq \xi \right\} \middle| X^{(n)}, A^{(n)} \right] \\ &\quad + E \left[ \left( \int_0^{\frac{\zeta_n}{\|F\|_{\tilde{P}_n}}} \sup_{\nu} \sqrt{\log N(\epsilon \|F\|_{\nu}, \mathcal{F}_{\theta_0, \infty}, L_2(\nu))} d\epsilon \right)^2 I \left\{ \frac{\zeta_n}{\|F\|_{\tilde{P}_n}} > \xi \right\} \middle| X^{(n)}, A^{(n)} \right] \\ &\lesssim \epsilon + P \left\{ \frac{\zeta_n}{\|F\|_{\tilde{P}_n}} > \xi \middle| X^{(n)}, A^{(n)} \right\} \xrightarrow{P} \epsilon , \end{aligned}$$

where the inequality follows from (47) and (49) and the convergence follows from (45). Because  $\epsilon > 0$  was arbitrary, (48) follows.

It thus follows from (42), (46), and (48) that

$$P \left\{ \sup_{\theta \in \Theta^*: \|\theta - \theta_0\| < \delta_n} |\mathbb{B}_n(\theta) - \mathbb{B}_n(\theta_0)| > \epsilon \middle| X^{(n)}, A^{(n)} \right\} \xrightarrow{P} 0 .$$

By the law of iterated expectations and the dominated convergence theorem we thus obtain

$$P \left\{ \sup_{\theta \in \Theta^*: \|\theta - \theta_0\| < \delta_n} |\mathbb{B}_n(\theta) - \mathbb{B}_n(\theta_0)| > \epsilon \right\} \rightarrow 0 ,$$

as desired. ■

### B.3 Auxiliary Lemmas

**Lemma B.5.** *Suppose (2) holds and  $\Pr\{A_i = 1|X_i = x\}$  as a function is identical across  $1 \leq i \leq n$ . Then,*

$$(R_i(1), R_i(0)) \perp\!\!\!\perp A_i | X_i . \quad (50)$$

Moreover,  $(X_i, A_i, R_i)$  is identically distributed across  $1 \leq i \leq n$ .

PROOF. Fix  $a \in \{0, 1\}$  and any Borel sets  $B \in \mathbf{R}^{d_r} \times \mathbf{R}^{d_r}$  and  $C \in \mathbf{R}^{d_x}$ .

$$\begin{aligned} & E[\Pr\{(R_i(1), R_i(0)) \in B, A_i = a | X_i\} I\{X_i \in C\}] \\ &= E[E[\Pr\{(R_i(1), R_i(0)) \in B, A_i = a | X^{(n)}\} | X_i] I\{X_i \in C\}] \\ &= E[E[\Pr\{(R_i(1), R_i(0)) \in B | X^{(n)}\} \Pr\{A_i = a | X^{(n)}\} | X_i] I\{X_i \in C\}] \\ &= E[\Pr\{(R_i(1), R_i(0)) \in B | X_i\} \Pr\{A_i = a | X_i\} I\{X_i \in C\}] , \end{aligned}$$

where the first equality follows from the law of iterated expectations, the second equality follows from (2), the third equality follows from the law of iterated expectation as well as the facts that  $Q_n = Q^n$  and  $\Pr\{A_i = 1|X_i = x\}$  as a function is identical across  $1 \leq i \leq n$ . The first statement of the lemma then follows from the definition of a conditional expectation.

To prove the second statement, fix units  $i$  and  $i'$ . Clearly  $X_i$  and  $X_{i'}$  are identically distributed. Conditional on  $X_i$ , for any Borel set  $C \in \mathbf{R}^{d_r}$  and  $a \in \{0, 1\}$ , it follows (a) that

$$\Pr\{R_i \in C, A_i = a | X_i\} = \Pr\{A_i = a | X_i\} \Pr\{R_i(a) \in C | X_i\} .$$

The conclusion then follows because  $\Pr\{A_i = 1|X_i = x\}$  is identical across  $1 \leq i \leq n$  and  $Q_n = Q^n$ . ■

**Lemma B.6.** *Suppose the treatment assignment mechanism satisfies Assumptions 3.1–3.2 and the moment functions satisfy Assumption 3.3. Then,  $\hat{\theta}_n \xrightarrow{P} \theta_0$ .*

PROOF OF LEMMA B.6. It follows from Assumption 3.3(a) and Theorem 5.9 in van der Vaart (1998) that we only need to establish for each  $1 \leq s \leq d_\theta$ ,

$$\sup_{\theta \in \Theta} \left| \frac{1}{n} \sum_{1 \leq i \leq n} m_s(X_i, A_i, R_i, \theta) - E_P[m_s(X_i, A_i, R_i, \theta)] \right| \xrightarrow{P} 0 . \quad (51)$$

To begin, note it follows from Assumption 3.3(d) and the dominated convergence theorem that if  $m_s(x, a, r, \theta_m) \rightarrow m_s(x, a, r, \theta)$  as  $m \rightarrow \infty$  for  $\{\theta_m\} \subset \Theta^*$ , then  $E_P[m_s(X_i, A_i, R_i, \theta_m)] \rightarrow E_P[m_s(X_i, A_i, R_i, \theta)]$ . Assumption 3.3(c) then implies

$$\sup_{\theta \in \Theta} \left| \frac{1}{n} \sum_{1 \leq i \leq n} m_s(X_i, A_i, R_i, \theta) - E_P[m_s(X_i, A_i, R_i, \theta)] \right|$$

$$= \sup_{\theta \in \Theta^*} \left| \frac{1}{n} \sum_{1 \leq i \leq n} m_s(X_i, A_i, R_i, \theta) - E_P[m_s(X_i, A_i, R_i, \theta)] \right|, \quad (52)$$

which is measurable. Next, note that

$$m(X_i, A_i, R_i, \theta) = A_i m(X_i, 1, R_i(1), \theta) + (1 - A_i) m(X_i, 0, R_i(0), \theta). \quad (53)$$

and it follows from Lemma B.5 that

$$E_P[m(X_i, A_i, R_i, \theta)] = \frac{\ell}{k} E_Q[m(X_i, 1, R_i(1), \theta)] + \frac{k - \ell}{k} E_Q[m(X_i, 0, R_i(0), \theta)], \quad (54)$$

which implies that

$$\begin{aligned} & \sup_{\theta \in \Theta^*} \left| \frac{1}{n} \sum_{1 \leq i \leq n} m_s(X_i, A_i, R_i, \theta) - E_P[m_s(X_i, A_i, R_i, \theta)] \right| \\ & \leq \sup_{\theta \in \Theta^*} \left| \frac{1}{n} \sum_{1 \leq i \leq n} A_i m_s(X_i, 1, R_i(1), \theta) - \frac{\ell}{k} E[m_s(X_i, 1, R_i(1), \theta)] \right| \\ & \quad + \sup_{\theta \in \Theta^*} \left| \frac{1}{n} \sum_{1 \leq i \leq n} (1 - A_i) m_s(X_i, 0, R_i(0), \theta) - \frac{k - \ell}{k} E[m_s(X_i, 0, R_i(0), \theta)] \right|. \end{aligned}$$

We study the first term on the right-hand side and similar arguments apply to the second term.

$$\begin{aligned} & \sup_{\theta \in \Theta^*} \left| \frac{1}{n} \sum_{1 \leq i \leq n} A_i m_s(X_i, 1, R_i(1), \theta) - \frac{\ell}{k} E[m_s(X_i, 1, R_i(1), \theta)] \right| \\ & \leq \sup_{\theta \in \Theta^*} \left| \frac{1}{n} \sum_{1 \leq i \leq n} A_i (m_s(X_i, 1, R_i(1), \theta) - E[m_s(X_i, 1, R_i(1), \theta) | X_i]) \right| \end{aligned} \quad (55)$$

$$+ \sup_{\theta \in \Theta^*} \left| \frac{1}{n} \sum_{1 \leq i \leq n} (A_i - \frac{\ell}{k}) E[m_s(X_i, 1, R_i(1), \theta) | X_i] \right| \quad (56)$$

$$+ \sup_{\theta \in \Theta^*} \left| \frac{\ell}{kn} \sum_{1 \leq i \leq n} (E[m_s(X_i, 1, R_i(1), \theta) | X_i] - E[m_s(X_i, 1, R_i(1), \theta)]) \right|. \quad (57)$$

We study each term separately. First note (56) is bounded by

$$\sup_{\theta \in \Theta^*} \frac{1}{n} \sum_{1 \leq j \leq n/k} |\zeta_j^*(\theta)|,$$

where

$$\zeta_j^*(\theta) = \sup_{I \subset \lambda_j} \left\{ \frac{k - \ell}{k} \sum_{i \in I} E[m_s(X_i, 1, R_i(1), \theta) | X_i] - \frac{\ell}{k} \sum_{i \in \lambda_j \setminus I} E[m_s(X_i, 1, R_i(1), \theta) | X_i] : |I| = \ell \right\}.$$

By Assumption 3.3(e) and Assumption 3.2 we then obtain

$$\sup_{\theta \in \Theta^*} \frac{1}{n} \sum_{1 \leq j \leq n/k} |\zeta_j^*(\theta)| \lesssim \frac{1}{n} \sum_{1 \leq j \leq n/k} \max_{i, i' \in \lambda_j} \|X_i - X_{i'}\| \xrightarrow{P} 0.$$

For (57), note the class of functions

$$\{E[m_s(X, 1, R(1), \theta)|X = x] : \theta \in \Theta^*\}$$

are Lipschitz continuous in  $x$  with a uniform Lipschitz constant. It therefore follows from Corollary 4.1 in van der Vaart (1994), applied with  $I_j$ s being hypercubes, that (57) converges in probability to 0.

To analyze (55), we apply the arguments in the proof of Theorem 2.4.3 in van der Vaart and Wellner (1996) conditional on  $X^{(n)}, A^{(n)}$ . Define  $F = \sup_{\theta \in \Theta^*} m_s(X, 1, R(1), \theta)$ , which is measurable because  $\Theta^*$  is countable by Assumption 3.3(c). Define for any  $K > 0$

$$\mathcal{F}_s^K(1) = \{m_s(X, 1, R(1), \theta)I\{F \leq K\} : \theta \in \Theta^*\},$$

and

$$\mathcal{F}_s(1) = \{m_s(X, 1, R(1), \theta) : \theta \in \Theta^*\}.$$

Next, let  $\tau_i, 1 \leq i \leq n$  be a sequence of i.i.d. Rademacher random variables independent of all other variables. It follows from Markov's inequality and the symmetrization Lemma 6.3 in Ledoux and Talagrand (1991) applied conditional on  $X^{(n)}, A^{(n)}$  for the distribution

$$\bigotimes_{1 \leq i \leq n: A_i=1} P\{X_i, R_i(1)|X_i\}$$

that for every  $\epsilon > 0$ ,

$$\begin{aligned} & P\left\{ \sup_{\theta \in \Theta^*} \left| \frac{1}{n} \sum_{1 \leq i \leq n} A_i(m_s(X_i, 1, R_i(1), \theta) - E[m_s(X_i, 1, R_i(1), \theta)|X_i]) \right| > \epsilon \middle| X^{(n)}, A^{(n)} \right\} \\ & \leq \frac{1}{\epsilon} E \left[ \sup_{\theta \in \Theta^*} \left| \frac{1}{n} \sum_{1 \leq i \leq n} A_i(m_s(X_i, 1, R_i(1), \theta) - E[m_s(X_i, 1, R_i(1), \theta)|X_i]) \right| \middle| X^{(n)}, A^{(n)} \right] \\ & \leq \frac{2}{\epsilon} E_P \left[ E_\tau \left[ \sup_{\theta \in \Theta^*} \left| \frac{1}{n} \sum_{1 \leq i \leq n} \tau_i A_i m_s(X_i, 1, R_i(1), \theta) \right| \middle| X^{(n)}, A^{(n)} \right] \right] \\ & \leq \frac{2}{\epsilon} E_P \left[ E_\tau \left[ \sup_{\theta \in \Theta^*} \left| \frac{1}{n} \sum_{1 \leq i \leq n} \tau_i A_i \min\{m_s(X_i, 1, R_i(1), \theta), K\} \right| \middle| X^{(n)}, A^{(n)} \right] \right] \\ & \quad + 2E[FI\{F > K\}], \end{aligned} \tag{58}$$

where  $E_\tau[\cdot]$  should be understood as the expectation with respect to  $(\tau_i, 1 \leq i \leq n)$ , holding all else fixed. The last term could be made as small as possible by choosing  $K$  large because of Assumption 3.3(d). Next, define

$$\tilde{P}_n = \frac{1}{\eta n} \sum_{1 \leq i \leq n: A_i=1} \delta_{(X_i, R_i(1))},$$

where  $\delta$  denotes the Dirac measure. Note that  $\mathcal{F}_s(1)$  is a VC class by Assumption 3.3(b) and so  $\mathcal{F}_s^K(1)$  is a VC-class by Lemma 2.6.18(vi) in van der Vaart and Wellner (1996), and thus both totally bounded in  $L_1(\tilde{P}_n)$  by Theorem 2.6.7 in van der Vaart and Wellner (1996) (note that if  $\|F\|_{L_1(\tilde{P}_n)} = 0$  then the conditional

expectation immediately below is trivially zero, so we can without loss of generality assume  $\|F\|_{L_1(\tilde{P}_n)} > 0$ . Accordingly, define  $\mathcal{G}$  to be an  $\epsilon$ -net in  $L_1(\tilde{P}_n)$  for  $\mathcal{F}_s^K(1)$  with cardinality  $N(\epsilon, \mathcal{F}_s^K(1), L_1(\tilde{P}_n))$ . We have

$$\begin{aligned}
& E_\tau \left[ \sup_{\theta \in \Theta^*} \left| \frac{1}{\eta n} \sum_{1 \leq i \leq n} \tau_i A_i \min\{m_s(X_i, 1, R_i(1), \theta), K\} \right| \right] \\
& \leq E_\tau \left[ \sup_{f \in \mathcal{G}} \left| \frac{1}{\eta n} \sum_{1 \leq i \leq n} \tau_i A_i f(X_i, R_i(1)) \right| \right] + \epsilon \\
& \leq \sqrt{1 + \log N(\epsilon, \mathcal{F}_s^K(1), L_1(\tilde{P}_n))} \sqrt{\frac{12}{n} \sup_{f \in \mathcal{G}} \left( \int f^2 d\tilde{P}_n \right)^{1/2}} + \epsilon \\
& \leq (\sqrt{1 + \log N(\epsilon, \mathcal{F}_s^K(1), L_1(\tilde{P}_n))}) \sqrt{\frac{12}{n}} K + \epsilon \\
& \leq (\sqrt{1 + \log N(\epsilon, \mathcal{F}_s(1), L_1(\tilde{P}_n))}) \sqrt{\frac{12}{n}} K + \epsilon \\
& \lesssim \left( \sqrt{1 + (V-1) \log \left( \frac{\|F\|_{L_1(\tilde{P}_n)}}{\epsilon} \right)} \sqrt{\frac{12}{n}} K + \epsilon \right) \wedge K, \tag{60}
\end{aligned}$$

where the second inequality follows from Lemma 2.2.2 in [van der Vaart and Wellner \(1996\)](#) applied with  $\exp(x^2) - 1$  and Hoeffding's lemma, the third follows because  $|f| \leq K$  for  $f \in \mathcal{G}$ , the fourth inequality follows from the fact that  $N(\epsilon, \mathcal{F}_s^K(1), L_1(\tilde{P}_n)) \leq N(\epsilon, \mathcal{F}_s(1), L_1(\tilde{P}_n))$  because

$$\int |f_1 I\{F \leq K\} - f_2 I\{F \leq K\}| d\tilde{P}_n = \int |f_1 - f_2| I\{F \leq K\} d\tilde{P}_n \leq \int |f_1 - f_2| d\tilde{P}_n, \tag{61}$$

and the last inequality follows from Theorem 2.6.7 in [van der Vaart and Wellner \(1996\)](#) and Assumption [3.3\(b\)](#).

Note it follows from Assumptions [3.1-3.2](#), [3.3\(d\)-\(f\)](#), and similar arguments to those in the first part of the proof of Lemma S.1.5 in [Bai et al. \(2022\)](#) that

$$E[\|F\|_{L_1(\tilde{P}_n)} | X^{(n)}, A^{(n)}] = \frac{1}{\eta n} \sum_{1 \leq i \leq n} I\{A_i = 1\} E[|F| | X_i] \xrightarrow{P} E[|F|]. \tag{62}$$

We can assume without loss of generality  $E[|F|] > 0$  because otherwise  $m_s(x, 1, r(1), \theta) \equiv 0$ . Therefore,

$$P \left\{ E[\|F\|_{L_1(\tilde{P}_n)} | X^{(n)}, A^{(n)}] > E[|F|] + \frac{E[|F|]}{2} \right\} \rightarrow 0. \tag{63}$$

On the other hand, Assumptions [3.1-3.2](#), [3.3\(d\)-\(f\)](#) and similar arguments to those in the last part of the proof of Lemma S.1.5 in [Bai et al. \(2022\)](#) that

$$P \left\{ \left| \|F\|_{L_1(\tilde{P}_n)} - E[\|F\|_{L_1(\tilde{P}_n)} | X^{(n)}, A^{(n)}] \right| > \frac{E[|F|]}{2} \mid X^{(n)}, A^{(n)} \right\} \xrightarrow{P} 0. \tag{64}$$



Let

$$\mathbb{L}_n = P \left\{ \sup_{\theta \in \Theta^*} \left| \frac{1}{n} \sum_{1 \leq i \leq n} A_i(m_s(X_i, 1, R_i(1), \theta) - E[m_s(X_i, 1, R_i(1), \theta) | X_i]) \right| > \epsilon \mid X^{(n)}, A^{(n)} \right\}.$$

To conclude the proof, note for every  $\eta > 0$ , for  $n$  large enough,

$$\begin{aligned} & P \{ \mathbb{L}_n > \eta \} \\ & \leq P \left\{ E[\|F\|_{L_1(\bar{P}_n)} | X^{(n)}, A^{(n)}] > \frac{3}{2} E\|F\| \right\} \\ & \quad + P \left\{ \mathbb{L}_n > \eta, E[\|F\|_{L_1(\bar{P}_n)} | X^{(n)}, A^{(n)}] \leq \frac{3}{2} E\|F\| \right\} \\ & \leq P \left\{ E[\|F\|_{L_1(\bar{P}_n)} | X^{(n)}, A^{(n)}] > \frac{3}{2} E\|F\| \right\} \\ & \quad + P \left\{ E \left[ \left( \sqrt{1 + (V-1) \log \left( \frac{\|F\|_{L_1(\bar{P}_n)}}{\epsilon} \right)} \sqrt{\frac{12}{n}} K + \epsilon \right) \wedge K \mid X^{(n)}, A^{(n)} \right] > \eta' , \right. \\ & \quad \left. E[\|F\|_{L_1(\bar{P}_n)} | X^{(n)}, A^{(n)}] \leq \frac{3}{2} E\|F\| \right\} \\ & \leq P \left\{ E[\|F\|_{L_1(\bar{P}_n)} | X^{(n)}, A^{(n)}] > \frac{3}{2} E\|F\| \right\} \\ & \quad + P \left\{ P \left[ \left| \|F\|_{L_1(\bar{P}_n)} - E[\|F\|_{L_1(\bar{P}_n)} | X^{(n)}, A^{(n)}] \right| > \frac{E\|F\|}{2} \mid X^{(n)}, A^{(n)} \right] > \eta'' , \right. \\ & \quad \left. E[\|F\|_{L_1(\bar{P}_n)} | X^{(n)}, A^{(n)}] \leq \frac{3}{2} E\|F\| \right\} \xrightarrow{P} 0, \end{aligned}$$

where  $\eta', \eta''$  are suitably chosen constants, the last line follows from the law of total expectation combined with the fact that the quantity in the expectation is bounded by  $K$ , and the convergence follows from (63)–(64). ■

**Lemma B.7.** For  $f : \mathbf{R}^{d_x} \rightarrow \mathbf{R}$ , define

$$\begin{aligned} h_f(x_1, \dots, x_k) = & \sup_{I \subset \{1, \dots, k\}, |I|=\ell} \left( \frac{1}{\ell} \sum_{i \in I} f(x_i) - \frac{1}{k-\ell} \sum_{i \in \{1, \dots, k\} \setminus I} f(x_i) \right) \\ & - \inf_{I \subset \{1, \dots, k\}, |I|=\ell} \left( \frac{1}{\ell} \sum_{i \in I} f(x_i) - \frac{1}{k-\ell} \sum_{i \in \{1, \dots, k\} \setminus I} f(x_i) \right). \end{aligned}$$

Then,

$$|h_f(x_1, \dots, x_k) - h_g(x_1, \dots, x_k)|^2 \lesssim \sum_{1 \leq i \leq k} |f(x_i) - g(x_i)|^2$$

PROOF. Suppose the supremum and infimum in the definition of  $h_f$  are attained at  $I^*$  and  $I_*$ . Then,

$$\begin{aligned} h_f - h_g \leq & \left( \frac{1}{\ell} \sum_{i \in I^*} f(x_i) - \frac{1}{k-\ell} \sum_{i \in \{1, \dots, k\} \setminus I^*} f(x_i) \right) - \left( \frac{1}{\ell} \sum_{i \in I^*} g(x_i) - \frac{1}{k-\ell} \sum_{i \in \{1, \dots, k\} \setminus I^*} g(x_i) \right) \\ & + \left( \frac{1}{\ell} \sum_{i \in I_*} f(x_i) - \frac{1}{k-\ell} \sum_{i \in \{1, \dots, k\} \setminus I_*} f(x_i) \right) - \left( \frac{1}{\ell} \sum_{i \in I_*} g(x_i) - \frac{1}{k-\ell} \sum_{i \in \{1, \dots, k\} \setminus I_*} g(x_i) \right), \end{aligned}$$

and the result follows from repeated applications of the inequality  $(a + b)^2 \leq 2(a^2 + b^2)$ . ■

**Lemma B.8.** *If  $X_n \xrightarrow{P} 0$  and  $|X_n| \leq X$  with  $E[X] < \infty$ , then  $E[X_n] \rightarrow 0$ .*

PROOF. Suppose not. Then along a subsequence  $\{n_k\}$ ,  $E[|X_{n_k}|] \rightarrow \delta > 0$ . Because  $X_n \xrightarrow{P} 0$ , there exists a further subsequence along which  $X_{n_{k_\ell}} \rightarrow 0$  with probability one, and by the dominated convergence theorem  $E[X_{n_{k_\ell}}] \rightarrow 0$ , a contradiction. ■

## References

- ABADIE, A. and IMBENS, G. W. (2008). Estimation of the Conditional Variance in Paired Experiments. *Annales d'Économie et de Statistique*, **1** 175–187.
- ABADIE, A. and IMBENS, G. W. (2016). Matching on the Estimated Propensity Score. *Econometrica*, **84** 781–807.
- 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.
- ARMSTRONG, T. B. (2022). Asymptotic Efficiency Bounds for a Class of Experimental Designs. ArXiv:2205.02726 [stat], URL <http://arxiv.org/abs/2205.02726>.
- BAI, Y. (2022). Optimality of Matched-Pair Designs in Randomized Controlled Trials. *American Economic Review*, **112** 3911–3940.
- BAI, Y., GUO, H., SHAIKH, A. M. and TABORD-MEEHAN, M. (2024a). Inference in Experiments with Matched Pairs and Imperfect Compliance. ArXiv:2307.13094 [econ, math, stat], URL <http://arxiv.org/abs/2307.13094>.
- BAI, Y., JIANG, L., ROMANO, J. P., SHAIKH, A. M. and ZHANG, Y. (2024b). Covariate adjustment in experiments with matched pairs. *Journal of Econometrics*, **241** 105740.
- BAI, Y., LIU, J., SHAIKH, A. M. and TABORD-MEEHAN, M. (2024c). Inference in Cluster Randomized Trials with Matched Pairs. ArXiv:2211.14903 [econ, stat], URL <http://arxiv.org/abs/2211.14903>.
- BAI, Y., LIU, J. and TABORD-MEEHAN, M. (2024d). Inference for Matched Tuples and Fully Blocked Factorial Designs. *Quantitative Economics*, **15** 279–330.
- BAI, Y., ROMANO, J. P. and SHAIKH, A. M. (2022). Inference in Experiments With Matched Pairs. *Journal of the American Statistical Association*, **117** 1726–1737.
- 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.
- BELLONI, A., CHERNOZHUKOV, V., FERNANDEZ-VAL, I. and HANSEN, C. (2017). Program evaluation and causal inference with high-dimensional data. *Econometrica*, **85** 233–298.
- 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.
- BUGNI, F., CANAY, I., SHAIKH, A. and TABORD-MEEHAN, M. (2022). Inference for Cluster Randomized Experiments with Non-ignorable Cluster Sizes. ArXiv:2204.08356 [econ, stat], URL <http://arxiv.org/abs/2204.08356>.

- BUGNI, F. A., CANAY, I. A. and SHAIKH, A. M. (2019). Inference under covariate-adaptive randomization with multiple treatments. *Quantitative Economics*, **10** 1747–1785.
- CAI, Y. and RAFI, A. (2022). On the Performance of the Neyman Allocation with Small Pilots. ArXiv:2206.04643 [econ], URL <http://arxiv.org/abs/2206.04643>.
- CASABURI, L. and REED, T. (2022). Using Individual-Level Randomized Treatment to Learn about Market Structure. *American Economic Journal: Applied Economics*, **14** 58–90.
- CHEN, X., HONG, H. and TAROZZI, A. (2008). Semiparametric Efficiency in GMM Models with Auxiliary Data. *The Annals of Statistics*, **36** 808–843.
- CHEN, X. and SANTOS, A. (2018). Overidentification in Regular Models. *Econometrica*, **86** 1771–1817.
- CHERNOZHUKOV, V., CHETVERIKOV, D., DEMIRER, M., DUFLO, E., HANSEN, C. and NEWEY, W. (2017). Double/Debiased/Neyman Machine Learning of Treatment Effects. *American Economic Review*, **107** 261–265.
- CHERNOZHUKOV, V., CHETVERIKOV, D., DEMIRER, M., DUFLO, E., HANSEN, C., NEWEY, W. and ROBINS, J. (2018). Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal*, **21** C1–C68.
- CHERNOZHUKOV, V., CINELLI, C., NEWEY, W., SHARMA, A. and SYRGGANIS, V. (2024). Long Story Short: Omitted Variable Bias in Causal Machine Learning. ArXiv:2112.13398 [cs, econ, stat], URL <http://arxiv.org/abs/2112.13398>.
- CYTRYNBAUM, M. (2023a). Covariate Adjustment in Stratified Experiments. ArXiv:2302.03687 [econ, stat], URL <http://arxiv.org/abs/2302.03687>.
- CYTRYNBAUM, M. (2023b). Designing representative and balanced experiments by local randomization. *arXiv preprint arXiv:2111.08157*.
- DIZON-ROSS, R. (2019). Parents’ beliefs about their children’s academic ability: Implications for educational investments. *American Economic Review*, **109** 2728–65.
- DUDLEY, R. M. (2014). *Uniform Central Limit Theorems*. 2nd ed. Cambridge Studies in Advanced Mathematics, Cambridge University Press, Cambridge.
- DUFLO, E., DUPAS, P. and KREMER, M. (2015). Education, hiv, and early fertility: Experimental evidence from kenya. *American Economic Review*, **105** 2757–2797.
- DUFLO, E., GLENNERSTER, R. and KREMER, M. (2007). Using randomization in development economics research: A toolkit. *Handbook of development economics*, **4** 3895–3962.
- FARRELL, M. H. (2015). Robust inference on average treatment effects with possibly more covariates than observations. *Journal of Econometrics*, **189** 1–23.
- FIRPO, S. (2007). Efficient Semiparametric Estimation of Quantile Treatment Effects. *Econometrica*, **75** 259–276.

- FREEDMAN, D. A. (2008). On regression adjustments to experimental data. *Advances in Applied Mathematics*, **40** 180–193.
- FROLICH, M. (2007). Nonparametric IV estimation of local average treatment effects with covariates. *Journal of Econometrics*, **139** 35–75.
- HAHN, J. (1998). On the Role of the Propensity Score in Efficient Semiparametric Estimation of Average Treatment Effects. *Econometrica*, **66** 315–331.
- HECKMAN, J. J., ICHIMURA, H. and TODD, P. (1998). Matching as an econometric evaluation estimator. *The review of economic studies*, **65** 261–294.
- HIRANO, K. and IMBENS, G. W. (2001). Estimation of Causal Effects using Propensity Score Weighting: An Application to Data on Right Heart Catheterization. *Health Services and Outcomes Research Methodology*, **2** 259–278.
- HIRANO, K., IMBENS, G. W. and RIDDER, G. (2003). Efficient Estimation of Average Treatment Effects Using the Estimated Propensity Score. *Econometrica*, **71** 1161–1189.
- IMBENS, G., NEWEY, W. and RIDDER, G. (2007). Mean-squared-error calculations for average treatment effects. *working paper*.
- IMBENS, G. W. (2004). Nonparametric Estimation of Average Treatment Effects Under Exogeneity: A Review. *The Review of Economics and Statistics*, **86** 4–29.
- IMBENS, G. W. and ANGRIST, J. D. (1994). Identification and Estimation of Local Average Treatment Effects. *Econometrica*, **62** 467–475.
- JIANG, L., LINTON, O. B., TANG, H. and ZHANG, Y. (2022a). Improving Estimation Efficiency via Regression-Adjustment in Covariate-Adaptive Randomizations with Imperfect Compliance. ArXiv:2201.13004 [econ, stat], URL <http://arxiv.org/abs/2201.13004>.
- JIANG, L., LIU, X., PHILLIPS, P. C. and ZHANG, Y. (2021). Bootstrap inference for quantile treatment effects in randomized experiments with matched pairs. *Review of Economics and Statistics* 1–47.
- JIANG, L., PHILLIPS, P. C. B., TAO, Y. and ZHANG, Y. (2022b). Regression-Adjusted Estimation of Quantile Treatment Effects under Covariate-Adaptive Randomizations. ArXiv:2105.14752 [econ, stat], URL <http://arxiv.org/abs/2105.14752>.
- KOSOROK, M. R. (2008). *Introduction to Empirical Processes and Semiparametric Inference*. Springer Series in Statistics, Springer-Verlag, New York.
- KRANTZ, S. G. and PARKS, H. R. (2013). *The Implicit Function Theorem: History, Theory, and Applications*. Springer, New York, NY.
- LEDOUX, M. and TALAGRAND, M. (1991). *Probability in Banach Spaces: Isoperimetry and Processes*. Classics in Mathematics, Springer-Verlag, Berlin Heidelberg.

- LEHMANN, E. and ROMANO, J. P. (2022). *Testing Statistical Hypotheses*. Springer Texts in Statistics, Springer International Publishing, Cham.
- LIN, W. (2013). Agnostic notes on regression adjustments to experimental data: Reexamining Freedman's critique. *Annals of Applied Statistics*, **7** 295–318.
- NEWBY, W. K. (1994). The Asymptotic Variance of Semiparametric Estimators. *Econometrica*, **62** 1349–1382.
- NEWBY, W. K. and MCFADDEN, D. (1994). Large sample estimation and hypothesis testing. *Handbook of econometrics*, **4** 2111–2245.
- RAFI, A. (2023). Efficient semiparametric estimation of average treatment effects under covariate adaptive randomization. *arXiv preprint arXiv:2305.08340*.
- ROBINS, J. M., ROTNITZKY, A. and ZHAO, L. P. (1995). Analysis of Semiparametric Regression Models for Repeated Outcomes in the Presence of Missing Data. *Journal of the American Statistical Association*, **90** 106–121.
- ROSENBERGER, W. F. and LACHIN, J. M. (2015). *Randomization in clinical trials: theory and practice*. John Wiley & Sons.
- TSIATIS, A. A. (2006). *Semiparametric Theory and Missing Data*. Springer Series in Statistics, Springer, New York, NY.
- TSIATIS, A. A., DAVIDIAN, M., ZHANG, M. and LU, X. (2008). Covariate adjustment for two-sample treatment comparisons in randomized clinical trials: a principled yet flexible approach. *Statistics in Medicine*, **27** 4658–4677.
- VAN DER VAART, A. (1994). Bracketing smooth functions. *Stochastic Processes and their Applications*, **52** 93–105.
- VAN DER VAART, A. W. (1998). *Asymptotic statistics*, vol. 3 of *Cambridge Series in Statistical and Probabilistic Mathematics*. Cambridge University Press, Cambridge.
- VAN DER VAART, A. W. and WELLNER, J. (1996). *Weak Convergence and Empirical Processes: With Applications to Statistics*. Springer Series in Statistics, Springer-Verlag, New York.
- VERSHYNIN, R. (2018). *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge Series in Statistical and Probabilistic Mathematics, Cambridge University Press, Cambridge.
- ZHANG, M., TSIATIS, A. A. and DAVIDIAN, M. (2008). Improving efficiency of inferences in randomized clinical trials using auxiliary covariates. *Biometrics*, **64** 707–715.