

# Revisiting Randomization with the Cube Method

Laurent Davezies\*    Guillaume Hollard†    Pedro Vergara Merino‡

## Abstract

We propose a novel randomization approach for randomized controlled trials (RCTs), named the cube method. The cube method allows for the selection of balanced samples across various covariate types, ensuring consistent adherence to balance tests and, whence, substantial precision gains when estimating treatment effects. We establish several statistical properties for the population and sample average treatment effects (PATE and SATE, respectively) under randomization using the cube method. The relevance of the cube method is particularly striking when comparing the behavior of prevailing methods employed for treatment allocation when the number of covariates to balance is increasing. We formally derive and compare bounds of balancing adjustments depending on the number of units  $n$  and the number of covariates  $p$  and show that our randomization approach outperforms methods proposed in the literature when  $p$  is large and  $p/n$  tends to 0. We run simulation studies to illustrate the substantial gains from the cube method for a large set of covariates.

**Keywords:** Causal inference, covariate balance, experimental design, treatment effects.

**JEL Codes:** C13, C21

---

\*CREST - ENSAE - IP Paris, laurent.davezies@ensae.fr

†CREST - Ecole Polytechnique - IP Paris, guillaume.hollard@polytechnique.edu

‡CREST - ENSAE - IP Paris, pedro.vergaramerino@ensae.fr

We would like to thank Riccardo d’Adamo, Isaiah Andrews, Russell Davidson, Xavier d’Haultfoeuille, Xinran Li, Pauline Rossi, Felix Schlee, Sami Stouli, and many other colleagues in workshops and conferences whose comments and suggestions helped improve this paper. In particular, we would like to thank everybody at the 2023 Bristol Econometric Study Group, the 2023 Advances with Field Experiments Conference, the 18th IZA & 5th IZA/CREST Conference, the 2023 European Winter Meeting of the Econometric Society, the AMSE Big Data and Econometrics Seminar, the 2024 RCEA International Conference in Economics, Econometrics, and Finance, the 2024 BSE Summer Forum Workshop on Microeconometrics and Policy Evaluation, the 2024 AFSE Annual Conference, and the 2024 Annual Conference of the International Association for Applied Econometrics.

# 1 Introduction

Randomized controlled trials (RCTs) substantially differ in the number of available covariates used at the randomization stage. Considering RCTs published in top-five journals during the past five years, 34% do not use covariates for randomization, i.e., completely randomized designs, even if some covariates were available. Most studies, 54%, employ a stratified design for allocation, choosing a set of  $p$  baseline covariates to be “balanced” between treatment and control. For a large sample size  $n$ , researchers would like to stratify using all available covariates. For  $n$  fixed, however, “imbalances” become more frequent as  $p$  grows, limiting the relevance of stratified designs with large  $p$ . This curse of dimensionality in  $p$  arises in all stratification methods, including the matched-pairs design (Bai et al., 2022; Bai, 2022) and its recent generalization by Cytrynbaum (2023).

To attenuate this curse of dimensionality, we here extend the “cube method”, initially developed by Deville and Tillé (2004) for sampling, to the experimental framework. As we will show, the cube method improves randomization in RCTs on several key dimensions.

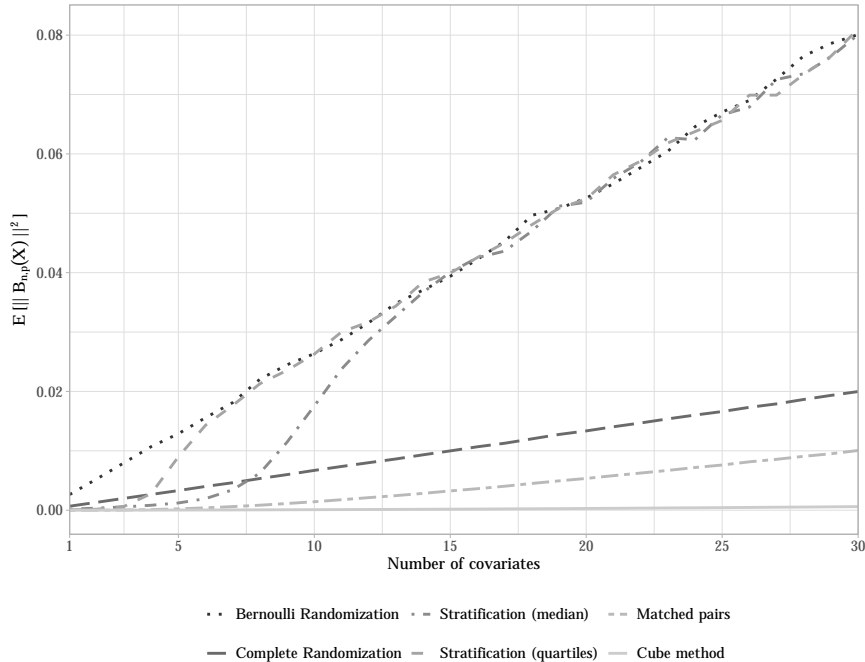
Let us start with a motivating example to illustrate how the number of covariates affects imbalance and how the cube method stands out. Consider an empiricist who wants to balance  $p$  covariates  $X_i = (X_{1i}, \dots, X_{pi})'$  between a treated and control groups for a given sample of units  $i = 1, \dots, n$ . We consider the simple case where  $(X_i)_{i=1, \dots, n}$  are drawn independently from uniform distributions on  $[0; 1]^p$ , and every individual is treated with probability  $1/2$ . The degree of imbalance between control group  $C$  and treatment group  $T$  is here measured by the squared Euclidean norm of  $B_{n,p}(X) = \frac{2}{n} \sum_{i \in T} X_i - \frac{2}{n} \sum_{i \in C} X_i$ .<sup>1</sup> For different randomization methods, we compute Monte Carlo estimates of  $\mathbb{E}[\|B_{n,p}(X)\|^2]$  for  $p = 1, \dots, 30$  and  $n = 500$ .

Figure 1 shows that stratification indeed limits imbalances for small values of  $p$  but imbalances become more frequent as  $p$  grows. For  $p > 15$  stratification provides worse results than complete randomization. Recent developments, such as the matched-pairs design, limit imbalances because of their ability to mitigate the small strata issue (i.e., situations in which strata are empty or contain a single observation). The cube method allows spectacular improvements and exhibits a different dynamic with a linear growth of the squared Euclidean norm. The cube method exhibits a different dynamic when  $p$  grows because it relies on a different notion of “balancing”. Indeed, stratification and matched-pairs design aim at making

---

<sup>1</sup>We here use the Horvitz-Thompson estimator to compute average differences between control and treatment. When the number of units in each group is fixed, as for complete randomization, matched-pairs design, or the cube method, this estimator is equivalent to a difference-in-means.

Figure 1: Balance Quality and Number of Covariates



the *joint probability distribution* of selected covariates as similar as possible between treatment and control. In sharp contrast, the cube method balances a set of selected *moments*. Balancing numerous covariates is crucial as precision gains in treatment estimates depend on picking out the most explicative ones. Empiricists are typically unsure which covariates to include, especially if they are interested in several treatment outcomes. Selecting more covariates, without facing the small strata issue, will help empiricists increase precision and thus reduce the cost of RCTs as the number of required units to achieve the same precision will be lower.

We here contribute to four streams of research. The cube method is a sampling method designed by Deville and Tillé (2004) to obtain balanced samples so that Horvitz–Thompson estimators of the population totals of a set of covariates equal the known totals of these variables. The cube method is a notable achievement that selects approximately balanced samples with equal or unequal inclusion probabilities and uses numerous covariates. The cube method is routinely used by national statistical institutes (see Tillé, 2011, for a review of applications of the cube method). Our technical contribution is to extend the scope of the cube method beyond sampling for estimating treatment effects in RCTs. Our contribution includes, among others, deriving asymptotic properties for estimators of the population average treatment effect (PATE) and sample average treatment effect (SATE). For the PATE, we

show how to perform inference and that the semiparametric efficiency bound in Hahn (1998) is attained. We also provide an explicit formula for asymptotic precision gains, showing that they are higher when balancing covariates more correlated to potential outcomes. We thus formally motivate the interest of increasing the number of covariates used for randomization.

Our second contribution is to provide a more general comparison of the behavior of randomization methods when  $p$ , the number of covariates used for balancing, increases. Asymptotic properties of randomization methods when the number of units  $n$  is getting larger are well studied already. In contrast, the asymptotic behavior when  $p$  grows has not been surveyed yet, to the best of our knowledge. Our introductory example provided already insights into this crucial question, illustrating an essential difference between the cube and other methods. We show that observed patterns are de facto generic and, thus, not specific to a particular example. In particular, we formally derive the relation between the number of units  $n$  and the number of covariates  $p$ . Our main result is to uncover three areas where stratification behaves differently. For a small number of covariates (i.e.,  $p \ll \ln(n)$ ), stratification improves balancing compared to complete randomization. When  $p \approx \ln(n)$ , there is a critical regime where the balancing quality deteriorates quickly. Last, for  $p \gg \ln(n)$ , because of the small strata issue, stratification is similar to Bernoulli randomization and worse than complete randomization. In sum, stratification exhibits different balancing properties when  $p$  varies. In sharp contrast, the cube method grants the balance of (the moments of) selected covariates with no critical change when  $p$  grows, allowing balancing on a large set of covariates. By construction, the cube method further guarantees the balance between treatment and control on selected covariates *with probability one* (up to an asymptotically negligible rounding term). Complete randomization only guarantees that balancing occurs *on average*, which means that in some instances, some imbalances occur as  $p$  grows. Some RCTs will thus not pass balance checks, which is shown to increase rejection probability by economic journals, creating a publication bias (see Snyder and Zhuo, 2024, for evidence and discussion). A side effect of using the cube method could thus be to reduce publication bias by avoiding large imbalances.

We also contribute to a stream of recent papers that propose new randomization techniques to increase precision gains when estimating treatment effects. Bai et al. (2024) propose a stimulating review of recent developments. An essential question in this literature is how randomization methods relate to the precision of treatment effects estimates. For instance, there are several ways of using stratification to create a control and a treatment group. Bugni et al. (2018) provides several estimators with associated exact confidence intervals, allowing

exploitation of precision gains obtained by stratification. Bai (2022) shows that a specific matched-pairs design achieves the maximum statistical precision for estimating the average treatment effect using the difference in means estimator. Figure 1 shows an example of the sizeable improvements of matched-pairs designs. Cytrynbaum (2023) extends this result and shows that precision gains are more substantial when using covariates more predictive of treatment effect heterogeneity. A crucial question in this literature is thus to identify the most relevant covariates to use at the randomization stage. A pivotal advantage of the cube method is to allow the inclusion of a large set of covariates, thus maximizing the probability of including the relevant ones (i.e., the ones that correlate to potential outcomes). Put differently, by providing precision gains, the cube achieves the same precision with a smaller sample, which diminishes the cost of RCTs by reducing the number of surveyed units. We gauge the size of precision gains using both Monte Carlo simulations and an empirical example using real world data.

When choosing a particular randomization method and estimating treatment effects, experimenters now have the choice among several methods and associated estimators. Building on the seminal work of Bruhn and McKenzie (2009), Athey and Imbens (2017) that provide a systematic review of the pros and cons of each method, Bai et al. (2024) review the most recent developments. Our last and modest contribution emphasizes five properties of the cube method that experimenters may want to consider when choosing a randomization method. First, the cube method removes most of the bad luck that may arise from sampling errors, as the most unfavorable samples have a null probability of being selected. As explained, the same desirable property is granted by stratification designs but only for a “small” number of covariates. Second, the cube method requires no particular assumption regarding the probability of being allocated to the treatment group (e.g., probabilities do not need to be equal to  $1/2$ ), and probabilities may differ across units (e.g., inclusion probabilities may depend on some covariates). Freedom of choosing assignment probabilities is desirable when the experimenter wants to target a specific group (for instance, oversampling a group because of anticipated attrition). Third, the cube method does not require the choice of any tuning parameter (unlike, for example, the Gram-Schmidt walk design). The experimenter only needs to choose the set of moments on which she wishes to balance control and treatment groups. This feature avoids ad-hoc choices regarding parameters. Fourth, the cube method is computationally simple and can rapidly run on any computer, even for large samples and numerous covariates. Five, there already exist packages allowing the interested reader to run the cube algorithm. Six, we show through simulations that balancing using the cube method

provides very similar estimates as the double-lasso regression (Belloni et al., 2014), which uses regularization techniques to estimate treatment effects after running the experiment. While both methods provide very similar estimates- there is a reason to prefer the cube over the double lasso. As shown by Kolesár et al. (2024), estimation of treatment effect using lasso regressions is prone to instability arising from seemingly innocuous choices (e.g., centering variables around mean vs median, choosing a reference category). The cube method does not involve any data formatting, thus avoiding a potential source of instability.

Section 2 introduces the potential outcome framework and covariate balancing. Section 3 presents the cube algorithm and its application to RCTs. Section 4 gives the balancing properties of the cube method, compares imbalances to other methods, and provides novel asymptotic expressions for the variance of average treatment effect estimators. We then specify two ways of performing inference based on asymptotic normality and the randomization mechanism. Section 5 uses simulated and experimental data to show our precision gains and how the cube method might be less constrained by the curse of dimensionality. Finally, Section 6 reviews current practices in RCTs and discusses practical considerations of the cube method.

## 2 Setup

This section presents the potential outcomes framework, provides assumptions on the data-generating process, and formally defines covariate balancing.

### 2.1 Data Generating Process and Assignment Design

We consider the standard Neyman-Rubin framework of potential outcomes where  $Y_i(0)$  is the outcome of unit  $i$  when untreated and  $Y_i(1)$  is the outcome when treated. We consider  $X_i$  a vector of  $p$  covariates. According to the literature, we assume i.i.d.ness and the existence of second-order moments for all these variables.

#### Assumption 1

$(Y_i(0), Y_i(1), X_i)$  are *i.i.d.* across  $i$  and  $\mathbb{E}(Y(0)^2 + Y(1)^2 + \|X\|^2) < \infty$

The empiricist observes  $(X_1, \dots, X_n)$  for a finite sample of size  $n$ . She wants to randomly allocate these  $n$  units to treatment according to a design  $\Pi$ , i.e., a distribution on the set of

the possible treatment allocations  $\{0, 1\}^n$ . If the design  $\Pi$  does not depend on the potential outcomes, it balances potential outcomes in the treatment and control groups *in average*, avoiding selection bias. The design  $\Pi$  could depend on  $(X_1, \dots, X_n)$ . For instance, the treatment probability of a unit  $i$  could depend on  $X_i$  for various reasons, such as efficiency, cost of the treatment depending on  $X_i$ , or subpopulations of particular interest. In the following,  $D_i$  is the dummy variable indicating if  $i$  is treated or untreated. Empiricists have to choose not only each individual selection probability  $\mathbb{P}_\Pi(D_i = 1|X_1, \dots, X_n)$  but the full design  $\Pi$  that determines  $\mathbb{P}_\Pi(\cap_{i=1, \dots, n} D_i = d_i|X_1, \dots, X_n)$  for any potential allocations  $(d_i)_{i=1, \dots, n} \in \{0, 1\}^n$ . A major issue is exploiting the knowledge of  $(X_1, \dots, X_n)$  to define a “good” design  $\Pi$  to go beyond the balancing of potential outcomes in average. To study this question, let us formulate the assumption on the class of design we consider in the following.

**Assumption 2** *The empiricist observes a sample  $(X_i)_{i=1, \dots, n}$  of size  $n$  and generates a random assignment  $(D_i)_{i=1, \dots, n}$  according to a randomization design  $\Pi$  such that:*

$$(D_1, \dots, D_n) \perp\!\!\!\perp (Y_1(0), Y_1(1), \dots, Y_n(0), Y_n(1)) | X_1, \dots, X_n \quad (1)$$

and for any  $i = 1, \dots, n$

$$\mathbb{P}_\Pi(D_i = 1 | X_1, \dots, X_n) = p(X_i) \in [c, 1 - c], \quad (2)$$

where  $p$  is a function chosen by the empiricist and  $c$  is a positive constant.

Equation (1) means that assignment is independent of the unknown potential outcomes, conditional on the auxiliary information  $X$ . Equation (2) specifies that the assignment probability of unit  $i$  could depend on  $X_i$  but not on  $X_j$  for  $j \neq i$ . It also states that the propensity score  $p(X_i)$  fulfills a common support condition. This restriction is usual and necessary in the literature on treatment effects estimation. In many RCTs,  $p(X_i) = 1/2$  for all  $i$ , and in that case, there is (on average) the same number of treated and untreated units. But in some cases, for theoretical reasons (e.g., efficiency, population of interest) or practical reasons (e.g., budget constraints), only a smaller fraction of units could be treated ( $p(X_i) < 1/2$ ) or  $p(X_i)$  could be heterogeneous across  $i$  ( $\mathbb{V}(p(X_i)) > 0$ ). In the following, we denote the propensity score  $p(X_i)$  as  $\pi_i$ . Our proposition of design accommodates any propensity score type, offering complete flexibility to empiricists concerning its definition.

After the experiment, the empiricist observes  $Y_i = Y_i(1) \times D_i + Y_i(0) \times (1 - D_i)$ . She will thus never observe both potential outcomes for the same unit.

Empiricists are generally interested in estimating the sample and population average treat-

ment effects given by

$$\text{SATE : } \theta_0 = \frac{1}{n} \sum_{i=1}^n Y_i(1) - Y_i(0) \quad (3)$$

and

$$\text{PATE : } \theta_0^* = \mathbb{E}[Y_i(1) - Y_i(0)], \quad (4)$$

respectively.<sup>2</sup>

In this paper, we will focus on the Horvitz-Thompson estimator (HT) and the Hájek estimator (H), which are of central interest in RCTs. The Horvitz-Thompson estimator is

$$\hat{\theta}_{HT} = \frac{1}{n} \sum_{i=1}^n \left( \frac{Y_i D_i}{\pi_i} - \frac{Y_i(1 - D_i)}{1 - \pi_i} \right) \quad (5)$$

which is unbiased for both the SATE and the PATE and is the difference between the inverse probability weighting estimators on the treated and the control group.

The Hájek estimator is

$$\hat{\theta}_H = \frac{1}{\sum_{i=1}^n \frac{D_i}{\pi_i}} \sum_{i=1}^n \frac{Y_i D_i}{\pi_i} - \frac{1}{\sum_{i=1}^n \frac{1-D_i}{1-\pi_i}} \sum_{i=1}^n \frac{Y_i(1 - D_i)}{1 - \pi_i} \quad (6)$$

and corresponds as well to the inverse probability weighting OLS estimator

$$\hat{\theta}_H = \arg \min_{\theta} \min_a \sum_{i=1}^n w_i (Y_i - a - \theta D_i)^2$$

for  $w_i = \frac{1}{\pi_i}$  if  $D_i = 1$  and  $w_i = \frac{1}{1-\pi_i}$  if  $D_i = 0$ . Let  $n_T$  denote the number of treated units and  $n_C$  the number of control units. When  $\pi_i$  is constant,  $\hat{\theta}_H = \frac{1}{n_T} \sum_{i:D_i=1} Y_i - \frac{1}{n_C} \sum_{i:D_i=0} Y_i$  is the difference between the average on the treated group and the control group whereas  $\hat{\theta}_{HT} = \frac{1}{\mathbb{E}(n_T)} \sum_{i:D_i=1} Y_i - \frac{1}{\mathbb{E}(n_C)} \sum_{i:D_i=0} Y_i$  is a slight modification of this difference of averages. If, additionally,  $n_T$  and  $n_C$  are fixed, both estimators are identical to the difference-in-means estimator.

## 2.2 Balancing Constraints

Under Assumption 2, as soon as  $\mathbb{E}(D_i | (X_{i'}, Y_{i'}(0), Y_{i'}(1))_{i'=1, \dots, n}) = \pi_i$  and we have average balance for the potential outcomes, i.e.,

$$\mathbb{E} \left( \frac{1}{n} \sum_{i:D_i=1} \frac{Y_i(1)}{\pi_i} \middle| (X_{i'})_{i'=1, \dots, n} \right) = \frac{1}{n} \sum_{i=1}^n Y_i(1) \text{ and}$$

---

<sup>2</sup>In some cases, they are interested in similar parameters for some subpopulations:  $\frac{1}{\sum_{i=1}^n \mathbb{1}\{X_i \in \mathcal{X}\}} \sum_{i=1}^n (Y_i(1) - Y_i(0)) \mathbb{1}\{X_i \in \mathcal{X}\}$  or  $\mathbb{E}[Y_i(1) - Y_i(0) | X_i \in \mathcal{X}]$ . Estimators of these quantities are defined restricting the sample to units such that  $X_i \in \mathcal{X}$  and the asymptotic properties of these estimators follow from a straightforward adaptation of what is presented in the following.



$$\mathbb{E} \left( \frac{1}{n} \sum_{i:D_i=0} \frac{Y_i(0)}{1 - \pi_i} \middle| (X_{i'})_{i'=1, \dots, n} \right) = \frac{1}{n} \sum_{i=1}^n Y_i(0),$$

for any covariates  $X_j$  for  $j = 1, \dots, p$ :

$$\begin{aligned} \mathbb{E} \left( \frac{1}{n} \sum_{i:D_i=1} \frac{X_{ji}}{\pi_i} \middle| (X_{i'}, Y_{i'}(0), Y_{i'}(1))_{i'=1, \dots, n} \right) &= \mathbb{E} \left( \frac{1}{n} \sum_{i:D_i=0} \frac{X_{ji}}{1 - \pi_i} \middle| (X_{i'} Y_{i'}(0), Y_{i'}(1))_{i'=1, \dots, n} \right) \\ &= \frac{1}{n} \sum_{i=1}^n X_{ji}. \end{aligned}$$

As explained in Section 1, to go beyond the balancing of potential outcomes on average, empiricists can take advantage of the observation of covariates  $X$  before the experiment. A long and natural idea (Fisher, 1926) is to balance these covariates not only in average but also almost surely. Let us define more precisely a perfectly-balanced design.

**Definition 1 (Perfectly-balancing Design)**

A design  $\Pi$  is perfectly-balancing over  $X = (X_1, \dots, X_p)'$  if for  $(D_i)_{i=1, \dots, n}$  sampled in  $\Pi$  we always have for any  $j = 1, \dots, p$ :

$$\frac{1}{n} \sum_{i=1}^n \frac{X_{ji} D_i}{\pi_i} = \frac{1}{n} \sum_{i=1}^n X_{ji} \quad (7)$$

and

$$\frac{1}{n} \sum_{i=1}^n \frac{X_{ji}(1 - D_i)}{1 - \pi_i} = \frac{1}{n} \sum_{i=1}^n X_{ji} \quad (8)$$

Equation (7) describes equality between the covariate sample mean and the estimated mean in the treatment group, whereas equation (8) ensures perfect balance for the control group. A perfectly-balanced assignment eliminates any allocation to the treatment that does not balance perfectly the covariates between treatment and control groups. Note that when  $\pi_i$  is constant, conditions (7) and (8) are equivalent. But this is not the case if the  $\pi_i$  are heterogeneous. A common practice in experiments is to form treatment and control groups of fixed sizes,  $n_T$ , and  $n_C = n - n_T$ , respectively. This is equivalent to setting the constraint in (7) with  $X_{ji} = \pi_i$ . Indeed, for any possible allocation  $(d_1, \dots, d_n)$ :

$$n_T = \sum_{i=1}^n d_i = \sum_{i=1}^n \pi_i = \mathbb{E}(n_T) \quad (9)$$

In that case we also have:  $n_C = \sum_{i=1}^n (1 - d_i) = \sum_{i=1}^n (1 - \pi_i) = \mathbb{E}(n_C)$ . As recommended by Deville and Tillé (2004), we also balance a constant ( $X_{ji} = 1$ ), for the treatment and control groups:

$$\frac{1}{n} \sum_{i=1}^n \frac{D_i}{\pi_i} = \frac{1}{n} \sum_{i=1}^n 1 = 1 \quad (10)$$

$$\frac{1}{n} \sum_{i=1}^n \frac{(1 - D_i)}{1 - \pi_i} = \frac{1}{n} \sum_{i=1}^n 1 = 1. \quad (11)$$

Under such assignment  $\widehat{\theta}_{HT}$  defined in (5) is equal to  $\widehat{\theta}_H$  defined in (6). Notice that, as in Tillé and Favre (2004), we can rewrite (7), (8), (9), (10), (11) as

$$\frac{1}{n} \sum_{i=1}^n \frac{Z_{1i} D_i}{\pi_i} = \frac{1}{n} \sum_{i=1}^n Z_{1i} \Leftrightarrow \frac{1}{n} \sum_{i=1}^n \frac{Z_{0i} (1 - D_i)}{1 - \pi_i} = \frac{1}{n} \sum_{i=1}^n Z_{0i} \quad (12)$$

with  $Z_{1i} = (1, \frac{\pi_i}{1-\pi_i}, \pi_i, X_i', \frac{\pi_i}{1-\pi_i} X_i')'$  and  $Z_{0i} = (\frac{1-\pi_i}{\pi_i}, 1, 1 - \pi_i, \frac{1-\pi_i}{\pi_i} X_i', X_i')' = \frac{1-\pi_i}{\pi_i} Z_{1i}$ . If assignment probabilities are homogeneous (i.e.,  $\pi_i = \pi$ ), the balancing covariates are reduced to  $Z_{1i} = Z_{0i} = (1, X_i')$  due to perfect multicollinearity, but this is no more the case if the  $\pi_i$  are heterogeneous.

The notion of perfect balance is closely related to the balance tests produced by empiricists after randomization. These tests check ex-post that randomization has balanced or not treatment and control on the covariates not only in average but for a particular sampling according to the design  $\Pi$ . A perfectly-balancing design integrates ex-ante the information contained in the observation of  $X$  to ensure balance ex-post and improve the balance in the potential outcomes. For these reasons, debates on the conciliation of randomization and balancing have a long history in statistical sciences (see for instance Fisher, 1926). It is worth noticing that perfect balance is not always attainable: for instance, if  $n = 101$  and  $\pi_i = 1/2$ . Imposing (9) implies  $n_T = 50.5$ , which is simply impossible. But statistical analysis ensures that balancing up to a  $o_p\left(\frac{1}{\sqrt{n}}\right)$  is sufficient to take full advantage of the auxiliary information  $X_i, \pi_i$ . Empiricists are then reduced to find a design  $\Pi$  such that for  $(D_i)_{i=1, \dots, n} \sim \Pi$ ,

$$\frac{1}{n} \sum_{i=1}^n \frac{Z_{1i} D_i}{\pi_i} = \frac{1}{n} \sum_{i=1}^n Z_{1i} + o_p\left(\frac{1}{\sqrt{n}}\right) \text{ for } Z_{1i} = \left(1, \frac{\pi_i}{1 - \pi_i}, \pi_i, X_i', \frac{\pi_i}{1 - \pi_i} X_i'\right)' \quad (13)$$

Various randomization strategies have been proposed and used in the literature to achieve Equation (13). Some of the oldest and widest-used strategies to do so are stratified designs or blocking (Fisher, 1926), re-randomization (Student, 1938; Morgan and Rubin, 2012), matched-pairs design (Imai et al., 2009; Greevy et al., 2004; Ball et al., 1973; Bai et al., 2022). We briefly discussed their limitations in Section 1. In this paper, we advocate the cube method that achieves simultaneously many desirable properties of the various usual balancing designs with a large number of covariates  $Z_1$  that could be as large as a  $O(n^{1/2-1/r})$  if covariates admit moments of order  $r$  and as large as  $o(n^{1/2})$  if covariates have bounded supports. Let us briefly present the cube method before detailing its theoretical properties, the consequence of its use on the estimators  $\widehat{\theta}^H$  and  $\widehat{\theta}^{HT}$  and on the inference about PATE and SATE.

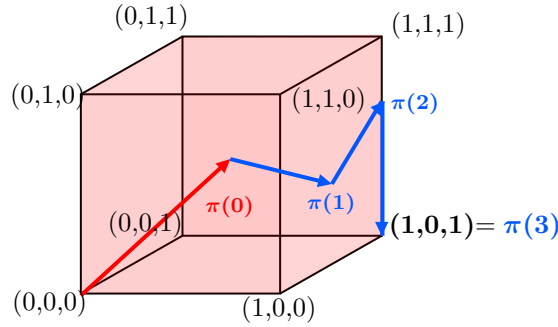
### 3 Balancing Covariates with the Cube Method

Deville and Tillé (2004) first introduced the cube method to produce samples balanced to the population. The cube method consists of an algorithm in two steps: the *flight* and *landing* phases. The technique gets its name from the graphical representation of a sampling problem. Equation (12) ensures that balancing treatment and control groups in an experimental setting for some covariates is equivalent to balancing the treatment group to the entire sample. Let us consider the  $n$ -cube  $C = [0, 1]^n$ . Each vertex of  $C$  (from  $2^n$  possibilities) represents a possible allocation: for instance,  $(1, 1, \dots, 1)$  corresponds to the situation where all units are allocated to treatment,  $(1, 0, 1, 0, \dots, 1, 0)$  corresponds to the case where the treatment group is  $\{i : i \text{ odd}\}$ . A sampling design  $\Pi$  corresponds to how a vertex is selected. Recall that we consider a framework where empiricists impose that Equation (2) holds for  $\Pi$  and a vector  $(\pi_i)_{i=1, \dots, n}$ .

We will first describe the cube algorithm without balancing constraints before moving to the more interesting case where the balancing constraints in Equation (12) are considered. Whatever the set of balancing constraints, the cube method is a discrete martingale that moves in (at most)  $n + 1$  steps from the interior point  $\boldsymbol{\pi}(\mathbf{0}) = (\pi_i)_{i=0}^n$  to  $\boldsymbol{\pi}(\mathbf{n} + \mathbf{1}) = (D_i)_{i=0}^n$  a vertex of  $C$ . Let us consider the case without constraints. At the first step, one chooses a random direction for  $\boldsymbol{\pi}(\mathbf{1}) - \boldsymbol{\pi}(\mathbf{0})$  and a step size such that  $\boldsymbol{\pi}(\mathbf{1})$  belongs to a facet of  $C$  and that  $\mathbb{E}[\boldsymbol{\pi}(\mathbf{1}) | \boldsymbol{\pi}(\mathbf{0})] = \boldsymbol{\pi}(\mathbf{0})$ . After this step, because  $\boldsymbol{\pi}(\mathbf{1})$  belongs to a facet of  $C$ , one component  $i_0$  of  $\boldsymbol{\pi}(\mathbf{1})$  is equal to 0 or 1, selecting  $D_{i_0} = \pi_{i_0}(1)$  one has thus assigned a first unit to either treatment or control. Because a facet of a  $n$ -cube is a  $(n - 1)$ -cube, one can then repeat the process in a  $(n - 1)$ -cube, and so on, until landing in a vertex of  $C$ . At the final step  $n$ , one will have  $(D_i)_{i=1, \dots, n} = \boldsymbol{\pi}(\mathbf{n}) \in \{0, 1\}^n$  and  $\mathbb{E}[D_i] = \pi_i$  (i.e., every unit is allocated to the treatment group with the probability specified by the empiricist). These successive steps are the *flight phase* and for the cube method without balancing constraint, allocation  $(D_i)_{i=1, \dots, n}$  is always determined at the end of this phase. Figure 2 illustrates graphically the method.

In Figure 2, all vertices of the  $n$ -cube can be selected, meaning that all individuals could be allocated to the control group. We now consider that the empiricist wants to allocate a fixed number  $n_T$  of units to the treatment and  $n_c$  units to the control. This can be achieved with the cube method as soon as  $\sum_{i=1}^n \pi_i = n_T$ . The condition that exactly  $n_T$  units are assigned to the treatment can be expressed as a balancing constraint. Indeed, because  $n_T = \sum_i D_i$  and  $\sum_i \pi_i = n_T$ , the fixed size condition is equivalent to  $\sum_i \frac{Z_i D_i}{\pi_i} = \sum_i Z_i$

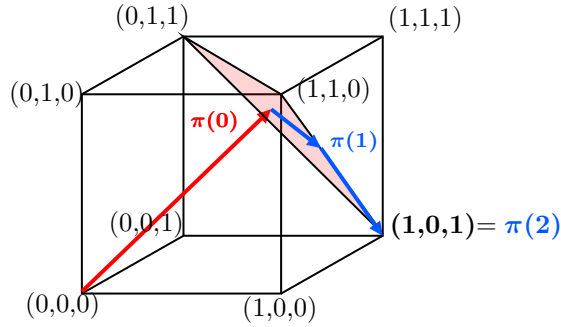
Figure 2: Cube Method without Balancing Constraints



This figure depicts an example of the cube method with  $n = 3$  when no balancing constraint is imposed. The red arrow represents the initial treatment probabilities  $(\pi_i)_{i=1,\dots,n}$ . Then, every blue arrow is a step of the *flight phase*. In this example, the first unit is initially assigned to the treatment group. Then, the third unit is assigned to the treatment group. Last, the second unit is assigned the control group. Therefore, the final allocation – in bold – is  $(1, 0, 1)$ .

for  $Z_i = \pi_i$ . Let  $K$  the set of vectors  $s$  in the  $n$ -cube  $C$  such that  $\sum_i s_i = n_T$ .  $K$  is a closed convex set and its extreme points are vertices of  $C$ , that is the set of allocations respecting the fixed-size constraints.  $K$  is contained in an affine subspace of dimension  $n - 1$  of direction  $V := \{v : \sum_{i=1}^n v_i = 0\}$ , we have  $K = C \cap \{\pi(\mathbf{0}) + v : \sum_i v_i = 0\}$ . The cube method selects randomly an element of  $V$  for the direction of  $\pi(\mathbf{1}) - \pi(\mathbf{0})$  and fixes the step size such that  $\pi(\mathbf{1})$  is a border point of  $K$  and that  $E(\pi(\mathbf{1})|\pi(\mathbf{0})) = \pi(\mathbf{0})$ . After this first step,  $\pi(\mathbf{1})$  belongs to a facet of  $C$  and a unit  $i_1$  is assigned either to the treatment either to the control group. Units  $i \neq i_0$  remain unassigned and we have  $\sum_{i:i \neq i_0} \pi_i(1) = n_T - D_{i_0}$ . We can then replicate the first step after replacing  $n_T$  by  $n_T - D_{i_1}$  the sample  $\{1, \dots, n\}$  by  $\{1, \dots, n\} \setminus \{i_0\}$  and to allocate a second unit and to update assignment probability as  $\pi(\mathbf{2})$ . At step  $n - 1$ ,  $\pi(\mathbf{n} - \mathbf{1})$  belongs to the extreme points of  $K$ , this ends the flight phase. If, for instance,  $\pi_i = 1/2$  and  $n$  is even, the extreme points of  $K$  are some vertices of  $C$ , so the assignment is achieved. Now imagine that one has 101 units to assign with equal probability to the treatment and control groups. Perfect balancing on the two group sizes is not possible: 101 is an odd integer and it is not feasible to assign 50.5 units to the treatment. A popular solution is to consider  $\pi_i = 50/101$  or  $\pi_i = 51/101$  and to sample randomly 50 (or 51) elements among the 101 units. However, this strategy does not accommodate easily with heterogeneous probabilities of assignment and does not generalize to take into account many balancing constraints. With the cube method described above, for each step  $t$  of the flight phase we have  $\sum_i \pi_i(t) = 50.5$  and the extreme points of  $K$  are not anymore vertices

Figure 3: Cube Method with Fixed Sample Size



This figure depicts an example of the cube method with  $n = 3$  when imposing the constraint  $n_T = 2$  and  $\sum_{i=1}^3 \pi_i = 2$ . The red area depicts the points  $(s_1, s_2, s_3)$  in the cube satisfying the equation  $\sum_{i=1}^3 s_i = 2$ . This condition is equivalent to imposing the balancing constraint  $\sum_i \frac{Z_i s_i}{\pi_i} = \sum_i Z_i$  with  $Z_i = \pi_i$ . The red arrow represents the initial treatment probabilities. Then, every blue arrow is a step of the *flight phase*. In this example, the first unit is initially assigned to the treatment group. Then, since  $n_T = 2$ , only one unit among the second and third units can be assigned to the treatment group. In this case, the second blue arrow shows that, in the same step, the second unit is assigned to the control group and the third one to the treatment group. Therefore, the last allocation – in bold – is  $(1, 0, 1)$ .

of  $C$ . In that case, at the end of the flight phase,  $n - 1 = 100$  units are assigned at the end of the flight phase with  $(n - 1)/2 = 50$  units to the treatment and  $(n - 1)/2 = 50$  units to the control. The cube method can be completed with a last phase that randomly assigns to the treatment of the control the remaining unit ensuring that  $n_T = 50$  or  $51$  and  $E(n_T) = 50.5$ . In that case, the sizes of treatment and control groups are not exactly fixed but almost fixed (in fact as fixed as possible as soon as we respect the initial assignment probabilities  $\pi_i = 1/2$ ). This second phase is called *landing phase*. These two phases, the flight phase and the landing phase, can be generalized to the case where the empiricist wants to impose several balancing constraints and heterogeneous probabilities of treatment.

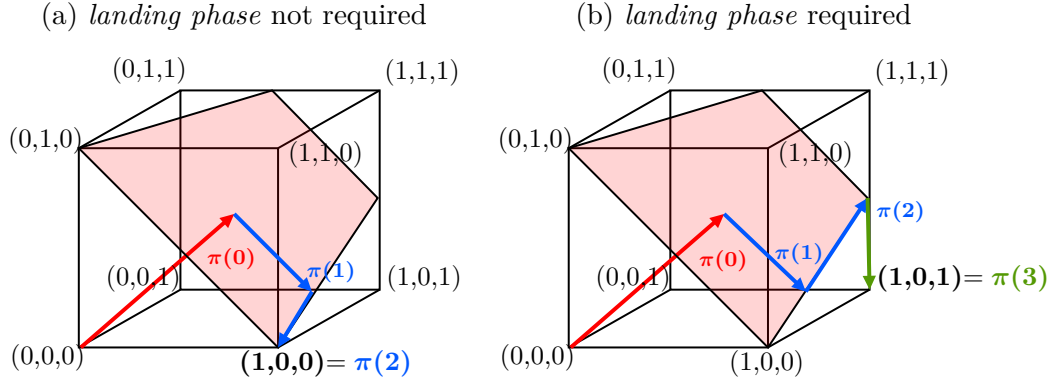
Let us describe the cube method with an arbitrary number of balancing constraints defined by  $q$  variables  $Z_i$ . A point  $\mathbf{s} \in C$  will satisfy an equation analog to (12) if

$$\sum_{i=1}^n \frac{Z_i s_i}{\pi_i} = \sum_{i=1}^n Z_i. \quad (14)$$

Let  $A_i = \frac{Z_i}{\pi_i}$  and  $A = (A_1, \dots, A_n)$  the matrix of size  $q \times n$ . Then (14) is equivalent to

$$\sum_{i=1}^n A_i s_i = \sum_{i=1}^n A_i \pi_i$$

Figure 4: Cube Method with One Balancing Constraint



This figure depicts an example of the cube method with  $n = 3$  where we do not always get perfectly-balanced allocations. We consider the initial treatment probabilities in (2) to be  $\pi_1 = \pi_2 = \pi_3 = \frac{2}{3}$ . The red area depicts the points  $(s_1, s_2, s_3)$  in the cube satisfying the equation  $\sum_{i=1}^3 s_i + s_2 - \frac{1}{2}s_3 = 1$ . This is equivalent to imposing the constraint in (14) with  $Z_1 = Z_2 = \frac{2}{3}$  and  $Z_3 = -\frac{1}{3}$ . The red arrow represents the initial treatment probabilities. Since not every vertex of the plane is a cube vertex, we cannot always satisfy the constraint. In both panels, the algorithm assigns the first unit to the treatment group (first blue arrow). The second blue arrow corresponds to the assignment of the third unit. If the algorithm assigns the third unit to the control group (panel a), it automatically assigns the second one to treatment. However, if the algorithm assigns the third unit to the treatment group (panel b), the second unit is in neither group, even if we attain a plane vertex. In the *landing phase*, the cube method will proceed by randomly allocating the second unit. In this example, the green arrow shows that the *landing phase* allocates the second unit to the control group.

$$\Leftrightarrow A\mathbf{s} = A\boldsymbol{\pi}(\mathbf{0})$$

$$\Leftrightarrow \mathbf{s} \in Q := \boldsymbol{\pi}(\mathbf{0}) + \ker(A).$$

$K = C \cap Q$  is, therefore, the  $(n - q)$ -polytope that contains all the points in  $C$  such that (14) holds. At the first step, one chooses a random direction in  $v \in \ker(A)$  and we select the unique  $\lambda > 0$  such that  $\boldsymbol{\pi}(\mathbf{1}) := \boldsymbol{\pi}(\mathbf{0}) + \lambda v$  is on a facet of  $K$  and that  $\mathbb{E}[\boldsymbol{\pi}(\mathbf{1})|\boldsymbol{\pi}(\mathbf{0})] = \boldsymbol{\pi}(\mathbf{0})$ . Because any facet of  $K$  is the intersection of a facet of  $C$  with  $Q$ , a component  $i_0$  of  $\boldsymbol{\pi}(\mathbf{1})$  is 0 or 1 and defining  $D_{i_0} = \pi_{i_0}(1)$  one has assigned a first unit. Next, one applies a similar step for the facet of  $K$  instead of  $K$  and  $\boldsymbol{\pi}(\mathbf{1})$  as a starting point instead of  $\boldsymbol{\pi}(\mathbf{0})$ . After  $n - q$  steps, one has reached a vertex of  $K$ . This process corresponds to the flight phase in Deville and Tillé (2004). If this vertex of  $K$  is also a vertex of  $C$ , the flight phase allocates every unit, and the two groups are perfectly balanced (see Figures 3 and 4a). But in many

cases, the vertex of  $K$  is not a vertex of  $C$ , and there remain at most  $q$  units to assign during the landing phase (according to the wording of Deville and Tillé, 2004) (see Figure 4b).

Say that at the end of the flight phase, one has not assigned  $r \leq q$  units and let  $\boldsymbol{\pi}^* = \boldsymbol{\pi}(\mathbf{n} - \mathbf{q})$  be the updated treatment probabilities at this stage. The landing phase of the cube method assigns the  $r$  missing units such that  $\mathbb{E}[D_i | \boldsymbol{\pi}^*] = \boldsymbol{\pi}^*$ . Grafström and Tillé (2013) describe two methods for the landing phase (these are also the options used in sampling packages): (i) Linear programming: one considers all the  $2^r$  allocations for these units and assigns probabilities to each allocation to minimize a cost function and satisfy  $\mathbb{E}[D_i | \boldsymbol{\pi}^*] = \boldsymbol{\pi}^*$ . Sampling probabilities are chosen to minimize

$$\mathbb{E} \left( \sum_{i \notin S} Z_i' (D_i - \pi_i^*) M \sum_{i \notin S} Z_i (D_i - \pi_i^*) | W \right),$$

where  $S$  is the set of units allocated at the flight phase,  $W = (S, (D_i)_{i \in S}, (\pi_i^*)_{i \notin S}, (Z_i)_{i=1, \dots, n})$ , and  $M$  is a symmetric positive-definite matrix  $q \times q$ . Common choices for  $M$  are the identity matrix or the inverse-covariance matrix. After solving this minimization problem, the empiricist randomly draws an allocation using these probabilities. (ii) Suppression of variables: if  $r > 20$ , solving a linear problem becomes computationally difficult. In that case, at the end of the flight phase, one can drop a covariate and continue with the flight phase. One can thus successively drop variables until attaining a vertex of  $C$ . This method, however, implies that the empiricist has to define an order to drop the covariates, ideally from the least to the most important.

## 4 Statistical Properties of the Cube Method

This Section shows how the cube method allows obtaining an “almost-perfect balance” between the treatment and control groups and relates this balance to gains in precision for treatment effect estimators.

### 4.1 Balancing Approximations

#### 4.1.1 Balance Checks

As explained above, designing an allocation mechanism that always produces perfectly-balanced groups is impossible. However, we here prove that the cube method is successful,

under certain conditions, in creating almost-perfectly-balanced samples in the sense of Equation (13).

To check balance properties after allocating individuals according to the design  $\Pi$ , empiricists are interested in computing the difference

$$\Delta_{j,n}^{\Pi} = \frac{1}{n} \sum_{i=1}^n \frac{X_{ji} D_i}{\pi_i} - \frac{X_{ji}(1 - D_i)}{1 - \pi_i}.$$

Because  $\mathbb{P}_{\Pi}(D_i = 1 | (X_{i'})_{i'=1, \dots, n}) = \pi_i$ , we have  $\mathbb{E}(\Delta_{j,n}^{\Pi}) = 0$  and under weak conditions on  $\Pi$ , we have

$$\sqrt{n} \Delta_{j,n}^{\Pi} \xrightarrow{d} \mathcal{N}(0, \mathbb{V}(\Delta_{j,n}^{\Pi})), \quad (15)$$

where  $\mathbb{V}(\Delta_{j,n}^{\Pi})$  is an asymptotic variance depending on  $\Pi$  and the distribution of  $X$ .

For the so-called baseline balance tests, empiricists often consider the  $t$ -statistic

$$t_{j,n}^{\Pi} = \sqrt{n} \frac{\Delta_{j,n}^{\Pi}}{\sqrt{\widehat{\mathbb{V}}(\Delta_{j,n}^{\Pi})}}$$

where  $\widehat{\mathbb{V}}(\Delta_{j,n}^{\Pi})$  is a consistent estimator of the asymptotic variance of  $\Delta_{j,n}^{\Pi}$  to test the null hypothesis of perfect balance.  $t_{j,n}^{\Pi}$  is then associated to a  $p$ -value  $p_{j,n}^{\Pi}$  which take values between 0 and 1. As explained by Snyder and Zhuo (2024), when creating balance tests for RCTs, small  $p$ -values (below 0.15) are considered problematic.

Let us first consider a naive mechanism that does not use baseline information to assign units. Such situations correspond to the case where the design  $\Pi$  is a Poisson design, ie a design where each unit  $i$  is allocated to the treatment independently of the allocation of other units:

$$\mathbb{P}_{\Pi} \left( \bigcap_{i=1}^n \{D_i = d_i\} | (X_i)_{i=1}^n \right) = \prod_{i=1}^n \pi_i^{d_i} (1 - \pi_i)^{1-d_i}.$$

A Poisson design does not balance any variable nor group sizes.

When  $\pi_i = \frac{n_T}{n}$  for any  $i$ , another popular design is sampling without replacement of  $n_T$  treated units, also known as complete randomization:

$$\mathbb{P}_{\Pi} \left( \bigcap_{i=1}^n \{D_i = d_i\} | (X_i)_{i=1}^n \right) = \binom{n}{n_T}^{-1} \mathbb{1} \left\{ \sum_{i=1}^n d_i = n_T \right\}.$$

Complete randomization only balances constant variables. In that case, the sample of treated and control groups are fixed, and the design is also balanced on the constant  $\sum_{i=1}^n D_i = n_T$ ,  $\sum_{i=1}^n (1 - D_i) = n - n_T$  and  $\sum_{i=1}^n \frac{D_i}{\pi_i} = n$ .



Under such assignments and Assumptions 1 and 2, and more generally for any design  $\Pi$  such that (15) holds with  $\mathbb{V}(\Delta_{j,n}^\Pi) > 0$ , we have  $\Delta_{j,n}^\Pi = O_p\left(\frac{1}{\sqrt{n}}\right)$ ,  $t_{j,n}^\Pi \xrightarrow{d} \mathcal{N}(0, 1)$  and  $p_{j,n}^\Pi \xrightarrow{d} \mathcal{U}(0, 1)$ . This result means that if one randomizes naively, control and treatment groups will present imbalances with a strictly-positive probability. Moreover, for a confidence level of  $100(1 - \alpha)\%$ , there exists always  $100\alpha\%$  chance of obtaining significant differences. If an empiricist evaluates the balance of 10 independent covariates at the 85% confidence level (a level over which rejection is considered problematic for publication as shown by Snyder and Zhuo, 2024), there is more than 80% chance of having at least one significant difference. This magnitude questions the mere implementation of such widely used tests. Even if a multiple F-test with a confidence level of 85% mitigates this rejection rate, the null hypothesis of simultaneously balanced covariates is rejected by construction with a 15% chance.

The cube method ensures that these tests are unnecessary since we can balance control and treatment groups in any covariate  $(X_j)_{j=1,\dots,p}$ . This is achieved because  $\mathbb{V}(\Delta_{j,n}^\Pi) = 0$  for any  $j = 1, \dots, p$  in (15). Performing these tests would not make sense since we never reject the null hypothesis by construction. However, one might report them if the editor worries about empiricists randomizing badly. Usual balancing strategies are stratified or matched-pairs designs. These methods ensure  $\mathbb{V}(\Delta_{j,n}^\Pi) = 0$  if the covariates  $(X_j)_{j=1,\dots,p}$  are all discrete but will always generate imbalances for continuous ones since the empiricist needs to discretize or aggregate them before randomizing.

The following proposition explains how the balancing approximations are satisfied with the cube method. Because the number  $q$  of balancing constraints in equation (13) could be large with the cube method, we are also explicit on how  $q$  affects balancing approximations to allow us to consider a framework where  $q$  tends to  $\infty$ .

**Proposition 1 (Balancing Approximations with the Cube Method)**

*If Assumptions 1 and 2 hold, then*

$$\Delta_{j,n}^{cube} = o_p\left(\frac{q}{\sqrt{n}}\right).$$

*Moreover if  $\mathbb{E}[|X_{j1}|^r] < \infty$  for  $r \geq 2$ , then  $\Delta_{j,n}^{cube} = o_p\left(\frac{q}{n^{1-1/r}}\right)$ , if  $X_{j1}$  is sub-Gaussian, then  $\Delta_{j,n}^{cube} = O_p\left(\frac{q\sqrt{\ln(n)}}{n}\right)$ , and if  $X_{j1}$  has a bounded support, then  $|\Delta_{j,n}^{cube}| < \frac{Kq}{cn}$  for  $K$  such that  $|X_{j1}| < K$ . As soon as  $\sqrt{n}\Delta_{j,n}^{cube} = o_p(1)$ , we have  $t_{j,n}^{cube} \xrightarrow{P} 0$ , and  $p_{j,n}^{cube} \xrightarrow{P} 1$ .*

Proposition 1 shows that, as  $n$  grows, the cube method ensures the balancing equation (13) as soon as the second-order moments of  $X$  exist. Furthermore, if moments of order  $r > 2$

exist for  $X$ , (13) holds as soon as  $q = O\left(n^{\frac{1}{2}-\frac{1}{r}}\right)$ .  $q$  can even be  $o\left(\sqrt{\frac{n}{\ln(n)}}\right)$  if the covariates  $X$  are all sub-Gaussian or  $o(\sqrt{n})$  if they are bounded. This means that with probability tending to one, the  $p$ -values of balance tests tend to 1. Balance is thus never rejected for large  $n$  contrary to randomization under a design  $\Pi$  such that (13) does not hold.

#### 4.1.2 Comparison with Other Methods

We here compare the balancing properties of the cube method with other randomization methods. For the sake of simplicity, we fix  $\pi_i = 1/2$ . We check imbalances by using the Horvitz-Thompson estimators for the average difference between the control and treatment groups  $B_{n,p}(X) = \frac{2}{n} \sum_{i=1}^n X_i D_i - X_i(1 - D_i)$  and looking at their squared Euclidean norm  $\|B_{n,p}(X)\|^2 = \frac{4}{n^2} \sum_{j=1}^p (\sum_{i=1}^n X_{ji} D_i - X_{ji}(1 - D_i))^2$ .

**Assumption 3**  $n$  is an even positive number and  $X_i = (X_{1i}, \dots, X_{pi})'$  are some independent and identically distributed random vectors of dimension  $p$  for  $i = 1, \dots, n$ .  $X_1$  admits a density  $f_X$  with respect to the Lebesgue measure on  $[0; 1]^p$  such that there exists some positive constants  $\underline{C}$  and  $\overline{C}$  (independent of  $p$ ) such that for any  $x \in [0; 1]^p$ ,  $\underline{C} \leq f_X(x) \leq \overline{C}$ .

Assumption 3 imposes mild conditions over the baseline covariates. In particular, the components of the vector  $X_i$  are not assumed to be independent. Figure 1 illustrates our main results for a simple case where this assumption holds:  $(X_{ji})_{j=1, \dots, p, i=1, \dots, n}$  are independent and follow a uniform distribution on  $[0, 1]$ . In this case, we have  $V(X_{ji}) = 1/12$ ,  $\mathbb{E}(X_{ji}^2) = 1/3$ , and  $\underline{C} = \overline{C} = 1$ .

To create comparable treatment and control groups, empiricists sometimes use “naive randomization” methods, meaning they do not use any information previously available to design the assignment mechanism. The simplest naïve method is Bernoulli randomization, where empiricists allocate units to treatment independently and with equal probabilities. This method corresponds to assigning units to treatment using an independent coin toss for each unit. Bernoulli randomization ensures that treatment and control groups are balanced on average (i.e.,  $\mathbb{E}[B_{n,p}] = 0$ ). However, we can still have imbalances between groups for a given allocation. Additionally, Bernoulli randomization will often generate different sizes between the treatment and control groups, meaning that it fails to balance on a constant. Proposition A.1 in Appendix A shows that for Bernoulli randomization,  $\frac{4\underline{C}}{3} \frac{p}{n} \leq \mathbb{E}[\|B_{n,p}(X)\|^2] \leq \frac{4\overline{C}}{3} \frac{p}{n}$ . For the case illustrated in Figure 1, we thus have  $\mathbb{E}[\|B_{n,p}(X)\|^2] = \frac{4p}{3n}$ .

Another naïve method that improves from Bernoulli sampling is complete randomization. In this method, the researcher fixes the group sizes. Since we here assume  $\pi_i = 1/2$ , the

empiricist randomly chooses an allocation among those having an equal number of treated and untreated units. The empiricist still does not use any information on the baseline covariates to refine the randomization process but manages to reduce the imbalances due to different group sizes. Indeed, under Assumption 3 and complete randomization  $\frac{C}{3} \frac{p}{n} \leq \mathbb{E}[|B_{n,p}(X)|^2] \leq \frac{\bar{C}}{3} \frac{p}{n}$ , so bounds reduce by four, relative to Bernoulli randomization. For the example in Figure 1, we have  $\mathbb{E}[|B_{n,p}(X)|^2] = \frac{p}{3n}$ . This result clearly shows the advantages of using designs with fixed sample sizes such as complete randomization, but also the cube method with the constraints in Section 2 or matched-pairs design.<sup>3</sup>

In sharp contrast with naïve methods, covariate-adaptive randomization uses baseline information to improve balance between treatment and control. Stratified designs are the most used and studied covariate-adaptive method. Stratification has a long tradition in RCTs (Fisher, 1935; Higgins et al., 2016). Stratified designs are the most popular assignment mechanisms used in RCTs as they are simple to grasp and can produce balanced samples. This method consists of using one or several baseline variables to create blocks or strata and then using complete randomization inside each stratum. A common practice in experiments is to block on gender, meaning that randomization is performed independently amongst male and female units, generating the same proportion of men and women in each treatment arm. When using dummy variables to define the strata, stratified or blocked randomization allows almost perfect balancing of the variables used to create them. Athey and Imbens (2017) recommend balancing on small strata since this method generates substantial precision gains. However, stratified designs do not come without any limitations. Notably, the type and number of covariates that one wants to balance can impose some difficulties. Facing continuous covariates, such as income or grades, makes it impossible to stratify without the empiricist deciding how to create the strata. Assumption 3 imposes continuous covariates. We thus will focus on two ways of generating (possibly-)small strata, discretization and matched pairs.

First, the empiricist can discretize continuous variables using  $\ell$ -quantiles for each covariate, generating thus  $\ell^p$  strata. Stratifying will produce balance gains as long as the number of units remains large compared to the number of strata. In particular, we show in Proposition A.1 that whenever  $n\ell^{-p} \rightarrow \infty$ , stratified designs through discretization outperform complete randomization. Discretizing baseline covariates, however, does not ensure fixed sizes for each

---

<sup>3</sup>One can show that the gains from fixed group sizes are not present if one uses a difference-in-means estimator instead. Moreover, in the case of  $n$  even and  $\pi_i = 1/2$ , this estimator is equivalent to the Horvitz-Thompson for complete randomization, matched-pairs design, and the cube method. However, it can lead to more precise estimates for Bernoulli randomization or stratified designs.

stratum. In particular, if the number of strata is big compared to the sample size, there is a big chance of having some strata with only one unit. In the limit case where  $n\ell^{-p} \rightarrow 0$ , every non-empty strata has one unit with probability one, and stratifying through discretization performs strictly worse than complete randomization and approximates Bernoulli randomization. We also show that, in both limit cases, imbalances grow at a rate of  $p/n$ . We observe this behavior in Figure 1 for  $\ell = 2, 4$ . In this example, the stratified designs perform better than complete randomization whenever  $l^p \leq n/2$  and similarly to Bernoulli randomization for  $l^p \geq 32n$ . Balancing deterioration can thus occur quite rapidly when stratifying is done by discretizing many continuous variables. It is worth noticing that this issue is not exclusive to continuous variables, as it arises when stratifying using many categorical variables. In that case, the number of strata equates to the product of covariate support cardinalities.

To eliminate the issue of single-unit strata, empiricists may use a more sophisticated way of creating their strata: matched-pair designs. Following Greevy et al. (2004); Bai et al. (2022); Bai (2022), the empiricist can create  $n/2$  strata of two units to minimize the average intra-strata distance. The empiricist thus creates pairs of two units that resemble each other. After constructing these strata, the empiricist randomly allocates one to treatment. By doing so, she creates control and treatment groups that are very similar. We show in Proposition A.1 that this design always outperforms complete randomization. However, this type of strata construction works by trying to have a similar joint distribution of  $X$  between treatment and control. This approach to balancing implies that for large  $p$ , it becomes more difficult to find pairs of units close to each other. In particular, we show that under this design and Assumption 3,  $\mathbb{E}[|B_{n,p}(X)|^2] \geq \frac{p}{n} \left( \frac{1}{3} - \sqrt{\frac{2\ln(n-1)+4\ln\bar{C}}{p}} \right)$ . This result entails that the number of balancing covariates  $p$  is large relative to  $\ln(n)$ , balance gains shrink, and imbalances increase at the rate of  $p/n$ . Figure 1 illustrates this effect. Indeed, when  $p$  becomes larger than  $\ln(500) \approx 6$ , the matched-pairs design performs better than complete randomization, but its relative gains quickly reduce.

For all these randomization methods, imbalances grow at the rate of  $p/n$ . We now show that the cube method is less concerned by this curse of dimensionality since imbalances grow at a rate  $p^2/n^2$ . If  $p$  remains smaller than  $n$ , then this result implies a much slower balance deterioration than for the methods described above. Proposition 2 gives this upper bound for  $\mathbb{E}[|B_{n,p}(X)|^2]$  when using the cube method. This proposition is also stated and proved in Appendix A.

**Proposition 2 (Imbalance under the Cube Method)**

*Suppose Assumption 3 holds. Under the cube method using linear programming with positive-*

definite matrix  $M$  for the landing phase, we have

$$\mathbb{E} [\|B_{n,p}(X)\|^2] \leq 4 \frac{(p+1)^2}{n^2} \frac{\lambda_{max}(M)}{\lambda_{min}(M)}$$

for  $\lambda_{max}(M)$  and  $\lambda_{min}(M)$  the largest and the smallest eigenvalues of  $M$ .

The upper bound depends on the matrix  $M$ , described in equation Equation (7) in Deville and Tillé (2004), used during the landing phase. Notably, one can take  $M$  the identity matrix, and we have  $\frac{\lambda_{max}(M)}{\lambda_{min}(M)} = 1$ . Then, we see that the cube method outperforms other methods that grow at a rate of  $p/n$ . This is clearly illustrated in Figure 1, where we see that imbalances increase only very lightly on the number of covariates when randomizing with the cube method. The main difference between the cube and other designs is that it balances selected moments of the covariates instead of balancing the whole joint distribution of  $X$ . It thus reduces the burden of balancing a higher number of covariates. Balancing moments can also be achieved through other methods. In particular, we can perform re-randomization such that the stopping criterion requires balancing moments of  $X$  or perform a Gram-Schmidt walk design (Harshaw et al., 2024).

Re-randomization is another method that allows obtaining balance between covariates that has gained focus in the last decades (Morgan and Rubin, 2012; Li et al., 2018; Imbens, 2011). The main idea of re-randomization is to completely randomize repeatedly until the obtention of balanced groups. Some empiricists perform re-randomization without pre-specifying it. This repetition affects treatment probabilities in an unknown manner, which induces invalid inference (Bruhn and McKenzie, 2009; Athey and Imbens, 2017). There are, however, several ways of performing re-randomization that allow valid inferences to some extent. Most of them rely on the simulation of the distribution under the re-randomization procedure used to assign units. This implies that empiricist should draw a large number  $N$  of balanced samples. One can keep randomizing until  $\|B_{n,p}(X)\|^2 \leq 4 \frac{(p+1)^2}{n^2}$  to achieve the upper bound of Proposition 2 (with  $M = Id$ ). However, this upper bound is not sharp and to compare re-randomization with the cube method, we counted how many times an empiricist should sample with naive randomization to get  $N = 1000$  samples that are balanced as well as  $B^{*2} = \mathbb{E} (\|B_{n,p}(X)\|^2)$ , where the previous expectation is computed for the cube method through simulations. Under the design used in Figure 1, and for complete randomization,  $12/4 \times n \times \|B_{n,p}(X)\|^2$  converges in distribution to  $\chi^2(p)$ . Next, the probability to achieve balancing as good as the cube is  $F_{\chi^2(p)}(12/4 * n * B^{*2})$ . To have  $N = 1000$  samples balanced as well as the cube, empiricists thus have to sample approximately  $1000/F_{\chi^2(p)}(3 * n * B^{*2})$ . For  $p = 3$ , the empiricist have to sample more than  $10^6$  samples and for  $p = 10$  this is

more than  $9,98 \times 10^{11}$ . The probability of getting a sample that has the same properties as the cube method becomes small very quickly, so it becomes demanding computationally, in particular, if one wants several allocations to perform randomization-based inference.

Harshaw et al. (2024) recently developed the Gram-Schmidt walk design to obtain balanced groups in RCTs. As the authors consider a tradeoff between balance and robustness, they impose a choice of parameter  $\phi \in [0, 1]$ . For  $\phi = 1$ , the algorithm from Harshaw et al. (2024) reduces to Bernoulli randomization and next  $\frac{4Cp}{3n} \leq \mathbb{E}(\|B_{n,p}(X)\|^2) \leq \frac{4\bar{C}p}{3n}$  under Assumption 3. For  $\phi \in (0; 1)$ , we conjecture that  $\underline{K}_1\phi\frac{p}{n} + \underline{K}_0(1 - \phi)\frac{p^2}{n^2} \leq \mathbb{E}(\|B_{n,p}(X)\|^2) \leq \bar{K}_1\phi\frac{p}{n} + \bar{K}_0(1 - \phi)\frac{p^2}{n^2}$  for some constant  $\underline{K}_0, \underline{K}_1, \bar{K}_0, \bar{K}_1$ . The balance of the Gram-Schmidt method would be as good as the cube method if  $\phi = 0$ . However, theoretical results in Harshaw et al. (2024) and implementation of the Julia package depend on bounding  $\phi$  above 0.

## 4.2 Variance Reduction

The balance between covariates in the control and treatment groups is also beneficial if these variables are related to the potential outcomes. In this case, using the cube method will also reduce the variance of the Horvitz-Thompson and Hájek estimators.

### Assumption 4

For  $d \in \{0, 1\}$ ,

$$Y_i(d) = \beta_d Z_{di} + \varepsilon_i(d), \text{ with } \mathbb{E}[\varepsilon_i(d)|Z_{di}] = 0$$

.

Assumption 4 states that potential outcomes are linearly related to observable covariates. However, we allow heterogeneity in treatment effects by specifying different equations for control and treatment groups.

### Conjecture 1 (Poisson Approximation)

For any  $k \in \mathbb{N}^*$  we have with probability one:

$$\lim_{n \rightarrow \infty} \sup_{i_1, \dots, i_k} \left| \mathbb{E} \left( \prod_{j=1}^k (D_{i_j} - \pi_{i_j}) \mid X_1, \dots, X_n \right) \right| = 0$$

This conjecture establishes that as  $n$  increases, the cube method tends to Poisson sampling. As  $n$  goes to infinity, the dependence between the assignment of a finite number of

individuals disappears. We draw this conjecture from results in Deville and Tillé (2005) and simulations that confirm it.

To have a benchmark for the gains in variance decline, we compare the cube method with Poisson randomization, i.e., an unconstrained sampling with heterogeneous treatment probabilities. The results also hold for Bernoulli randomization, that is, with homogeneous treatment probabilities.

**Proposition 3 (Asymptotic Normality)**

Let  $\theta_0$  be the SATE defined in (3),  $\theta_0^*$  the PATE in (4) and  $\hat{\theta}$  the HT or H estimator in (5) and (6). If Assumptions 1, 2 and 4, and Conjecture 1 hold, and if  $\Pi$  is a balancing sampling using the cube method we have:

$$\sqrt{n} (\hat{\theta} - \theta_0) \xrightarrow{d} \mathcal{N}(0, V_0)$$

and

$$\sqrt{n} (\hat{\theta} - \theta_0^*) \xrightarrow{d} \mathcal{N}(0, V_0^*).$$

for  $V_0 = \mathbb{E} \left( \pi_i(1 - \pi_i) \left( \frac{\varepsilon_i(1)}{\pi_i} + \frac{\varepsilon_i(0)}{1 - \pi_i} \right)^2 \right)$  and  $V_0^* = \mathbb{V}(Z'_{1i}\beta_1 - Z'_{0i}\beta_0) + \mathbb{E} \left[ \frac{\varepsilon_i(1)^2}{\pi_i} \right] + \mathbb{E} \left[ \frac{\varepsilon_i(0)^2}{1 - \pi_i} \right]$ . If Poisson randomization is used instead, we have:

$$\sqrt{n} (\hat{\theta} - \theta_0) \xrightarrow{d} \mathcal{N}(0, V_0 + \Sigma_0)$$

and

$$\sqrt{n} (\hat{\theta} - \theta_0^*) \xrightarrow{d} \mathcal{N}(0, V_0^* + \Sigma_0).$$

with  $\Sigma_0 = \mathbb{E} \left[ \frac{\pi_i}{1 - \pi_i} (Z'_{i0}(\beta_1 + \beta_0))^2 \right] \geq 0$ .

Proposition 3 shows the gain in asymptotic variance from balancing covariates using the cube method. The reduction is more substantial when  $X$  explains more of the potential outcomes. Estimates of the ATE are thus more precise when using the cube method. This reduction can represent significantly lower costs when conducting an RCT. Notice that under the same set of assumptions,  $V_0^*$  corresponds to the semiparametric efficiency bound in Hahn (1998). Simulations in Sections 5 illustrate these gains.

### 4.3 Inference

This section provides properties of the cube algorithm and methods to perform inference. We elicit two main techniques of conducting inference, one based on the asymptotic properties of the HT estimator and the other based on the randomization mechanism.

### 4.3.1 Asymptotics-based Inference

Some methods, such as re-randomization, alter the inclusion probabilities in a manner that is unclear to the empiricist (Imbens, 2011). When the criterion for selection is known and behaves in a known way, such as the Mahalanobis distance, one can perform conservative inference. However, balance is imperfect for numerous covariates. Since the cube method assigns treatment only once, we can perform asymptotic-based inference. We here give the asymptotic properties and propose an easy way to construct exact confidence intervals.

To construct a confidence interval, one would like to estimate either  $V_0$  or  $V_0^*$ . Estimating  $V_0/n$  is impossible without making assumptions on the relation between  $\varepsilon_i(1)$  and  $\varepsilon_i(0)$ . This issue is common in RCTs. We can, nonetheless, easily construct an unbiased estimator  $\widehat{V}$  for  $V_0^*/n$ . Let  $\widehat{\beta}_d$  and  $\widehat{\varepsilon}_i(d)$  be the estimated coefficients and residuals, respectively, of a regression of  $Y_i(d)$  on  $Z_{di}$ , for  $d \in \{0, 1\}$ . We then have

$$\widehat{V} = \frac{1}{n} \left[ \widehat{V}(Z'_{1i}\widehat{\beta}_1 - Z'_{0i}\widehat{\beta}_0) + \frac{1}{n} \sum_{i=1}^n \frac{\widehat{\varepsilon}_i(1)^2 D_i}{\pi_i^2} + \frac{1}{n} \sum_{i=1}^n \frac{\widehat{\varepsilon}_i(0)^2 (1 - D_i)}{(1 - \pi_i)^2} \right]. \quad (16)$$

Then, we can test the weak hypothesis

$$H_0 : \theta_0^* = 0, \quad (17)$$

and construct the confidence interval based on

$$\widehat{\theta} \pm \Phi^{-1}(1 - \alpha/2) \sqrt{\widehat{V}} \quad (18)$$

In Section 5.2, we perform simulations that confirm the exact coverage rate of this confidence interval when  $n$  is big enough.

### 4.3.2 Randomized-based Inference

We here study the properties of randomization-based inference when permuting treatment status while satisfying balancing constraints. For these tests, we consider the stronger null hypothesis:

$$H_0 : (Y_i(1), X_i) \stackrel{d}{=} (Y_i(0), X_i). \quad (19)$$

Notice that testing this hypothesis, under Assumptions 1 and 2 is equivalent to testing  $(Y_i)_{i=1}^n \perp\!\!\!\perp (D_i)_{i=1}^n | X_1, \dots, X_n$  (Proof in Appendix B.4).

To explain the test, we introduce some new notation. Let  $G_n$  be the set of all possible  $2^n$



assignments. Then, we can define the set of assignments  $G_n^{cube} \subseteq G_n$  satisfying the constraints imposed by the cube method. That is, with Assumptions 1 and 2,

$$G_n^{cube} = \left\{ g \in G_n : \Delta_{j,n} = o_p \left( \frac{q}{\sqrt{n}} \right) \text{ for } 1 \leq j \leq p \right\}.$$

We note  $\mathbf{P}_n = (Y_i, D_i, X_i)_{i=1}^n$  the observed values, and  $\mathbf{P}_n^{(g)} = (Y_i, D_i^{(g)}, X_i)_{i=1}^n$ , the new data where we have reassigned treatment according to  $g \in G_n^{cube}$ . For computational facility, we can replace  $G_n^{cube}$  by  $G_n^B = \{g_1, \dots, g_B\}$ , such that  $g_1$  is the assignment really obtained and  $(g_i)_{i=1}^B$  are drawn independently from a uniform distribution on  $G_n^{cube}$ .

Then, for a given test statistic  $T_n(\mathbf{P}_n)$ , we consider the test

$$\phi_n^{rand}(\mathbf{P}_n) = \mathbb{1} \{T_n(\mathbf{P}_n) > c_n(\mathbf{P}_n, 1 - \alpha)\}$$

with

$$c_n(\mathbf{P}_n, 1 - \alpha) = \inf \left\{ t \in \mathbb{R} : \frac{1}{B} \sum_{g \in G_n^B} \mathbb{1} \{T_n(\mathbf{P}_n^{(g)}) \leq t\} \geq 1 - \alpha \right\}.$$

**Proposition 4 (Randomization-based Inference)**

*Under Assumptions 1 and 2, and the null hypothesis in (19),*

$$\mathbb{E} [\phi_n^{rand}(\mathbf{P}_n)] \leq \alpha.$$

Proposition 4 indicates that if  $T_n(\mathbf{P}_n) > c_n(\mathbf{P}_n, 1 - \alpha)$ , we reject the null hypothesis (19) at the  $\alpha$  level. The proof is similar to previous results on other covariate-adaptive assignment mechanisms (Heckman et al., 2010, 2011; Lee and Shaikh, 2014; Bai et al., 2022), but it is presented for completeness. This proposition ensures that we can compute Fisher’s  $p$ -values by comparing our test statistic with those produced by other assignments made by the Cube method.

## 5 Simulations

This section compares the cube method to other randomization methods by performing Monte Carlo simulations. We are interested in examining the impact of introducing new covariates in the variance of treatment effect estimates. For this purpose, we evaluate different randomization methods using one example from data following a simple DGP in the spirit of Figure 1 and another using data-driven methods from an empirical application.

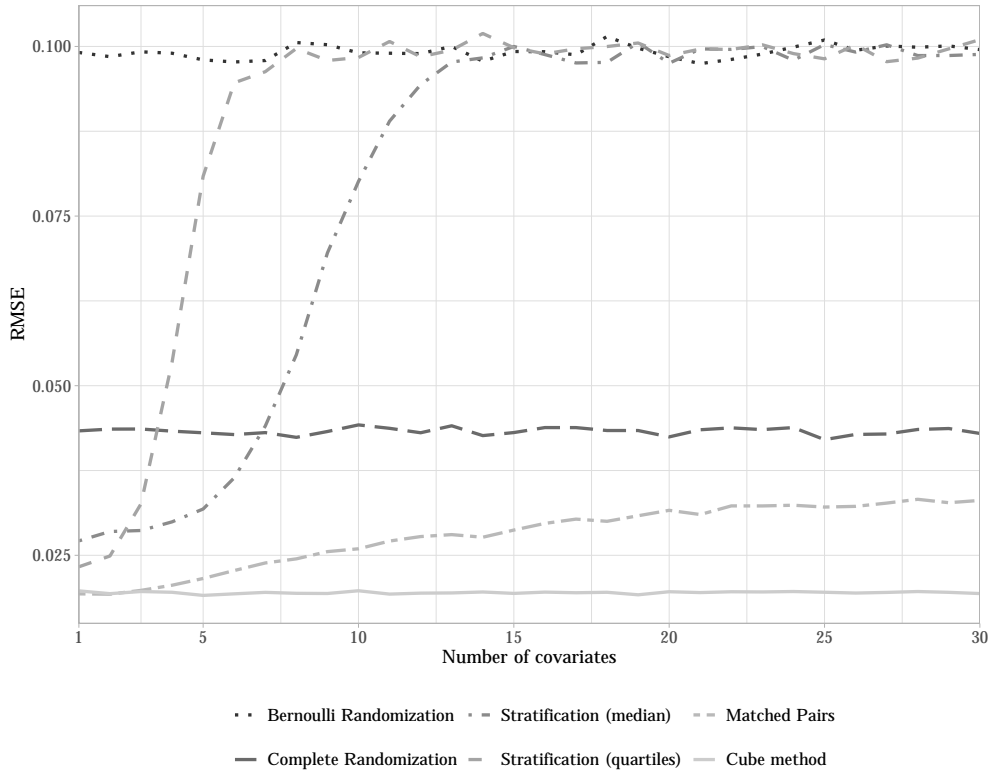
## 5.1 Simple DGP

For  $k = 1, \dots, K$  the number of iterations,  $j = 1, \dots, p$  the number of covariates, and  $i = 1, \dots, n$ , the number of observations, we independently draw  $X_{jik} \sim \mathcal{U}(0, 1)$  and  $\varepsilon_{ik}(d) \sim \mathcal{N}(0, 1)$ , for  $d = 0, 1$ . We then generate the potential outcomes  $Y_{ik}(0) = 1 + (X_{ik} - 1/2)'\beta_0 + \varepsilon_{ik}(0)$  and  $Y_{ik}(1) = 1 + X_{ik}'\beta_1 + (X_{ik} - 1/2)'A(X_{ik} - 1/2) + \varepsilon_{ik}(1)$ , with  $A = (1/20) \times (\mathbb{1}\mathbb{1}' - \text{diag}(1))$ . Notice that, in this example,  $\theta_0^* = 0$ .

We consider  $n = 500$ ,  $p = 30$ ,  $\beta_0 = (1, \mathbf{0})'$ ,  $\beta_1 = 2\beta_0$ , so only one covariate explains variations in potential outcomes. We assume that the empiricist knows she should always balance this covariate. Still, she does not have previous information about the (un)informativeness of the 29 other covariates. In these simulations, the empiricist has to choose which simulation method she uses and how many covariates to include. That choice corresponds to an assignment design  $\Pi$  and generates treatment statuses  $D_{ik}^\Pi$ . We estimate the PATE using the HT estimator  $\hat{\theta}_{HT,k}^\Pi$ . To evaluate the precision entailed by the assignment design, we perform  $K = 5,000$  simulations and compute the standard deviation of the estimator over the simulations. Since the PATE is null, this is equivalent to estimating the root mean square error (RMSE).

Figure 5 shows the RMSE of the HT estimator by number of covariates and randomization method. The simulations show that the cube method is always competitive. Since including more covariates deteriorates balancing only very lightly, precision gains are maintained even when  $p = 30$ . This behavior is not present for other randomization methods. Indeed, stratification using the median (quartile) leads to worse precision than complete randomization as soon as  $p > 6$  ( $p > 3$ ) and converges to Bernoulli randomization for  $p > 12$  ( $p > 6$ ). Moreover, using a matched-pairs design improves from complete randomization but, for  $p > 3$ , underperforms compared to the cube method: when  $p$  increases, precision for matched-pairs design worsens, whereas it remains the same when using the cube. By allowing an abundant set of covariates, the cube method improves the exploitation of balancing gains, even when the empiricist chooses to balance covariates that are not explicative of potential outcomes. This behavior could arise if the empiricist is interested in several treatment outcomes and collects their pre-treatment values. Then, she would ideally want to balance them all, even if only one covariate is explicative of one outcome. As described through these simulations, the cube method ensures precision gains for a particular outcome variable, even when balancing another 29 baseline variables.

Figure 5: Precision of HT Estimator by Randomization Design



This figure depicts the precision of the HT estimator for each design (randomization method + number of covariates). The DGP is as described in Section 5.1 and such that only one covariate is informative of potential outcomes, which is always included by the empiricist when performing covariate-adaptive randomization. We compute the RMSE of the HT estimator by taking its standard deviation of 5,000 Monte Carlo iterations.

## 5.2 Empirical Data

We further illustrate the properties of the cube method by using experimental data from Gerber et al. (2020). The authors investigate how informing potential voters about the closeness of an election affects their beliefs and voting behavior. Since the experimental data only represents one of many possible samplings, we proceed by generating a *superpopulation*. We create a large dictionary with baseline outcomes, covariates, and demographics. We consider all possible interactions and second-order polynomials. We thus generate a dataset of 6,424 observations and 7,381 covariates (hereon denoted by  $X$ ), with 3,193 individuals in the treatment group. We consider beliefs about the closeness of the election as the main outcome  $Y$ . We ran two lasso regressions separately for treated and control units to train two

models,  $f_1$  and  $f_0$ . We then estimate  $s_1^2 = \widehat{\mathbb{V}}(Y - f_1(X)|D = 1)$  and  $s_0^2 = \widehat{\mathbb{V}}(Y - f_0(X)|D = 0)$ . To generate the superpopulation we draw  $N = 50,000$  individuals, with replacement and we generate  $Y_i(1) = f_1(X_i) + \varepsilon_i(1)$  and  $Y_i(0) = f_0(X_i) + \varepsilon_i(0)$  for  $i = 1, \dots, N$  with  $(\varepsilon_i(1); \varepsilon_i(0)) \sim \mathcal{N}((0, 0), (s_1^2 \ 0.5s_1s_0, 0.5s_1s_0 \ s_0^2))$ . We thus obtain a superpopulation  $(X_i, Y_i(1), Y_i(0))_{i=1, \dots, N}$ .

We then run  $K = 10,000$  Monte Carlo simulations, where for every iteration, we draw  $n \in (100, 256, 500, 1000)$ <sup>4</sup> individuals, allocate them according to five treatment allocation methods: complete randomization, stratified randomization using median values for continuous variables, matched-pairs design using the Mahalanobis distance when balancing multiple covariates, the cube method with the two first moments per variable, and complete randomization with ex-post double-lasso selection. For stratified designs, matched-pairs design, and the cube method, we balance between 1 and 12 covariates. When balancing only one, we use the pre-treatment value of  $Y$ . For the 12 covariates, we consider five pre-treatment outcomes and seven baseline covariates. We always prioritize pre-treatment outcomes as they are likely the most explicative variable for their post-treatment counterpart. We set  $\pi_i = \frac{1}{2}$ . For complete randomization, matched pairs, and the cube method, we compute the HT estimator. For these methods, we compute confidence intervals using, respectively, White standard errors, Equation (14) in Bai (2022), and Equation (16) in Section 4.3 above. For stratification, since  $\pi_i = \frac{1}{2}$ , we use an OLS regression with strata fixed-effects, which gives consistent estimators and exact inference as shown by Bugni et al. (2018).) For the double lasso, we use the same augmented dictionary as with imputation and proceed with the double selection method described in Belloni et al. (2014).

Table 1 reports estimators of the effective sample size  $ESS = \frac{\mathbb{V}(\hat{\theta}_{HT}^{\Pi})}{\mathbb{V}(\hat{\theta}_{HT}^{CR})} \times n$  of each design, based on the variance of the estimators across the  $K$  iterations. The ESS indicates, for every allocation design, the experimental sample size required to estimate the treatment effect with the same precision as with complete randomization. We see that almost every covariate-adaptive method does better than complete randomization, as they allow to reduce the sample size, often substantially. The only exception is stratification with many covariates. In general, if  $n < 2^p$ , stratification becomes worse than complete randomization. This phenomenon is due to the small strata issue and worsens with the strata fixed-effects estimator. Across different  $n$  sizes, we see that the matched-pairs design does better than the cube method when we balance a few covariates (up to three, in general). However, once

---

<sup>4</sup>We select 256 because exact inference methods for matched pair designs require a sample size divisible by four (Bai et al., 2022).

balancing more covariates, the cube method becomes more efficient. As expected, these relative gains to the matched-pairs design are more apparent for smaller  $n$  since the curse of dimensionality is more stringent. Notice there are clear gains from using the cube method even when allowing for non-linearities in the imputation of  $Y(1)$  and  $Y(0)$ . These results and those in Section 5 constitute suggestive evidence for a possible relaxation of Assumption 4.

Column 5 shows the results of using the double-lasso procedure. The method is quite effective in reducing the variance of the estimators, giving the same precision as the cube method or the matched-pairs designs with slightly fewer observations. However, for small  $n$ , the double-lasso procedure tends to give biased estimators (see Table D.2 in Appendix) and is unstable to different methods for creating the augmented dictionary (Kolesár et al., 2024).

Tables D.1-D.4 in the Appendix show additional results for each design: standard deviation and bias of the HT estimator, confidence coverage rate, and power for testing a null PATE. In particular, Table D.3 verifies that the coverage rate is exact for  $n$  large enough.

## 6 Practical Considerations about RCTs

### 6.1 Randomization and Baseline Information: How Do Researchers Randomize in Practice?

Bai (2022) indicates that among 5,000 RCTs in the AEA RCT Registry, more than 800 are stratified (i.e., about 16%). We complement this insight by gathering information from 104 randomized controlled trials (RCTs) published in top-5 and AEA journals<sup>5</sup> between 2019 and 2023. Specifically, we examined their collection of baseline information, baseline outcomes, and randomization methods. Figure 1 summarizes these details. Our findings indicate a lack of consensus among RCTs regarding the method used for allocating individuals to treatment. Most published papers (54%) employ a stratified design for allocation, followed by completely randomized designs (34%). A minority of researchers utilize alternative methods such as matching or re-randomization. However, there is substantial agreement regarding the collection of baseline data, with 90% of papers gathering information before treatment

---

<sup>5</sup>American Economic Review, AEJ: Applied Economics, AEJ: Economic Policy, AEJ: Macroeconomics, AEJ: Microeconomics, Econometrica, Journal of Political Economy, Quarterly Journal of Economics, and Review of Economic Studies

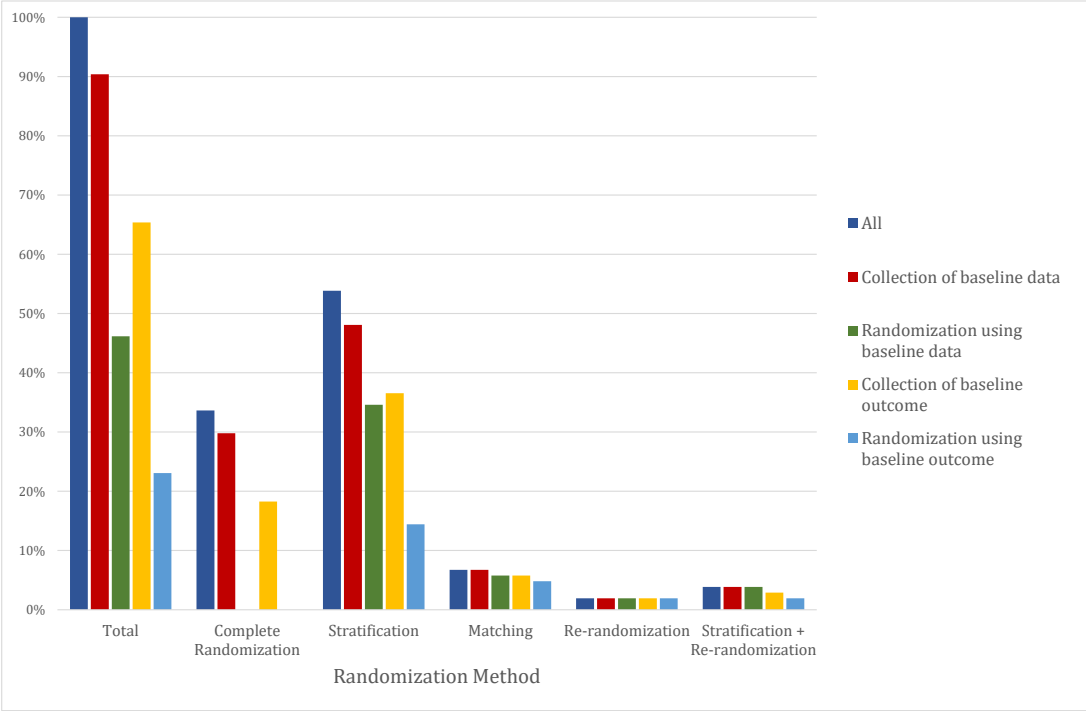
Table 1: Effect of number of covariates on effective sample size

Number of covariates	Complete Randomization	Stratified Randomization	Matched Pairs	Cube Method	Double Lasso	
	(1)	(2)	(3)	(4)	(5)	
$n = 100$	1	100.00	80.99	63.82	64.90	64.24
	2	–	81.66	63.57	67.89	–
	3	–	80.72	65.24	68.46	–
	5	–	90.57	67.28	67.37	–
	7	–	113.60	69.10	67.70	–
	9	–	160.25	73.14	65.31	–
	12	–	523.57	75.08	65.40	–
$n = 256$	1	256.00	215.49	162.69	167.76	163.12
	2	–	212.66	162.92	170.79	–
	3	–	215.07	163.50	169.53	–
	5	–	215.00	174.35	168.93	–
	7	–	240.92	180.31	171.15	–
	9	–	308.55	183.49	165.50	–
	12	–	610.22	194.93	174.00	–
$n = 500$	1	500.00	419.11	328.17	329.99	315.35
	2	–	408.31	328.15	325.72	–
	3	–	408.02	320.39	332.88	–
	5	–	418.08	326.86	331.19	–
	7	–	445.93	329.81	332.37	–
	9	–	514.09	349.28	334.98	–
	12	–	847.19	363.50	330.05	–
$n = 1000$	1	1000.00	828.85	646.01	669.16	642.44
	2	–	826.04	640.22	647.56	–
	3	–	830.47	641.96	638.66	–
	5	–	819.82	651.41	666.30	–
	7	–	845.31	657.62	663.35	–
	9	–	955.18	686.53	664.19	–
	12	–	1329.51	717.35	656.39	–

This table shows the effective sample size (ESS) for different allocation designs and experimental sample sizes. For each allocation design, the ESS corresponds to the number of observations needed to have the same precision as under complete randomization, equal to  $n \times (\widehat{\mathbb{V}}(\widehat{\theta}_{HT}^{\text{I}})/\widehat{\mathbb{V}}(\widehat{\theta}_{HT}^{\text{CR}}))$ , where we compute the variance estimates over 10,000 simulations. For columns 1, 3, and 4, the estimator used is the Horvitz-Thompson algorithm. Column 1 gives the size  $n$ , as the design used is complete randomization. For column 3, we assign treatment using the matched-pairs design, pairing individuals to the closest unit and using the Mahalanobis distance whenever more than one covariate is balanced. Column 4 shows the results for the cube method with two moments for each variable. For column 2, we run stratified randomization. We use median values for continuous variables and estimate the PATE using OLS with strata fixed-effects following Bugni et al. (2018). Finally, column 5 uses a double-lasso selection procedure as described in Belloni et al. (2014).

allocation. Nevertheless, this information is not always utilized during the randomization process, as only 46% of the papers leverage it to achieve covariate balance ex-ante. The remaining studies collect baseline data for balance tests, covariate adjustment in regression, and/or heterogeneity analysis of treatment effects. If outcomes of interest are relatively stable over time, empiricists should be interested in balancing their pre-treatment values, as they are highly likely to be correlated with potential outcomes. In our sample, 65% of researchers collect these variables, yet only 23% incorporate them into the allocation design, indicating an area for improvement in experimental design and inference. When various outcomes are considered in RCT, the curse of dimensionality arising in stratification, matched pair design, or re-randomization (where computational time could become prohibitive) may prevent empiricists from balancing on a large set of pre-treatment outcomes and sociodemographic covariates. The cube method could improve experimental randomization by allowing balancing on more variables than other methods.

Figure 6: Distribution of Randomization Methods for 104 RCTs in Top-5 + AEA Journals (2019-2023)



In sum, researchers forgo some precision gains by not using available covariates. However, we note that stratification and other covariate-adaptative methods seem to gain popularity or be increasingly favored by top journals. Furthermore, as explained in the next section, precision gains are a crucial reason, but there are others.

## 6.2 Removing the Back Luck

Researchers routinely provide “balance tables”, i.e., a comparison of the moments of available covariates between control and treatment. Due to bad luck, the empiricist can expect a certain proportion of imbalances. As explained, the likelihood of imbalances increases with the number of covariates used. There seems to be a grey area around the reporting of balance checks. In particular, pre-analysis plans often report no clear justification regarding the choice of the covariates to include in the balance table. Researchers are not very comfortable with reporting substantial unbalances. Snyder and Zhuo (2024) analyzing a large set of balance checks find that the editorial process removes an ample part of studies reporting imbalances, perhaps as much as 30%. Imbalances increase the risk of rejection by journals and, even worse, may provide incentives to engage in p-hacking (e.g., removing some covariates from the balance table).

How can “bad luck” be mitigated? Stratification obviously mitigates back luck by avoiding large imbalances. However, as explained, there is a limit to the number of covariates the researcher can include when using stratification. The cube method stands out in two respects. First, it allows using a large set of covariates, which is particularly relevant when the number of units is limited or when researchers do not know which covariates are the most important. But the second aspect is perhaps of greater interest. By construction, the cube will balance selected covariates with probability one. Put differently, all selected covariates will pass balance checks with probability one. Under stratification, imbalances are less likely than under completely randomized designs but would arise with a strictly positive probability. A more systematic use of recent methods, e.g., matched-pair designs or the cube method, will thus remove a clear publication bias and, as a result, make RCTs overall results less conservative.



### 6.3 Heterogenous Probabilities

The cube method can easily handle heterogeneous probabilities of assignment. This feature could be particularly relevant for various reasons. First, in view to minimize the variance of the estimator of the ATE, the optimal assignment probabilities corresponding to the so-called Neyman allocation are  $\pi(X) = \mathbb{V}(Y(1)|X)^{1/2} (\mathbb{V}(Y(1)|X)^{1/2} + \mathbb{V}(Y(0)|X)^{1/2})^{-1}$ . However, this optimal allocation is often not feasible without pilots and is irrelevant if the empiricist has various outcomes of interest. Even if  $\mathbb{V}(Y(1)|X) = \mathbb{V}(Y(0)|X)$  and if the empiricist wants to minimize the variance of estimates of  $\mathbb{E}(Y(1) - Y(0))$ , she could also adapt some assignment probabilities to heterogeneous costs  $c_1(X), c_0(X)$  of treatment and control to fulfill a budget constraint. Third, the empiricist could be interested in treatment effects on some subpopulations  $X = x$  that would not be precisely estimated if using a constant assignment probability. More generally, in an armed bandit perspective, empiricists may adapt assignment probabilities with respect to what they learned to maximize some objective, explore treatment effects on some subpopulations, and minimize regret. The cube method offers ample flexibility because empiricists could choose the probability of assignment of each unit.

### 6.4 Experimenter Choices and Tuning Parameters

When choosing a randomization method, researchers still have to make a few choices, the main one being which covariates to include. But there are also more subtle choices. For instance, stratification involves some choices when transforming continuous variables (e.g., income, age, test score, etc) into discrete ones (e.g., should one use median or quartile splits?). Similarly, generalization of matched-pairs design to heterogeneous probabilities involves approximation of probabilities assignment (as for instance in Cytrynbaum, 2023). The consequences of such choices are expected to be mild but may nonetheless cast doubts. In re-randomization strategies, the empiricist makes many choices (balancing criteria, stopping strategy, inferences). In the Gram-Schmidt random walk (Harshaw et al., 2024), the empiricist has to choose a tuning parameter. In contrast, choices when using the cube method boil down to selecting the set of covariates to be balanced. In that sense, the cube method is more transparent.

## 6.5 Computation and Packages

Packages in R for randomization are available here: <https://rdr.io/cran/BalancedSampling/>. Interestingly enough, the cube method is not computationally demanding. While running simulations included in the present paper, we systematically found the cube method to have an order of magnitude faster than the algorithms we used to implement other presented methods.

## 6.6 Ex-ante or Ex-post?

Economic intuition may suggest a set of covariates that might be important. However, there is some residual uncertainty regarding which covariates to include. Furthermore, researchers are often interested in several outcomes that might be affected by the treatment, expanding the set of relevant covariates. Ideally, one would resolve uncertainty regarding which covariates to include using the largest possible set of covariates, including available covariates, the interaction terms, and so on and forth. Obviously, the number of covariates to include may be large and even bigger than the number of observations. A possible strategy is thus to rely on techniques that allow researchers to infer treatment effects by *ex-post* controlling for a large set of covariates. For instance, (Belloni et al., 2014) propose a double-lasso approach to select relevant covariates among a very large set so as to estimate treatment effects. In sharp contrast, randomization *ex-ante* controls for some covariates. However, both approaches are not mutually exclusive, as it is possible to use some covariates at the randomization stage and implement still some ex-post estimation strategy. A more in-depth discussion would derive a more formal comparison between estimators based on ex-ante and ex-post balancing, respectively.

## 7 Conclusion

The cube method, first introduced for survey sampling by Deville and Tillé (2004), outperforms commonly used designs, especially when the number of covariates to be incorporated is large. We provide a set of results formalizing these gains in the RCT context. We tackle common issues empiricists face, such as balance, inference, sample size, and precision. Compared to other covariate-adaptive mechanisms, the cube method allows for better balance

and precision. More precise estimates significantly reduce the sample size needed for a minimum detectable effect and will make RCTs less conservative. By construction, the cube method guarantees the balance between treatment and control groups. The cube method thus makes balance tests unnecessary for covariates used at the randomization stage. By providing more balanced control and treatment, the cube method helps reduce publication bias and potential  $p$ -hacking.

A particular effort is devoted to comparing existing methods. We provide simulations and derive analytical results regarding the behavior of commonly used methods. We clarified how the number of covariates  $p$  to balance and the number of available units  $n$  affect balance for various randomization methods. To the best of our knowledge, for all the methods but the cube, imbalances grow as  $p/n$ . With the cube method, the imbalance is bounded by  $p^2/n^2$ . Such a systematic comparison might be relevant to both empiricists interested in comparing methods and researchers interested in the asymptotic behavior of randomization methods. For instance, there is a consensus to balance pre-treatment outcomes. Next, there is a need to simultaneously balance several variables when a single RCT is used to evaluate treatment effects on various outcomes. The cube method could be particularly relevant in this setup. Beyond the pre-treatment outcome, our review of pre-analysis plans of published papers in top economics journals shows no clear consensus on the selection of variables to balance. We understand this lack of consensus as a consequence of a trade-off between the balancing adjustment and the number of variables to balance for methods that are not well-suited to balance numerous variables. Because the cube method is more robust than the others when the number of balanced variables is large, this randomization method better addresses this trade-off.

At a more general level, we contribute to a recent stream of research on the usage of information available before randomization in RCTs. The information available to researchers comes in various forms. Some covariates are usually available, notably through the systematization of baseline surveys. As explained, the researcher can incorporate baseline covariates to improve balance and ensure some precision gains. Which covariates to select is often a matter of intuition or an “educated guess”. Sometimes, a pilot study is also available. Pilot studies allow insights into which covariates correlate the most with potential outcomes. Last, some RCTs use a design with repeated experiments, allowing an increasing gain of information across time. Reallocation of treated units among treatment arms via “armed-bandits” algorithms, for instance, speeds up the identification of treatment effects. The cube method is particularly well suited to take advantage of available information because

it permits incorporating a large set of covariates. However, its statistical properties are not yet fully explored for multiple treatment arms or repeated experiments, opening roads for future research.

## References

- Aaronson, J., Burton, R., Dehling, H., Gilat, D., Hill, T., and Weiss, B. (1996). Strong Laws for L- and U- Statistics. *Research Scholars in Residence*, 348.
- Athey, S. and Imbens, G. W. (2017). Chapter 3 - The Econometrics of Randomized Experiments. In Banerjee, A. V. and Duflo, E., editors, *Handbook of Economic Field Experiments*, volume 1 of *Handbook of Field Experiments*, pages 73–140. North-Holland.
- Bai, Y. (2022). Optimality of Matched-Pair Designs in Randomized Controlled Trials. *American Economic Review*, 112(12):3911–3940.
- Bai, Y., Romano, J. P., and Shaikh, A. M. (2022). Inference in Experiments With Matched Pairs. *Journal of the American Statistical Association*, 117(540):1726–1737. Publisher: Taylor & Francis \_eprint: <https://doi.org/10.1080/01621459.2021.1883437>.
- Bai, Y., Shaikh, A. M., and Tabord-Meehan, M. (2024). A Primer on the Analysis of Randomized Experiments and a Survey of some Recent Advances. arXiv:2405.03910 [econ, stat].
- Ball, S., Bogatz, G. A., Rubin, D. B., and Beaton, A. E. (1973). Reading with Television: An Evaluation of the Electric Company. A Report to the Children’s Television Workshop. Volumes 1 and 2. Technical Report PR-73-02, ETS Program Report.
- Belloni, A., Chernozhukov, V., and Hansen, C. (2014). Inference on Treatment Effects after Selection among High-Dimensional Controls†. *The Review of Economic Studies*, 81(2):608–650.
- Bruhn, M. and McKenzie, D. (2009). In Pursuit of Balance: Randomization in Practice in Development Field Experiments. *American Economic Journal: Applied Economics*, 1(4):200–232.
- Bugni, F. A., Canay, I. A., and Shaikh, A. M. (2018). Inference Under Covariate-Adaptive Randomization. *Journal of the American Statistical Association*, 113(524):1784–1796. Publisher: Taylor & Francis \_eprint: <https://doi.org/10.1080/01621459.2017.1375934>.
- Chen, J. and Rao, J. N. K. (2007). Asymptotic Normality Under Two-Phase Sampling Designs. *Statistica Sinica*, 17(3):1047–1064. Publisher: Institute of Statistical Science, Academia Sinica.

- Cytrynbaum, M. (2023). Optimal Stratification of Survey Experiments. arXiv:2111.08157 [econ, math, stat].
- Deville, J.-C. and Tillé, Y. (2004). Efficient balanced sampling: The cube method. *Biometrika*, 91(4):893–912.
- Deville, J.-C. and Tillé, Y. (2005). Variance approximation under balanced sampling. *Journal of Statistical Planning and Inference*, 128(2):569–591.
- Fisher, R. A. (1926). The arrangement of field experiments. *Journal of the Ministry of Agriculture*, 33:503–515. Publisher: Ministry of Agriculture and Fisheries.
- Fisher, S. R. A. (1935). *The Design of Experiments*. Oliver and Boyd.
- Gerber, A., Hoffman, M., Morgan, J., and Raymond, C. (2020). One in a Million: Field Experiments on Perceived Closeness of the Election and Voter Turnout. *American Economic Journal: Applied Economics*, 12(3):287–325.
- Grafström, A. and Tillé, Y. (2013). Doubly balanced spatial sampling with spreading and restitution of auxiliary totals: Doubly balanced spatial sampling. *Environmetrics*, 24(2):120–131.
- Greevy, R., Lu, B., Silber, J. H., and Rosenbaum, P. (2004). Optimal multivariate matching before randomization. *Biostatistics (Oxford, England)*, 5(2):263–275.
- Gut, A. (2013). *Probability: A Graduate Course*, volume 75 of *Springer Texts in Statistics*. Springer New York, New York, NY.
- Hahn, J. (1998). On the Role of the Propensity Score in Efficient Semiparametric Estimation of Average Treatment Effects. *Econometrica*, 66(2):315.
- Harshaw, C., Sävje, F., Spielman, D. A., and Zhang, P. (2024). Balancing Covariates in Randomized Experiments with the Gram–Schmidt Walk Design. *Journal of the American Statistical Association*, 0(0):1–13. Publisher: Taylor & Francis \_eprint: <https://doi.org/10.1080/01621459.2023.2285474>.
- Heckman, J., Moon, S. H., Pinto, R., Savelyev, P., and Yavitz, A. (2010). Analyzing social experiments as implemented: A reexamination of the evidence from the HighScope Perry Preschool Program. *Quantitative Economics*, 1(1):1–46. \_eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.3982/QE8>.

- Heckman, J. J., Pinto, R., Shaikh, A. M., and Yavitz, A. (2011). Inference with Imperfect Randomization: The Case of the Perry Preschool Program.
- Higgins, M. J., Sävje, F., and Sekhon, J. S. (2016). Improving massive experiments with threshold blocking. *Proceedings of the National Academy of Sciences*, 113(27):7369–7376. Publisher: Proceedings of the National Academy of Sciences.
- Imai, K., King, G., and Nall, C. (2009). The Essential Role of Pair Matching in Cluster-Randomized Experiments, with Application to the Mexican Universal Health Insurance Evaluation. *Statistical Science*, 24(1).
- Imbens, G. W. (2011). Experimental design for unit and cluster randomid trials. Technical report, Harvard University.
- Kolesár, M., Müller, U. K., and Roelsgaard, S. T. (2024). The Fragility of Sparsity. arXiv:2311.02299 [econ, stat].
- Lee, S. and Shaikh, A. M. (2014). Multiple Testing and Heterogeneous Treatment Effects: Re-Evaluating the Effect of Progresa on School Enrollment. *Journal of Applied Econometrics*, 29(4):612–626. \_eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/jae.2327>.
- Li, X., Ding, P., and Rubin, D. B. (2018). Asymptotic theory of rerandomization in treatment–control experiments. *Proceedings of the National Academy of Sciences*, 115(37):9157–9162. Publisher: Proceedings of the National Academy of Sciences.
- Morgan, K. L. and Rubin, D. B. (2012). Rerandomization to improve covariate balance in experiments. *The Annals of Statistics*, 40(2). arXiv:1207.5625 [math, stat].
- Snyder, C. M. and Zhuo, R. (2024). Examining Selection Pressures in the Publication Process through the Lens of Sniff Tests. *The Review of Economics and Statistics*, pages 1–45.
- Student (1938). Comparison Between Balanced and Random Arrangements of Field Plots. *Biometrika*, 29(3/4):363–378. Publisher: [Oxford University Press, Biometrika Trust].
- Takacs, L. (1991). A Moment Convergence Theorem. *The American Mathematical Monthly*, 98(8):742–746. Publisher: Mathematical Association of America.
- Tillé, Y. (2011). Ten years of balanced sampling with the cube method: An appraisal. *Survey Methodology*, 37(2):215–226.

Tillé, Y. and Favre, A.-C. (2004). Coordination, Combination and Extension of Balanced Samples. *Biometrika*, 91(4):913–927. Publisher: [Oxford University Press, Biometrika Trust].



# A Imbalances in Randomization Methods

In this appendix, we prove the theoretical results discussed in Section 4.1.2.

We consider  $(D_1, \dots, D_n) \in \{0; 1\}^n$  a random selection of some units among  $n$ . We will consider the following schemes:

1. Bernoulli randomization (BR): for any  $(d_1, \dots, d_n) \in \{0; 1\}^n$

$$\mathbb{P}((D_1, \dots, D_n) = (d_1, \dots, d_n) | (X_i)_{1 \leq i \leq n}) = 2^{-n}$$

2. Complete randomization (CR): for any  $(d_1, \dots, d_n) \in \{0; 1\}^n$

$$\mathbb{P}((D_1, \dots, D_n) = (d_1, \dots, d_n) | (X_i)_{1 \leq i \leq n}) = \binom{n}{n/2}^{-1} \mathbf{1} \left\{ \sum_{i=1}^n d_i = n/2 \right\}$$

3. Stratification (S- $\ell$ ): we consider  $q_p^k$  the empirical  $p$ -quantile of  $(X_{ki})_{i=1, \dots, n}$  and the partition  $H_n = \left\{ \prod_{k=1}^p [q_{j_k/\ell}^k; q_{(j_k+1)/\ell}^k] : \forall k \in \{1, \dots, p\}, j_k \in \{0, 1, \dots, \ell - 1\} \right\}$  of  $[0; 1]^k$  for some  $\ell \geq 2$  and  $\mathcal{S}_n = \{ \{i : X_i \in h\} : h \in H_n \text{ such that } \exists i : X_i \in h \}$  a partition of  $\{1, \dots, n\}$ . For any  $(d_1, \dots, d_n) \in \{0; 1\}^n$

$$\mathbb{P}((D_1, \dots, D_n) = (d_1, \dots, d_n) | (X_i)_{1 \leq i \leq n}) = \prod_{s \in \mathcal{S}_n} \binom{|s|}{\lfloor |s|/2 \rfloor}^{-1} \mathbf{1} \{ \lfloor |s|/2 \rfloor \leq \sum_{i \in s} d_i \leq \lfloor (|s|+1)/2 \rfloor \}.$$

4. Matched-pairs design (MP): we consider a partition  $\mathcal{S}_n$  of  $\{1, \dots, n\}$  such that

$$\mathcal{S}_n := \arg \min_{\mathcal{S}: |s|=2 \text{ for any } s \in \mathcal{S}} \sum_{s \in \mathcal{S}} \left\| \sum_{i, j \in s^2, i \neq j} X_i - X_j \right\|^2 \quad (\text{A.1})$$

5. Cube method with landing phase using linear programming (CM) as described in Section 3.

In all the previous designs, we have  $\mathbb{P}(D_i = 1 | (X_i)_{i=1, \dots, n}) = 1/2$  and next, the estimator  $B_{n,p}(X) = \frac{2}{n} \sum_{i=1}^n X_i D_i - \frac{2}{n} \sum_{i=1}^n X_i (1 - D_i)$  has expectation  $\mathbb{E}(X | D = 1) - \mathbb{E}(X | D = 0) = 0$ , meaning that these assignments balance  $X$  on average. To compare the balancing of these assignments, we will bound  $\mathbb{E}(\|B_{n,p}(X)\|^2)$  for each design.

**Proposition A.1**

Suppose that Assumption 3 holds.

1. Under assignment design (BR) we have:

$$\mathbb{E} (\|B_{n,p}(X)\|^2) = \frac{4}{n} \sum_{k=1}^p \mathbb{E} (X_{k1}^2) \quad (\text{A.2})$$

$$\text{and } \frac{4\underline{C}}{3} \frac{p}{n} \leq \mathbb{E} (\|B_{n,p}(X)\|^2) \leq \frac{4\overline{C}}{3} \frac{p}{n}. \quad (\text{A.3})$$

2. Under assignment design (CR) we have:

$$\mathbb{E} (\|B_{n,p}(X)\|^2) = \frac{4}{n} \sum_{k=1}^p \mathbb{V} (X_{k1}^2) \quad (\text{A.4})$$

$$\text{and } \frac{\underline{C}}{3} \frac{p}{n} \leq \mathbb{E} (\|B_{n,p}(X)\|^2) \leq \frac{\overline{C}}{3} \frac{p}{n}. \quad (\text{A.5})$$

3. Under assignment design (S- $\ell$ ), we have:

(a) if  $n\ell^{-p} \rightarrow \infty$ :

$$\|B_{n,p}(X)\|^2 = B_1^2 + o_p \left( \frac{p}{n} \right) \quad (\text{A.6})$$

$$\text{with } \frac{\underline{C}}{6\ell^2\overline{C}} \frac{p}{n} (1 - o(1)) \leq \mathbb{E} (B_1^2) \leq \frac{4}{n} \sum_{k=1}^p \mathbb{V} (X_{k1})$$

(b) if  $n\ell^{-p} \rightarrow 0$ :

$$\|B_{n,p}(X)\|^2 = B_2^2 + o_p \left( \frac{p}{n} \right). \quad (\text{A.7})$$

$$\text{with } \mathbb{E} (B_2^2) = \frac{4}{n} \sum_{k=1}^p \mathbb{E} (X_{k1}^2).$$

4. Under assignment design (MP), we have:

$$\frac{p}{n} \left( \frac{1}{3} - \sqrt{\frac{2 \ln(n-1) + 4 \ln \overline{C}}{p}} \right) \leq \mathbb{E} (\|B_{n,p}(X)\|^2) \leq \frac{4}{n} \sum_{k=1}^p \mathbb{V} (X_{1k}) \quad (\text{A.8})$$

5. Under assignment design (CM), we have:

$$\mathbb{E} (\|B_{n,p}(X)\|^2) \leq 4 \frac{(p+1)^2 \lambda_{\max}(M)}{n^2 \lambda_{\min}(M)}, \quad (\text{A.9})$$

for  $\lambda_{\max}(M)$  and  $\lambda_{\min}(M)$  the largest and the smallest eigenvalues of  $M$ .

Proof:

For all the assignment designs considered, we have

$$\begin{aligned}\mathbb{E}(\|B_{n,p}(X)\|^2) &= 4 \sum_{k=1}^p \mathbb{E} \left[ \left( \frac{1}{n} \sum_{i=1}^n X_{ki} (2D_i - 1) \right)^2 \right] \\ &= 4 \sum_{k=1}^p \mathbb{E} \left[ \left( \frac{1}{n} \sum_{i=1}^n \frac{X_{ki} D_i}{\mathbb{P}(D_i = 1 | (X_j)_{j=1, \dots, n})} - \bar{X}_k \right)^2 \right] \\ &= 4 \sum_{k=1}^p \mathbb{E} \left[ \mathbb{V} \left( \frac{1}{n} \sum_{i=1}^n \frac{X_{ki} D_i}{\mathbb{P}(D_i = 1 | (X_j)_{j=1, \dots, n})} \middle| (X_j)_{j=1, \dots, n} \right) \right].\end{aligned}$$

Under assignment design (BR), we have  $D_i \perp\!\!\!\perp D_j | (X_j)_{j=1, \dots, n}$  and  $\mathbb{E}(2D_i - 1 | (X_j)_{j=1, \dots, n}) = 0$ . Next, expanding the square we have

$$\mathbb{E} \left[ \left( \frac{1}{n} \sum_{i=1}^n X_{ki} (2D_i - 1) \right)^2 \middle| (X_j)_{j=1, \dots, n} \right] = \frac{1}{n^2} \sum_{i=1}^n X_{ki}^2.$$

Identical distribution across  $i$  ensures (A.2).

For any  $k \in \{1, \dots, p\}$  we have:

$$\begin{aligned}\frac{\underline{C}}{3} &\leq \frac{1}{3} \left( \underline{C} + \frac{(1 - \underline{C})^3}{(\bar{C} - \underline{C})^2} \right) = \int_0^1 u^2 \left( \bar{C} \mathbf{1}\{u \leq \frac{1 - \underline{C}}{\bar{C} - \underline{C}}\} + \underline{C} \mathbf{1}\{u \geq \frac{1 - \underline{C}}{\bar{C} - \underline{C}}\} \right) du \\ &\leq \mathbb{E}(X_{k1}^2) \leq \int_0^1 u^2 \left( \underline{C} \mathbf{1}\{u \leq \frac{\bar{C} - 1}{\bar{C} - \underline{C}}\} + \bar{C} \mathbf{1}\{u \geq \frac{\bar{C} - 1}{\bar{C} - \underline{C}}\} \right) du = \frac{1}{3} \left( \bar{C} - \frac{(\bar{C} - 1)^3}{(\bar{C} - \underline{C})^2} \right) \leq \frac{\bar{C}}{3},\end{aligned}$$

and next (A.3) follows.

Under assignment design (CR), we have:

$$\mathbb{V} \left( \frac{1}{n} \sum_{i=1}^n \frac{X_{ki} D_i}{\mathbb{P}(D_i = 1 | (X_j)_{j=1, \dots, n})} \middle| (X_j)_{j=1, \dots, n} \right) = \left(1 - \frac{1}{2}\right) \frac{2}{n} S_k^2$$

for  $S_k^2 = \frac{1}{n-1} \sum_{i=1}^n (X_{ki} - \bar{X}_k)^2$  and  $\bar{X}_k = \frac{1}{n} \sum_{i=1}^n X_{ki}$ . Because  $\mathbb{E}(S_k^2) = \mathbb{V}(X_{k1})$ , A.4 follows.

A.5 is ensured by

$$\begin{aligned}\frac{\underline{C}}{12} &\leq \frac{1}{12} \left( \underline{C} + \frac{(1 - \underline{C})^3}{(\bar{C} - \underline{C})^2} \right) \\ &= \int_0^1 \left( u - \frac{1}{2} \right)^2 \left( \bar{C} \mathbf{1} \left\{ \left| u - \frac{1}{2} \right| \leq \frac{1 - \underline{C}}{2(\bar{C} - \underline{C})} \right\} + \underline{C} \mathbf{1} \left\{ \left| u - \frac{1}{2} \right| \geq \frac{1 - \underline{C}}{2(\bar{C} - \underline{C})} \right\} \right) du \\ &\leq \mathbb{V}(X_{k1}) \leq \int_0^1 \left( u - \frac{1}{2} \right)^2 \left( \underline{C} \mathbf{1} \left\{ \left| u - \frac{1}{2} \right| \leq \frac{\bar{C} - 1}{2(\bar{C} - \underline{C})} \right\} + \bar{C} \mathbf{1} \left\{ \left| u - \frac{1}{2} \right| \geq \frac{\bar{C} - 1}{2(\bar{C} - \underline{C})} \right\} \right) du\end{aligned}$$

$$= \frac{1}{12} \left( \bar{C} - \frac{(\bar{C} - 1)^3}{(\bar{C} - \underline{C})^2} \right) \leq \frac{\bar{C}}{12}.$$

Under assignment design (S- $\ell$ ), for  $s \in \mathcal{S}_n$  let  $m_s$  the number of units assigned to the treatment in strata  $s$ . Similarly, for  $h \in \mathcal{H}_n$  let  $m_h$  the number of units assigned to the treatment in strata  $\{i : X_i \in h\}$  and  $m = \sum_{s \in \mathcal{S}_n} m_s = \sum_{h \in \mathcal{H}_n} m_h$ . Last,  $n_h$  denotes  $|\{i : X_i \in h\}|$  for  $h \in \mathcal{H}_n$ . We have

$$\begin{aligned} & \mathbb{V} \left( \frac{1}{n} \sum_{i=1}^n \frac{X_{ki} D_i}{\mathbb{P}(D_i = 1 | (X_j)_{j=1, \dots, n})} \middle| (X_j)_{j=1, \dots, n}, (m_s)_{s \in \mathcal{S}} \right) \\ &= \sum_{s \in \mathcal{S}_n} \left[ \left( \frac{m_s}{m} \right)^2 \left( 1 - \frac{m_s}{|s|} \right) \frac{1}{m_s} S_{ks}^2 \mathbb{1}\{|s| \geq 2\} + \frac{1}{n^2} \sum_{i \in s} X_{ki}^2 \mathbb{1}\{|s| = 1\} \right] \\ &= \sum_{h \in \mathcal{H}_n} \left[ \left( \frac{m_h}{m} \right)^2 \left( 1 - \frac{m_h}{n_h} \right) \frac{1}{m_h} S_{kh}^2 \mathbb{1}\{n_h \geq 2\} + \frac{1}{n^2} \left( \sum_{i \in \{j : X_j \in h\}} X_{ki}^2 \right) \mathbb{1}\{n_h = 1\} \right] \end{aligned}$$

with  $S_{ks}^2 = \frac{1}{|s|-1} \sum_{i \in s} (X_{ki} - \bar{X}_{ks})^2$  for  $\bar{X}_{ks} = \frac{1}{|s|} \sum_{i \in s} X_{ki}$  and  $S_{kh}^2 = \frac{1}{n_h-1} \sum_{i \in \{j : X_j \in h\}} (X_{ki} - \bar{X}_{kh})^2$  for  $\bar{X}_{kh} = \frac{1}{n_h} \sum_{i \in \{j : X_j \in h\}} X_{ki}$ .  $(n_h)_{h \in \mathcal{H}_n}$  follows a multinomial distribution of parameters  $(n, \ell^p, (p_h)_{h \in \mathcal{H}})$  with  $\frac{\underline{C}}{\bar{C}} \ell^{-p} \leq \frac{\underline{C}}{\underline{C} + \bar{C}(\ell^p - 1)} \leq p_h \leq \frac{\bar{C}}{\bar{C} + \underline{C}(\ell^p - 1)} \leq \frac{\bar{C}}{\underline{C}} \ell^{-p}$ . Then for any  $h \in \mathcal{H}_n$ , we have  $\mathbb{P}(n_h = 1) = np_h(1 - p_h)^{n-1}$  and  $\mathbb{P}(n_h \geq 2) = 1 - (1 - p_h)^n - np_h(1 - p_h)^{n-1}$ . Moreover,  $(m_h)_{h \in \mathcal{H}_n} \perp\!\!\!\perp (S_{kh'})_{h' \in \mathcal{H}_n, k=1, \dots, p} | (n_{h'})_{h' \in \mathcal{H}_n}$ ,  $\mathbb{E}(S_{kh}^2 | (n_{h'})_{h' \in \mathcal{H}_n}) = \mathbb{V}(X_{k1} | X_1 \in h) \in \left[ \frac{\underline{C}}{12\bar{C}\ell^2}; \frac{\bar{C}}{12\underline{C}\ell^2} \right]$  and, conditional on  $(n_h)_{h \in \mathcal{H}_n}$ ,  $m_h$  are independent across  $h$  with probability distribution  $\frac{1}{2} \delta_{\lfloor \frac{n_h}{2} \rfloor} + \frac{1}{2} \delta_{\lfloor \frac{n_h+1}{2} \rfloor}$ . It follows that  $\mathbb{E}(m_h | (n_{h'})_{h' \in \mathcal{H}_n}) = n_h/2$ ,  $\mathbb{V}(m_h | (n_{h'})_{h' \in \mathcal{H}_n}) = (\lfloor (n_h + 1)/2 \rfloor - \lfloor n_h/2 \rfloor)^2/4 \leq 1/4$ ,  $\mathbb{E}(m | (n_{h'})_{h' \in \mathcal{H}_n}) = n/2$ ,  $\mathbb{V}(m | (n_{h'})_{h' \in \mathcal{H}_n}) \leq \ell^p/4$ . Chebyshev's inequality implies  $\mathbb{P}(|m - n/2| \geq M) \leq \frac{\ell^p}{4M^2}$ . Next,  $m = \frac{n}{2} (1 + O_p(n^{-1}\ell^{p/2})) = \frac{n}{2} (1 + o_p(1))$  and

$$\sum_{h \in \mathcal{H}_n} \left( \frac{m_h}{m} \right)^2 \left( 1 - \frac{m_h}{n_h} \right) \frac{1}{m_h} S_{kh}^2 \mathbb{1}\{n_h \geq 2\} = \frac{4}{n^2} (1 + o_p(1)) \sum_{h \in \mathcal{H}_n} m_h \left( 1 - \frac{m_h}{n_h} \right) S_{kh}^2 \mathbb{1}\{n_h \geq 2\}$$

$$\text{and } \frac{n_h}{8} \mathbb{1}\{n_h \geq 2\} \leq \frac{n_h^2 - 1}{4n_h} \mathbb{1}\{n_h \geq 2\} \leq m_h \left( 1 - \frac{m_h}{n_h} \right) \mathbb{1}\{n_h \geq 2\} \leq \frac{n_h}{4} \mathbb{1}\{n_h \geq 2\}.$$

If  $n\ell^{-p} \rightarrow \infty$  as  $n$  and  $p$  increase, we have  $np_h \geq \frac{\underline{C}}{\bar{C}} n\ell^{-p} \rightarrow \infty$  and next  $\sup_{h \in \mathcal{H}_n} \mathbb{P}(n_h = 1) \leq \sup_{h \in \mathcal{H}_n} np_h e^{(n-1)\ln(1-p_h)} \leq \sup_{h \in \mathcal{H}_n} np_h e^{-(n-1)p_h} \leq n\ell^{-p} \frac{\bar{C}}{\underline{C}} \exp(-(n-1)\ell^{-p}\underline{C}/\bar{C})$  and because  $|\mathcal{H}_n| = \ell^p$  and  $|X_{ki}^2| \leq 1$ , we have:

$$\begin{aligned} \mathbb{E} \left[ \sum_{k=1}^p \sum_{h \in \mathcal{H}_n} \frac{1}{n^2} \sum_{i \in \{j : X_j \in h\}} X_{ki}^2 \mathbb{1}\{n_h = 1\} \right] &\leq \frac{p}{n} \frac{|\mathcal{H}_n|}{n} \sup_{h \in \mathcal{H}_n} \mathbb{P}(\mathbb{1}\{n_h = 1\}) \\ &\leq \frac{p}{n} \frac{\bar{C}}{\underline{C}} \exp\left(\frac{\underline{C}}{\ell^p \bar{C}}\right) \exp\left(-n\ell^{-p} \frac{\underline{C}}{\bar{C}}\right) = o\left(\frac{p}{n}\right). \end{aligned}$$

We have  $\mathbb{P}(|m - n/2| \geq M) \leq \frac{\ell^p}{4M^2}$ , next  $m = \frac{n}{2} (1 + O_p(1/(n\ell^{p/2})))$  and

$$\begin{aligned} \sum_{h \in \mathcal{H}_n} \left(\frac{m_h}{m}\right)^2 \left(1 - \frac{m_h}{n_h}\right) \frac{1}{m_h} S_{kh}^2 \mathbf{1}\{n_h \geq 2\} &= \frac{4}{n^2} (1 + o_p(1)) \sum_{h \in \mathcal{H}_n} m_h \left(1 - \frac{m_h}{n_h}\right) S_{kh}^2 \mathbf{1}\{n_h \geq 2\} \\ &\leq \frac{1 + o_p(1)}{n} \sum_{h \in \mathcal{H}_n} \frac{n_h}{n} S_{kh}^2 \mathbf{1}\{n_h \geq 2\} \\ &\leq \frac{1 + o_p(1)}{n} \sum_{h \in \mathcal{H}_n} \frac{n_h}{n} S_{kh}^2 \end{aligned}$$

Last, note that  $\mathbb{E}(S_{kh}^2 | (n_{h'})_{h' \in \mathcal{H}_n}) = \mathbb{V}(X_k | X \in h)$  and  $\mathbb{E}(n_h/n) = p_h$  ensuring that  $\mathbb{E}(\sum_{h \in \mathcal{H}_n} \frac{n_h}{n} S_{kh}^2) = \mathbb{E}(\mathbb{V}(X_k | (\mathbf{1}\{X \in h\})_{h \in \mathcal{H}_n})) \leq \mathbb{V}(X_k)$ . Next for  $B_1^2 = \frac{16}{n^2} \sum_{k=1}^p \sum_{h \in \mathcal{H}_n} m_h \left(1 - \frac{m_h}{n_h}\right) S_{kh}^2 \mathbf{1}\{n_h \geq 2\}$ , we have  $\mathbb{E}(B_1^2) \leq \frac{4}{n} \sum_{k=1}^p \mathbb{E}(\mathbb{V}(X_k | \mathbf{1}\{X_1 \in h\})_{h \in \mathcal{H}_n}) = O(\frac{p}{n})$ . It follows that  $o_p(B_1^2) = o_p(\frac{p}{n})$  and this prove (A.6) and the upper bound of  $\mathbb{E}(B_1^2)$ . It remains to show the lower bound for  $\mathbb{E}(B_1^2)$ .

$$\begin{aligned} \mathbb{E}(B_1^2) &\geq \frac{2p}{n^2} \frac{\underline{C}}{12\ell^2\overline{C}} \mathbb{E}\left(\sum_{h \in \mathcal{H}_n} n_h \mathbf{1}\{n_h \geq 2\}\right) \\ &\geq \frac{p}{n^2} \frac{\underline{C}}{6\ell^2\overline{C}} \mathbb{E}\left(\sum_{h \in \mathcal{H}_n} n_h (1 - \mathbf{1}\{n_h = 1\})\right) \\ &\geq \frac{p}{n} \frac{\underline{C}}{6\ell^2\overline{C}} \left(1 - \frac{\ell^p}{n} \sup_{h \in \mathcal{H}_n} \mathbb{P}(n_h = 1)\right) \\ &\geq \frac{p}{n} \frac{\underline{C}}{6\ell^2\overline{C}} \left(1 - \frac{\overline{C}}{\underline{C}} \exp\left(- (n-1)\ell^{-p} \frac{\underline{C}}{\overline{C}}\right)\right). \end{aligned}$$

If  $n\ell^{-p} \rightarrow 0$ , for any  $h \in \mathcal{H}_n$  we have  $\mathbb{P}(n_h \geq 2) \leq 1 - (1 - p_h)^n - np_h(1 - p_h)^{n-1} \leq 1 - (1 - np_h) - np_h(1 - (n-1)p_h) \leq n^2 p_h^2 \leq \frac{\overline{C}^2}{\underline{C}^2} (n\ell^{-p})^2$ . This ensures that  $\sup_{h \in \mathcal{H}_n} \mathbb{P}(n_h \geq 2) = o(1)$ .

Next, we have:

$$\begin{aligned} \mathbb{E}\left[\frac{1}{n} \sum_{h \in \mathcal{H}_n} n_h \mathbf{1}\{n_h \geq 2\}\right] &\leq n^{-1} \ell^p \sup_{h \in \mathcal{H}_n} \mathbb{E}(n_h \mathbf{1}\{n_h \geq 2\}) \\ &= n^{-1} \ell^p \sup_{h \in \mathcal{H}_n} (\mathbb{E}(n_h) - \mathbb{P}(n_h = 1)) \\ &= n^{-1} \ell^p \sup_{h \in \mathcal{H}_n} (np_h (1 - (1 - p_h)^{n-1})) \\ &\leq \ell^p (n-1) \sup_{h \in \mathcal{H}_n} p_h^2 \\ &\leq \frac{\overline{C}^2}{\underline{C}^2} n\ell^{-p} = o(1). \end{aligned}$$

We have  $S_{kh}^2 \mathbf{1}\{n_h \geq 2\} \leq 2\mathbf{1}\{n_h \geq 2\}$ ,  $m = \frac{n}{2}(1 + o_p(1))$  and  $m_h \left(1 - \frac{m_h}{n_h}\right) \mathbf{1}\{n_h \geq 2\} \leq \frac{n_h}{4} \mathbf{1}\{n_h \geq 2\}$  and next

$$\sup_{k=1, \dots, p} \sum_{h \in \mathcal{H}_n} \left(\frac{m_h}{m}\right)^2 \left(1 - \frac{m_h}{n_h}\right) \frac{1}{m_h} S_{kh}^2 \mathbf{1}\{n_h \geq 2\} \leq \frac{4}{n^2} (1 + o_p(1)) \sum_{h \in \mathcal{H}_n} \frac{n_h}{2} \mathbf{1}\{n_h \geq 2\}.$$

Then,  $4 \sum_{k=1}^p \sum_{h \in \mathcal{H}_n} \frac{m_h^2}{m^2} \left(1 - \frac{m_h}{n_h}\right) \frac{1}{m_h} S_{kh}^2 \mathbf{1}\{n_h \geq 2\} = o_p\left(\frac{p}{n}\right)$ . On the other hand,

$$\begin{aligned} \left| \frac{1}{n} \sum_{h \in \mathcal{H}_n} \left( \sum_{i \in \{j: X_j \in h\}} X_{ki}^2 \right) \mathbf{1}\{n_h = 1\} - \frac{1}{n} \sum_{i=1}^n X_{ki}^2 \right| &= \frac{1}{n} \sum_{h \in \mathcal{H}_n} \left( \sum_{i \in \{j: X_j \in h\}} X_{ki}^2 \right) \mathbf{1}\{n_h \geq 2\} \\ &\leq \frac{1}{n} \sum_{h \in \mathcal{H}_n} n_h \mathbf{1}\{n_h \geq 2\}. \end{aligned}$$

It follows that  $4 \sum_{k=1}^p \mathbb{E} \left( \frac{1}{n^2} \sum_{h \in \mathcal{H}_n} \left( \sum_{i \in \{j: X_j \in h\}} X_{ki}^2 \right) \mathbf{1}\{n_h = 1\} \right) = \frac{4}{n} \sum_{p=1}^k \mathbb{E}(X_{k1}^2) + o\left(\frac{p}{n}\right)$ . This proves (A.7).

Under assignment design (MP), partition  $\mathcal{S}_n$  is adaptive to the sample and  $\mathcal{S}_\setminus$  is  $(X_i)_{i=1, \dots, n}$ -measurable. All element of  $\mathcal{S}_n$  are of size 2 and in each strata  $s$  one unit over two is randomly assigned to the treatment. Next, the formula of the variance of the Horwitz-Thompson estimator of an average for a stratified sampling ensures

$$\begin{aligned} \mathbb{E} \left( \|B_{n,p}(X)\|^2 | (X_j)_{j=1, \dots, n} \right) &= 4 \sum_{k=1}^p \sum_{s \in \mathcal{S}_n} \left(\frac{2}{n}\right)^2 \left(1 - \frac{1}{2}\right) S_{ks}^2 \\ &= \frac{4}{n} \sum_{k=1}^p \frac{2}{n} \sum_{s \in \mathcal{S}_n} S_{ks}^2 \end{aligned}$$

with  $S_{ks}^2 = \left(X_{ki} - \frac{X_{ki} + X_{ki'}}{2}\right)^2 + \left(X_{ki'} - \frac{X_{ki} + X_{ki'}}{2}\right)^2 = \frac{1}{2} (X_{ki} - X_{ki'})^2$  for  $i$  and  $i'$  such that  $s = \{i, i'\}$ . According to Lemma 2 in Bai (2022), complete randomization could be implemented as a two-stage random process. In a first step, a random partition  $\mathcal{S}_n^*$  is selected with uniform probability among all the partitions such that  $|s| = 2$  for any  $s \in \mathcal{S}_n^*$ . In a second stage, stratified sampling is used in each strata  $s \in \mathcal{S}_n^*$ . According to (A.4), we have  $\frac{4}{n} \sum_{k=1}^p \mathbb{V}(X_{k1}) = \mathbb{E} \left( \frac{8}{n^2} \sum_{s \in \mathcal{S}_n^*} \sum_{k=1}^p S_{ks}^2 \right)$ . But program (A.1) is equivalent to  $\sum_{s \in \mathcal{S}_n} \sum_{k=1}^p S_{ks}^2 \leq \sum_{s \in \mathcal{S}_n^*} \sum_{k=1}^p S_{ks}^2$ . Next, the right hand side of (A.8) follows. To prove the left hand side of (A.8), note that  $\mathbb{E} \left( \|B_{n,p}(X)\|^2 | (X_j)_{j=1, \dots, n} \right) = \frac{4}{n^2} \sum_{s \in \mathcal{S}_n} \text{diam}^2(s, \mathbb{R}^p)$  for  $\text{diam}^2(s, \mathbb{R}^p) = \sup_{(i,j) \in s^2} \|X_i - X_j\|^2$ . Let  $d_i$  the distance in  $\mathbb{R}^p$  of the unit  $i$  to its nearest neighbor  $d_i = \min_{j \in \{1, \dots, n\} \setminus \{i\}} \|X_j - X_i\|$ . We have  $\text{diam}^2(\{i, i'\}, \mathbb{R}^p) \geq \max(d_i^2, d_{i'}^2) \geq \frac{1}{2}(d_i^2 + d_{i'}^2)$  and next  $\mathbb{E} \left( \|B_{n,p}(X)\|^2 | (X_j)_{j=1, \dots, n} \right) \geq \frac{2}{n^2} \sum_{i=1}^n d_i^2 = \frac{p}{3n} + \frac{2}{n^2} \sum_{i=1}^n \min_{j \in \{1, \dots, n\} \setminus \{i\}} \sum_{k=1}^p Z_{kij}$  for  $Z_{kij} = (X_{ki} - X_{kj})^2 - \frac{1}{6}$ . Let  $\tilde{Z}_{kij} = (U_{ki} - U_{kj})^2 - \frac{1}{6}$  for  $(U_{ki})_{k=1, \dots, p, i=1, \dots, n}$  some independent uniform variables on  $[0; 1]$ . For any  $k$  and any  $i \neq j$ , the  $\tilde{Z}_{kij}$  are zero mean variables

with bounded support  $[-1/6; 5/6]$ . A zero mean variable  $Z$  is sub-Gaussian with parameter  $\nu > 0$  if for any  $\lambda \in \mathbb{R}$ ,  $\mathbb{E}(e^{\lambda Z}) \leq e^{\lambda^2 \nu^2 / 2}$ . Hoeffding's lemma ensures that a zero mean variable with bounded support included in  $[a; b]$  is sub-Gaussian with parameter  $(b - a)/2$ . It follows that  $\tilde{Z}_{kij}$  and  $-\tilde{Z}_{kij}$  are sub-Gaussian with parameter  $1/2$ . Moreover, because  $f_X(x_{1i}, \dots, x_{pi})f_X(x_{1j}, \dots, x_{pj}) \leq \bar{C}^2$  and because  $(\tilde{Z}_{kij})_{k=1, \dots, p}$  are independent across  $k$ , we have for any  $i \neq j$ :

$$\begin{aligned} \mathbb{E}\left(e^{-\lambda \sum_{k=1}^p Z_{kij}}\right) &\leq \bar{C}^2 \mathbb{E}\left(e^{-\lambda \sum_{k=1}^p \tilde{Z}_{kij}}\right) \\ &= \bar{C}^2 \prod_{k=1}^p \mathbb{E}\left(e^{-\lambda \tilde{Z}_{kij}}\right) \\ &\leq \bar{C}^2 e^{\lambda^2 p / 8}. \end{aligned}$$

For any  $i \in \{1, \dots, n\}$ , by Jensen inequality, we have

$$\begin{aligned} e^{\lambda \mathbb{E}(\max_{j \neq i} - \sum_{k=1}^p Z_{kij})} &\leq \mathbb{E}\left(e^{\lambda \max_{j \neq i} - \sum_{k=1}^p Z_{kij}}\right) \\ &\leq \sum_{j \neq i} \mathbb{E}\left(e^{-\lambda \sum_{k=1}^p Z_{kij}}\right) \\ &\leq (n - 1) \bar{C}^2 e^{\lambda^2 p / 8}. \end{aligned}$$

Next,  $\mathbb{E}(\max_{j \neq i} (-\sum_{k=1}^p Z_{kij})) \leq \frac{\ln(\bar{C}^2(n-1))}{\lambda} + \lambda \frac{p}{8}$  for any  $\lambda > 0$ . Choosing  $\lambda = \sqrt{\frac{8 \ln(\bar{C}^2(n-1))}{p}}$ , we have  $\mathbb{E}(\max_{j \neq i} (-\sum_{k=1}^p Z_{kij})) \leq \sqrt{p \ln(\bar{C}^2(n-1))}/2$  or equivalently  $\mathbb{E}(\min_{j \neq i} \sum_{k=1}^p Z_{kij}) \geq -\sqrt{p \ln(\bar{C}^2(n-1))}/2$ . From what precedes, we have

$$\begin{aligned} \mathbb{E}(\|B_{n,p}(X)\|^2) &\geq \frac{p}{3n} + \frac{2}{n^2} \sum_{i=1}^n \mathbb{E}\left(\min_{j \in \{1, \dots, n\} \setminus \{i\}} \sum_{k=1}^p Z_{kij}\right) \\ &\geq \frac{p}{3n} - \frac{2}{n} \sqrt{p \ln(\bar{C}^2(n-1))}/2 \\ &= \frac{p}{n} \left(\frac{1}{3} - \sqrt{\frac{2 \ln(n-1) + 4 \ln \bar{C}}{p}}\right). \end{aligned}$$

This achieves the proof of (A.8).

Concerning assignment (CM), let  $Z_i = (1, X'_i)'$ . Then, we have  $\|B_{n,p}(X)\|^2 \leq \|B_{n,p}(Z)\|^2$ . At the end of the flying phase of the cube method, all units of a  $A \subset \{1, \dots, n\}$  such that  $|A| = n - \dim(Z) = n - p - 1$  have been allocated. This means that  $D_i$  has been drawn for any  $i \in A$ . For  $i \notin A$ , some random variables  $\pi_i^*$  have been generated such that  $\mathbb{E}(\pi_i^* | (Z_i)_{i=1, \dots, n}) = 1/2$  and

$$\sum_{i \in A} Z_i(2D_i - 1) + \sum_{i \notin A} Z_i(2\pi_i^* - 1) = 0,$$

or equivalently:

$$B_{n,p}(Z) = \frac{4}{n} \sum_{i \notin A} Z_i (D_i - \pi_i^*).$$

In the landing phase of the cube algorithm, given  $W = (A, (D_i)_{i \in A}, (\pi_i^*)_{i \notin A}, (Z_i)_{i=1, \dots, n})$ , the cube method samples  $(D_i)_{i \notin A}$  such that for any  $i \notin A$ ,  $D_i | W$  follows a Bernoulli of mean  $\pi_i^*$ . But the sampling of  $(D_i)_{i \notin A}$  is not independent across  $i \notin A$ . Indeed, sampling probabilities are correlated in view to minimize

$$\mathbb{E} \left( \sum_{i \notin A} Z_i (D_i - \pi_i^*) M \sum_{i \notin A} Z_i (D_i - \pi_i^*) | W \right)$$

where  $M$  is a symmetric-positive matrix  $\dim(Z) \times \dim(Z)$ . When  $M = Id$ , this means that sampling probabilities of  $(D_i)_{i \notin A}$  during the landing phase are coordinated to minimize  $\mathbb{E} (\| \sum_{i \notin A} Z_i (D_i - \pi_i^*) \|^2 | W)$ . Next, this means that

$$\mathbb{E} \left( \left\| \sum_{i \notin A} Z_i (D_i - \pi_i^*) \right\|^2 | W \right) \leq \mathbb{E} \left( \left\| \sum_{i \notin A} Z_i (\tilde{D}_i - \pi_i^*) \right\|^2 | W \right)$$

where  $(\tilde{D}_i)_{i \notin A} | W$  are sampled as some independent Bernoulli of mean  $\pi_i^*$ . If  $M \neq Id$ , we have

$$\begin{aligned} \mathbb{E} \left( \left\| \sum_{i \notin A} Z_i (D_i - \pi_i^*) \right\|^2 | W \right) &\leq \frac{1}{\lambda_{\min}(M)} \mathbb{E} \left( \left\| M^{1/2} \sum_{i \notin A} Z_i (D_i - \pi_i^*) \right\|^2 | W \right) \\ &\leq \frac{1}{\lambda_{\min}(M)} \mathbb{E} \left( \left\| M^{1/2} \sum_{i \notin A} Z_i (\tilde{D}_i - \pi_i^*) \right\|^2 | W \right) \\ &\leq \frac{\lambda_{\max}(M)}{\lambda_{\min}(M)} \mathbb{E} \left( \left\| \sum_{i \notin A} Z_i (\tilde{D}_i - \pi_i^*) \right\|^2 | W \right). \end{aligned}$$

To conclude, note that

$$\begin{aligned} \mathbb{E} \left( \left\| \sum_{i \notin A} Z_i (\tilde{D}_i - \pi_i^*) \right\|^2 | W \right) &= \sum_{k=1}^{p+1} \mathbb{E} \left[ \left( \sum_{i \notin A} Z_{ki} (\tilde{D}_i - \pi_i^*) \right)^2 | W \right] \\ &= \sum_{k=1}^{p+1} \mathbb{E} \left[ \sum_{i \notin A} Z_{ki}^2 \pi_i^* (1 - \pi_i^*) | W \right] \\ &\leq \frac{(p+1)|A|}{4} = \frac{(p+1)^2}{4}. \end{aligned}$$



## B Proofs of Propositions

### B.1 Proof of Balancing Approximations for the Cube Method (Proposition 1)

From Assumption 2 and Proposition 4 in Deville and Tillé (2004), we have:

$$\left| \frac{1}{n} \sum_{i=1}^n \frac{X_{ji} D_i}{\pi_i} - \frac{1}{n} \sum_{i=1}^n X_{ji} \right| \leq \frac{q}{n} \max_{i=1, \dots, n} \left| \frac{X_{ji}}{\pi_i} \right| \leq \frac{q \max_{i=1, \dots, n} |X_{ji}|}{cn}.$$

If Assumption 1 holds and if moments of order  $r$  exist for  $X_{j1}$ , from Proposition 1.5 and Theorem 2.1 in Chapter 6 in Gut (2013), we have  $\max_{i=1, \dots, n} |X_{ji}| = o_p(n^{1/r})$ . If  $X_{ji}$  sub-Gaussian  $\max_{i=1, \dots, n} |X_{ji}| = O_p(\sqrt{\ln(n)})$  and if  $X_{ji}$  is bounded by  $K$ ,  $\max_{i=1, \dots, n} |X_{ji}| \leq K$ .

### B.2 Proof of Asymptotic Normality for the SATE (First part of Proposition 3)

We want to prove

$$\sqrt{n} (\widehat{\theta}_{HT} - \theta_0) \xrightarrow{d} \mathcal{N}(0, V_0)$$

and

$$\sqrt{n} (\widehat{\theta}_H - \theta_0) \xrightarrow{d} \mathcal{N}(0, V_0)$$

where  $V_0 = \mathbb{E} \left[ \pi_i (1 - \pi_i) \left( \frac{\varepsilon_i(1)}{\pi_i} + \frac{\varepsilon_i(0)}{1 - \pi_i} \right)^2 \right]$ .

We first show that it is sufficient to prove asymptotic normality for one of the two estimators. Proposition 1 ensures that if the empiricist includes a constant in the set of covariates to balance, one has  $\sum_{i=1}^n \frac{D_i}{\pi_i} = n + o_p\left(\frac{1}{\sqrt{n}}\right)$  and  $\sum_{i=1}^n \frac{1 - D_i}{1 - \pi_i} = n + o_p\left(\frac{1}{\sqrt{n}}\right)$ . Then, the Hájek estimator in (6) is given by

$$\begin{aligned} \widehat{\theta}_H &= \frac{1}{n + o_p\left(\frac{1}{\sqrt{n}}\right)} \left( \sum_{i=1}^n \frac{Y_i D_i}{\pi_i} - \frac{Y_i (1 - D_i)}{1 - \pi_i} \right) \\ &= \frac{n}{n + o_p\left(\frac{1}{\sqrt{n}}\right)} \widehat{\theta}_{HT} \\ &= \left( 1 + o_p\left(\frac{1}{\sqrt{n}}\right) \right) \widehat{\theta}_{HT}. \end{aligned}$$

Then,

$$\sqrt{n}(\widehat{\theta}_H - \theta_0) = \sqrt{n}(\widehat{\theta}_{HT} - \theta_0) + o_p(1)(\widehat{\theta}_{HT} - \theta_0).$$

By Slutsky's theorem, if  $\sqrt{n}(\widehat{\theta}_{HT} - \theta_0) \xrightarrow{d} \mathcal{N}(0, V_0)$ , then  $\sqrt{n}(\widehat{\theta}_H - \theta_0) \xrightarrow{d} \mathcal{N}(0, V_0)$ . It is thus sufficient to prove asymptotic normality of the Horvitz-Thompson estimator, i.e.,

$$\sqrt{n}(\widehat{\theta}_{HT} - \theta_0) \xrightarrow{d} \mathcal{N}(0, V_0).$$

Under Assumptions 1 and 2, Proposition 1 ensures

$$\begin{aligned} \sqrt{n}(\widehat{\theta}_{HT} - \theta_0) &= \frac{1}{\sqrt{n}} \sum_{i=1}^n (D_i - \pi_i) \left( \frac{\varepsilon_i(1)}{\pi_i} + \frac{\varepsilon_i(0)}{1 - \pi_i} \right) + \sqrt{n} \left( \frac{1}{n} \sum_{i=1}^n \frac{D_i Z'_{1i}}{\pi_i} - \sum_{i=1}^n Z'_{1i} \right) \beta_1 \\ &\quad - \sqrt{n} \left( \frac{1}{n} \sum_{i=1}^n \frac{(1 - D_i) Z'_{0i}}{1 - \pi_i} - \sum_{i=1}^n Z'_{0i} \right) \beta_0 \\ &= \frac{1}{\sqrt{n}} \sum_{i=1}^n - \left( \frac{\varepsilon_i(1)}{\pi_i} + \frac{\varepsilon_i(0)}{1 - \pi_i} \right) \pi_i + \left( \frac{\varepsilon_i(1)}{\pi_i} + \frac{\varepsilon_i(0)}{1 - \pi_i} \right) D_i + o_p(1) \end{aligned}$$

Then, we have

$$\sqrt{n}(\widehat{\theta}_{HT} - \theta_0) = \frac{1}{\sqrt{n}} \sum_{i=1}^n f_i + g_i D_i + o_p(1)$$

with  $f_i := f(X_i, \varepsilon_i(0), \varepsilon_i(1)) = - \left( \frac{\varepsilon_i(1)}{\pi_i} + \frac{\varepsilon_i(0)}{1 - \pi_i} \right) \pi_i$  and  $g(X_i, \varepsilon_i(0), \varepsilon_i(1)) = \frac{\varepsilon_i(1)}{\pi_i} + \frac{\varepsilon_i(0)}{1 - \pi_i}$ . Slutsky's theorem ensures that we have to prove

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n f_i + g_i D_i \xrightarrow{d} \mathcal{N}(0, V_0).$$

By Assumption 4  $\mathbb{E}[f|X] = \mathbb{E}[g|X] = 0$ . Then, Conjecture 1 and Lemma C.2 give that, conditional on  $(X_i)_{i \geq 1}$ ,  $\frac{1}{\sqrt{n}} \sum_{i=1}^n f_i + g_i D_i \xrightarrow{d} \mathcal{N}(0, V_0)$ , with  $V_0 = \mathbb{E}[f_i^2 + (2g_i f_i + g_i^2) \pi_i] = \mathbb{E} \left[ \pi_i (1 - \pi_i) \left( \frac{\varepsilon_i(1)}{\pi_i} + \frac{\varepsilon_i(0)}{1 - \pi_i} \right)^2 \right]$ , in the sense of Definition C.1. Notice that  $V_0$  does not depend on  $(X_1)_{i \geq 1}$ , so convergence in distribution is unconditional. This concludes the proof.

### Comparison with Poisson randomization

If treated units are selected through Poisson randomization and Assumptions 1 and 2 hold, then  $(Y_i(1), Y_i(0), X_i, D_i)_{i=1, \dots, n}$  are i.i.d., and by the CLT,

$$\sqrt{n}(\widehat{\theta}_{HT} - \theta_0) \xrightarrow{d} \mathcal{N}(0, W_0)$$

with

$$W_0 = \mathbb{E} \left[ (D_i - \pi_i)^2 \left( \frac{Y_i(1)}{\pi_i} + \frac{Y_i(0)}{1 - \pi_i} \right)^2 \right]$$

$$\begin{aligned}
&= \mathbb{E} \left[ \pi_i(1 - \pi_i) \left( \frac{Y_i(1)}{\pi_i} + \frac{Y_i(0)}{1 - \pi_i} \right)^2 \right] \\
&= \mathbb{E} \left[ \pi_i(1 - \pi_i) \left( \frac{Z'_{1i}\beta_1}{\pi_i} + \frac{Z'_{0i}\beta_0}{1 - \pi_i} \right)^2 \right] + \mathbb{E} \left[ \pi_i(1 - \pi_i) \left( \frac{\varepsilon_i(1)}{\pi_i} + \frac{\varepsilon_i(0)}{1 - \pi_i} \right)^2 \right] \\
&= \mathbb{E} \left[ \pi_i(1 - \pi_i) \left( Z'_{0i} \frac{\beta_1 + \beta_0}{1 - \pi_i} \right)^2 \right] + V_0 = \Sigma_0 + V_0.
\end{aligned}$$

### B.3 Proof of Asymptotic Normality for the PATE (Second part of Propositions 3)

As shown in Proof B.2, if  $\sqrt{n}(\widehat{\theta}_{HT} - \theta_0^*) \xrightarrow{d} \mathcal{N}(0, V_0^*)$ , we have  $\sqrt{n}(\widehat{\theta}_H - \theta_0^*) \xrightarrow{d} \mathcal{N}(0, V_0^*)$ , we thus restrict ourselves to proving asymptotic normality for the Horvitz-Thompson estimator, i.e.,

$$\sqrt{n}(\widehat{\theta}_{HT} - \theta_0^*) \xrightarrow{d} \mathcal{N}(0, V_0^*)$$

where  $V_0^* = \mathbb{V}(Z'_{1i}\beta_1 - Z'_{0i}\beta_0) + \mathbb{E} \left[ \frac{\varepsilon_i(1)^2}{\pi_i} \right] + \mathbb{E} \left[ \frac{\varepsilon_i(0)^2}{1 - \pi_i} \right]$ .

Let us consider  $f_i := f(X_i, \varepsilon_i(1), \varepsilon_i(0)) = -\frac{\varepsilon_i(0)}{1 - \pi_i}$ ,  $g_i := g(X_i, \varepsilon_i(1), \varepsilon_i(0)) = \frac{\varepsilon_i(1)}{\pi_i} + \frac{\varepsilon_i(0)}{1 - \pi_i}$ , and  $h_i := h(X_i) = (Z_{1i} - \mathbb{E}[Z_{1i}])'\beta_1 - (Z_{0i} - \mathbb{E}[Z_{0i}])'\beta_0$ .

Under Assumptions 1 and 2, Proposition 1 ensures

$$\begin{aligned}
\sqrt{n}(\widehat{\theta}_{HT} - \theta_0^*) &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \left( \frac{Z'_{1i}D_i}{\pi_i} - E[Z_{1i}]' \right) \beta_1 \\
&\quad - \frac{1}{\sqrt{n}} \sum_{i=1}^n \left( \frac{Z'_{0i}(1 - D_i)}{1 - \pi_i} - E[Z_{0i}]' \right) \beta_0 \\
&\quad + \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\varepsilon_i(1)D_i}{\pi_i} - \frac{\varepsilon_i(0)(1 - D_i)}{1 - \pi_i} \\
&= \frac{1}{\sqrt{n}} \sum_{i=1}^n f_i + g_i D_i + \frac{1}{\sqrt{n}} \sum_{i=1}^n h_i + o_p(1).
\end{aligned}$$

Slutsky's theorem ensures that we have to prove

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n f_i + g_i D_i + \frac{1}{\sqrt{n}} \sum_{i=1}^n h_i \xrightarrow{d} \mathcal{N}(0, V_0^*). \quad (\text{B.1})$$

By Assumption 4,  $\mathbb{E}[f|X] = \mathbb{E}[g|X] = 0$ . Then, Conjecture 1 and Lemma C.2 give that, conditional on  $(X_i)_{i \geq 1}$ ,  $\sqrt{n}(\widehat{\theta}_{HT} - \theta_0^*) \xrightarrow{d} \mathcal{N}(0, V_{01}^*)$  with  $V_{01}^* = \mathbb{E}[f_i^2 + (2g_i f_i + g_i^2)\pi_i] = \mathbb{E} \left[ \frac{\varepsilon_i(1)^2}{\pi_i} \right] + \mathbb{E} \left[ \frac{\varepsilon_i(0)^2}{1 - \pi_i} \right]$ . Moreover, by the central limit theorem,  $\frac{1}{\sqrt{n}} \sum_{i=1}^n h_i \xrightarrow{d} \mathcal{N}(0, V_{02}^*)$  with  $V_{02}^* = \mathbb{V}(Z'_{1i}\beta_1 - Z'_{0i}\beta_0)$ . Theorem 2 in Chen and Rao (2007) ensures that  $\frac{1}{\sqrt{n}} \sum_{i=1}^n f_i + g_i D_i +$

$\frac{1}{\sqrt{n}} \sum_{i=1}^n h_i \xrightarrow{d} \mathcal{N}(0, V_{01}^* + V_{02}^*)$ . This concludes the proof.

### Comparison with Poisson randomization

If treated units are selected through Poisson randomization and Assumptions 1 and 2 hold, then  $(Y_i(1), Y_i(0), X_i, D_i)_{i=1, \dots, n}$  are i.i.d., and by the CLT,

$$\sqrt{n}(\widehat{\theta}_{HT} - \theta_0^*) \xrightarrow{d} \mathcal{N}(0, W_0^*)$$

with

$$\begin{aligned} W_0^* &= \mathbb{E} \left[ \left( \frac{Y_i(1)D_i}{\pi_i} - \frac{Y_i(0)(1-D_i)}{1-\pi_i} - \mathbb{E}[Y_i(1) - Y_i(0)] \right)^2 \right] \\ &= \mathbb{E} \left[ \frac{Y_i(1)^2}{\pi_i} \right] + \mathbb{E} \left[ \frac{Y_i(0)^2}{(1-\pi_i)} \right] - \mathbb{E}[Y_i(1) - Y_i(0)]^2 \\ &= \mathbb{E} \left[ \frac{(Z'_{1i}\beta_1)^2}{\pi_i} \right] + \mathbb{E} \left[ \frac{\varepsilon_i(1)^2}{\pi_i} \right] + \mathbb{E} \left[ \frac{(Z'_{0i}\beta_0)^2}{(1-\pi_i)} \right] + \mathbb{E} \left[ \frac{\varepsilon_i(0)^2}{1-\pi_i} \right] - \mathbb{E}[Z'_{1i}\beta_1 - Z'_{0i}\beta_0]^2 \\ &= \mathbb{E} \left[ \frac{(Z'_{1i}\beta_1)^2}{\pi_i} \right] + \mathbb{E} \left[ \frac{(Z'_{0i}\beta_0)^2}{(1-\pi_i)} \right] - \mathbb{E}[(Z'_{1i}\beta_1)^2] - \mathbb{E}[(Z'_{0i}\beta_0)^2] + 2\mathbb{E}[(Z'_{1i}\beta_1)(Z'_{0i}\beta_0)] + V_0^* \\ &= \mathbb{E} \left[ \frac{(Z'_{1i}\beta_1)^2(1-\pi_i)^2 + (Z'_{0i}\beta_0)^2\pi_i^2 + 2(Z'_{1i}\beta_1)(Z'_{0i}\beta_0)\pi_i(1-\pi_i)}{\pi_i(1-\pi_i)} \right] + V_0^* \\ &= \mathbb{E} \left[ \frac{\pi_i}{1-\pi_i} (Z'_{0i}(\beta_1 + \beta_0))^2 \right] + V_0^* = \Sigma_0 + V_0^*. \end{aligned}$$

## B.4 Proof of Randomized-based Inference (Proposition 4)

For completeness, we show first, as in Bai et al. (2022), that the strong null hypothesis (19)  $(Y_i(1), X_i) \stackrel{d}{=} (Y_i(0), X_i)$  is equivalent to stating  $Y_1, \dots, Y_n \perp\!\!\!\perp D_1, \dots, D_n | X_1, \dots, X_n$ .

Let us consider random allocations generated by the cube method  $d$  and  $d'$  in the support of  $D_1, \dots, D_n | X_1, \dots, X_n$  and any set  $A$ . Then we have,

$$\begin{aligned} &\mathbb{P}[(Y_1, \dots, Y_n) \in A | (D_1, \dots, D_n) = (d_1, \dots, d_n), X_1, \dots, X_n] \\ &= \mathbb{P}[(Y_1(d_1), \dots, Y_n(d_n)) \in A | (D_1, \dots, D_n) = (d_1, \dots, d_n), X_1, \dots, X_n] \\ &= \mathbb{P}[(Y_1(d_1), \dots, Y_n(d_n)) \in A | (D_1, \dots, D_n) = (d_1, \dots, d_n), X_1, \dots, X_n] \\ &= \mathbb{P}[(Y_1(d_1), \dots, Y_n(d_n)) \in A | X_1, \dots, X_n] \\ &= \mathbb{P}[(Y_1(d'_1), \dots, Y_n(d'_n)) \in A | X_1, \dots, X_n] \\ &= \mathbb{P}[(Y_1(d'_1), \dots, Y_n(d'_n)) \in A | (D_1, \dots, D_n) = (d'_1, \dots, d'_n), X_1, \dots, X_n] \end{aligned}$$

$$= \mathbb{P}[(Y_1, \dots, Y_n) \in A | (D_1, \dots, D_n) = (d'_1, \dots, d'_n), X_1, \dots, X_n],$$

so both hypothesis are equivalent. Then, under Assumptions 1 and 2, and the strong null hypothesis (19),

$$(Y_i, D_i, X_i)_{i \geq 1} \stackrel{d}{=} (Y_i, D_i^{(g)}, X_i)_{i \geq 1}.$$

We thus have

$$\begin{aligned} \mathbb{E} \left[ \sum_{g \in G_n^B} \phi_n^{rand}(\mathbf{P}_n^{(g)}) \right] &= \sum_{g \in G_n^B} \mathbb{E} \left[ \mathbb{E} \left[ \phi_n^{rand}(\mathbf{P}_n^{(g)}) | X_1, \dots, X_n \right] \right] \\ &= \sum_{g \in G_n^B} \mathbb{E} \left[ \mathbb{E} \left[ \phi_n^{rand}(\mathbf{P}_n) | X_1, \dots, X_n \right] \right] \\ &= B \mathbb{E} \left[ \phi_n^{rand}(\mathbf{P}_n) \right] \end{aligned} \tag{B.2}$$

Moreover,  $c_n(\mathbf{P}_n, 1 - \alpha) = c_n(\mathbf{P}_n^{(g)}, 1 - \alpha)$  for any  $g \in G_n^B$  ensures by definition of  $c_n(\mathbf{P}_n, 1 - \alpha)$

$$\sum_{g \in G_n^B} \phi_n^{rand}(\mathbf{P}_n^{(g)}) \leq B\alpha \tag{B.3}$$

Combining equations (B.2) and (B.3) we get  $\mathbb{E}[\phi_n^{rand}(\mathbf{P}_n)] \leq \alpha$ , which concludes the proof.

## C Lemmas for the Cube Method

### Lemma C.1 (Exchangeability)

For any permutation  $\sigma$  of  $\{1, \dots, n\}$  we have:

$$(D_{\sigma(i)}, \pi_{\sigma(i)}^* X_{\sigma(i)})_{i=1, \dots, n} \stackrel{d}{=} (D_i, X_i)_{i=1, \dots, n}$$

Proof:

For any value of  $n$ , the Cube algorithm ensures there exists a finite collection of independent uniform random variables  $(U_1, \dots, U_K)$  independent of  $(X_1, \dots, X_n)$  such that  $(D_1, \dots, D_n, \pi_1^*, \dots, \pi_n^*) = f(X_1, \dots, X_n, U_1, \dots, U_K)$ . Because the  $X$  are iid and independent of the  $U$ , we have:

$$(X_{\sigma(1)}, \dots, X_{\sigma(n)}, U_1, \dots, U_K) \stackrel{d}{=} (X_1, \dots, X_n, U_1, \dots, U_K).$$

The result follows.

### Definition C.1

$W_n \xrightarrow{d} \mathcal{N}(0, \sigma^2)$  conditional on  $(X_i)_{i \geq 1}$  if and only if for any  $h$  bounded Lipschitz  $\mathbb{E}(h(W_n) | (X_i)_{i \geq 1})$  converges almost surely to  $\int h(u) \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{u^2}{2\sigma^2}\right) du$ .

Usual criteria (e.g, Portmanteau's lemma or Lévy's continuity theorem) to prove convergence in distribution could be adapted to prove the convergence in distribution conditional on  $(X_i)_{i \geq 1}$  apply if the usual expectations and probabilities are replaced by conditional expectations and probabilities and usual convergence of sequences is replaced by almost sure convergence of random variables. More concretely, we will use the fact that if for any  $k \geq 1$ ,  $\mathbb{E}((W_n)^k | (X_i)_{i \geq 1})$  converges almost surely to the  $k$ th-raw moment of a Gaussian distribution of variance  $\sigma^2$  then  $W_n \xrightarrow{d} \mathcal{N}(0, \sigma^2)$  conditional on  $(X_i)_{i \geq 1}$ . This is an adaptation of the theorem of Takacs (1991) that states that if for any  $k \geq 1$ ,  $\mathbb{E}((W_n)^k)$  converges to the  $k$ th-raw moment of a Gaussian distribution of variance  $\sigma^2$  then  $W_n \xrightarrow{d} \mathcal{N}(0, \sigma^2)$ . Moreover,  $W_n \xrightarrow{d} \mathcal{N}(0, \sigma^2)$  conditional on  $(X_i)_{i \geq 1}$  if and only if  $\forall t \in \mathbb{R}$ ,  $\mathbb{P}(W_n \leq t | (X_i)_{i \geq 1})$  converges almost surely to  $\Phi\left(\frac{t}{\sqrt{\sigma^2}}\right)$  for  $\Phi$  the c.d.f. of the standard Gaussian.

### Lemma C.2 (Asymptotic normality)

Let  $f$  and  $g$  be two functions such that for  $f_i = f(\delta_i(1), \delta_i(0), X_i)$  and  $g_i = g(\delta_i(1), \delta_i(0), X_i)$  we have  $\mathbb{E}(f_i^2 + g_i^2) < \infty$  and  $\mathbb{E}[f_i | X_i] = \mathbb{E}[g_i | X_i] = 0$ .

If Assumptions 1 and 2 and Conjecture 1 hold. Then, conditional on  $(X_i)_{i \geq 1}$ ,

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n f_i + g_i D_i \xrightarrow{d} \mathcal{N}(0, V_0) \tag{C.1}$$

with  $V_0 = \mathbb{E}[f_1^2 + (2g_1 f_1 + g_1^2) \pi_1]$ .

Proof:

**First step:**  $|f_i| + |g_i|$  **bounded implies** (C.1)

Let us assume it exists  $K > 0$  such that  $|f_1| + |g_1| < K$  for any  $k \in \mathbb{N}$ . This ensures that all the moments of  $f_i + g_i D_i$  exist. Let  $M_{n,k} = \mathbb{E} \left[ \left( \frac{1}{\sqrt{n}} \sum_{i=1}^n f_i + g_i D_i \right)^k \middle| (X_i)_{i \geq 1} \right]$ .

We have

$$M_{n,k} = \mathbb{E} \left[ n^{-k/2} \sum_{1 \leq i_1, \dots, i_k \leq n} \prod_{\ell=1}^k (f_{i_\ell} + g_{i_\ell} D_{i_\ell}) \middle| (X_i)_{i \geq 1} \right].$$

Let us order the indices  $i_1, \dots, i_k$  as  $j_1, \dots, j_m$  for some  $1 \leq m \leq k$  with each  $j_\ell$  occurring with multiplicity  $a_\ell$ . Let  $A_{k,m} := \{a = (a_1, \dots, a_m) \in \mathbb{N}^{*m} : \sum_{\ell=1}^m a_\ell = k\}$  and for  $a \in A_{k,m}$ ,  $c_{k,a} = \frac{k!}{\prod_{\ell=1}^m a_\ell!}$ . We have:

$$M_{n,k} = \sum_{m=1}^k n^{-k/2} \sum_{1 \leq j_1 < \dots < j_m \leq n} \sum_{a \in A_{k,m}} c_{k,a} \mathbb{E} \left[ \prod_{\ell=1}^k (f_{j_\ell} + g_{j_\ell} D_{j_\ell})^{a_\ell} \middle| (X_i)_{i \geq 1} \right].$$

In order to prove the convergence of moments, we will focus on the summands

$$B_{n,k,m} = n^{-k/2} \sum_{1 \leq j_1 < \dots < j_m \leq n} \sum_{a \in A_{k,m}} c_{k,a} \mathbb{E} \left[ \prod_{\ell=1}^k (f_{j_\ell} + g_{j_\ell} D_{j_\ell})^{a_\ell} \middle| (X_i)_{i \geq 1} \right].$$

Notice that  $|B_{n,k,m}| \leq n^{-k/2} \binom{n}{m} \sum_{a \in A_{k,m}} c_{k,a} K^k = O(n^{m-k/2})$ . For  $m < k/2$ , we thus have  $\lim_n B_{n,k,m} = 0$ .

We focus now in the case  $m > k/2$ . For  $\mathcal{K} \subseteq \{1, \dots, m\}$ , we note  $\mathcal{K}^c = \{1, \dots, m\} \setminus \mathcal{K}$ . Then, the binomial theorem and the expansion  $\prod_{\ell=1}^m (x_\ell + y_\ell) = \sum_{\mathcal{K} \subseteq \{1, \dots, m\}} \prod_{\ell \in \mathcal{K}} x_\ell \prod_{\ell' \in \mathcal{K}^c} y_{\ell'}$  and identity  $D^a = D$  for  $a \geq 1$  ensure

$$\begin{aligned} & B_{n,k,m} \\ &= n^{-k/2} \sum_{1 \leq j_1 < \dots < j_m \leq n} \sum_{a \in A_{k,m}} c_{k,a} \mathbb{E} \left[ \prod_{\ell=1}^k (f_{j_\ell} + g_{j_\ell} D_{j_\ell})^{a_\ell} \middle| (X_i)_{i \geq 1} \right] \\ &= n^{-k/2} \sum_{1 \leq j_1 < \dots < j_m \leq n} \sum_{a \in A_{k,m}} c_{k,a} \mathbb{E} \left[ \prod_{\ell=1}^m \left[ f_{j_\ell}^{a_\ell} + \left( \sum_{r=1}^{a_\ell} \binom{a_\ell}{r} f_{j_\ell}^{a_\ell-r} g_{j_\ell}^r \right) D_{j_\ell} \right] \middle| (X_i)_{i \geq 1} \right] \\ &= n^{-k/2} \sum_{1 \leq j_1 < \dots < j_m \leq n} \sum_{a \in A_{k,m}} c_{k,a} \sum_{\mathcal{K} \subseteq \{1, \dots, m\}} \mathbb{E} \left[ \prod_{\ell \in \mathcal{K}} f_{j_\ell}^{a_\ell} \prod_{\ell' \in \mathcal{K}^c} \left( \sum_{r=1}^{a_{\ell'}} \binom{a_{\ell'}}{r} f_{j_{\ell'}}^{a_{\ell'}-r} g_{j_{\ell'}}^r \right) \prod_{\ell' \in \mathcal{K}^c} D_{j_{\ell'}} \middle| (X_i)_{i \geq 1} \right] \end{aligned}$$

Then, independence of  $(f_i, g_i)_{i \geq 1}$  across  $i$  and conditional independence  $(f_i, g_i) \perp\!\!\!\perp D_i \middle| (X_{i'})_{i' \geq 1}$  ensure

$$B_{n,k,m} = n^{-k/2} \sum_{1 \leq j_1 < \dots < j_m \leq n} \sum_{a \in A_{k,m}} c_{k,a} \sum_{\mathcal{K} \subseteq \{1, \dots, m\}} \prod_{\ell \in \mathcal{K}} \mathbb{E} [f_{j_\ell}^{a_\ell} | X_{j_\ell}] \prod_{\ell' \in \mathcal{K}^c} \left( \sum_{r=1}^{a_{\ell'}} \binom{a_{\ell'}}{r} \mathbb{E} [f_{j_{\ell'}}^{a_{\ell'}-r} g_{j_{\ell'}}^r | X_{j_{\ell'}}] \right)$$

$$\mathbb{E} \left[ \prod_{\ell'' \in \mathcal{K}^c} D_{j_{\ell''}} \mid (X_i)_{i \geq 1} \right].$$

Because  $m > k/2$ , for any  $a \in A_{k,m}$  there exists  $s$  such that  $a_s = 1$ . For any  $\mathcal{K}$ , if  $s \in \mathcal{K}$ , then  $\prod_{\ell \in \mathcal{K}} \mathbb{E} [f_{j_\ell}^{a_\ell} | X_{j_\ell}] = \mathbb{E} [f_{j_s} | X_{j_s}] \prod_{\ell \in \mathcal{K} \setminus \{s\}} \mathbb{E} [f_{j_\ell}^{a_\ell} | X_{j_\ell}] = 0$ , else  $s \in \mathcal{K}^c$  and  $\prod_{\ell' \in \mathcal{K}^c} \left( \sum_{r=1}^{a_{\ell'}} \binom{a_{\ell'}}{r} f_{j_{\ell'}}^{a_{\ell'}-r} g_{j_{\ell'}}^r \right) = \mathbb{E} (g_{j_s} | X_{j_s}) \prod_{\ell' \in \mathcal{K}^c \setminus \{s\}} \left( \sum_{r=1}^{a_{\ell'}} \binom{a_{\ell'}}{r} f_{j_{\ell'}}^{a_{\ell'}-r} g_{j_{\ell'}}^r \right) = 0$ . It follows that if  $m > k/2$  we have  $B_{n,k,m} = 0$ .

Let now consider the last case  $m = k/2$ . For  $a \in A_{k,k/2}$  either there exists  $s$  such that  $a_s = 1$  and by the previous reasoning, we have  $\prod_{\ell \in \mathcal{K}} \mathbb{E} [f_{j_\ell}^{a_\ell} | X_{j_\ell}] \prod_{\ell' \in \mathcal{K}^c} \left( \sum_{r=1}^{a_{\ell'}} \binom{a_{\ell'}}{r} \mathbb{E} [f_{j_{\ell'}}^{a_{\ell'}-r} g_{j_{\ell'}}^r | X_{j_{\ell'}}] \right) = 0$  for any  $\mathcal{K}$ , either  $a = (2, \dots, 2)$  and it follows

$$B_{n,k,k/2} = n^{-k/2} \sum_{1 \leq j_1 < \dots < j_{k/2} \leq n} \frac{k!}{2^{k/2}} \sum_{\mathcal{K} \subseteq \{1, \dots, k/2\}} \prod_{\ell \in \mathcal{K}} \mathbb{E} [f_{j_\ell}^2 | X_{j_\ell}] \prod_{\ell' \in \mathcal{K}^c} \mathbb{E} [2f_{j_{\ell'}} g_{j_{\ell'}} + g_{j_{\ell'}}^2 | X_{j_{\ell'}}] \mathbb{E} \left[ \prod_{\ell'' \in \mathcal{K}^c} D_{j_{\ell''}} \mid (X_i)_{i \geq 1} \right]$$

Conjecture 1 and the fact that  $\max(|f_i|^2, |2f_i g_i + g_i^2|) \leq 3K^2$  ensure

$$\begin{aligned} B_{n,k,k/2} &= n^{-k/2} \sum_{1 \leq j_1 < \dots < j_{k/2} \leq n} \frac{k!}{2^{k/2}} \sum_{\mathcal{K} \subseteq \{1, \dots, k/2\}} \prod_{\ell \in \mathcal{K}} \mathbb{E} [f_{j_\ell}^2 | X_{j_\ell}] \prod_{\ell' \in \mathcal{K}^c} \mathbb{E} [2f_{j_{\ell'}} g_{j_{\ell'}} + g_{j_{\ell'}}^2 | X_{j_{\ell'}}] \prod_{\ell'' \in \mathcal{K}^c} \pi_{j_{\ell''}} \\ &\quad + n^{-k/2} \binom{n}{k/2} \frac{k!}{2^{k/2}} 2^{k/2} (3K^2)^{k/2} o(1) \\ &= n^{-k/2} \sum_{1 \leq j_1 < \dots < j_{k/2} \leq n} \frac{k!}{2^{k/2}} \sum_{\mathcal{K} \subseteq \{1, \dots, k/2\}} \prod_{\ell \in \mathcal{K}} \mathbb{E} [f_{j_\ell}^2 | X_{j_\ell}] \prod_{\ell' \in \mathcal{K}^c} \mathbb{E} [(2f_{j_{\ell'}} g_{j_{\ell'}} + g_{j_{\ell'}}^2) \pi_{j_{\ell'}} | X_{j_{\ell'}}] \\ &\quad + o(1) \end{aligned}$$

Factorization formula  $\sum_{\mathcal{K} \subseteq \{1, \dots, m\}} \prod_{\ell \in \mathcal{K}} x_\ell \prod_{\ell' \in \mathcal{K}^c} y_{\ell'} = \prod_{\ell=1}^m (x_\ell + y_\ell)$  ensures

$$\begin{aligned} B_{n,k,k/2} &= \frac{k!}{2^{k/2}} n^{-k/2} \sum_{1 \leq j_1 < \dots < j_{k/2} \leq n} \prod_{\ell=1}^{k/2} \mathbb{E} [f_{j_\ell}^2 + (2f_{j_\ell} g_{j_\ell} + g_{j_\ell}^2) \pi_{j_\ell} | X_{j_\ell}] + o(1) \\ &= \frac{k!}{2^{k/2}} n^{-k/2} \binom{n}{k/2} \binom{n}{k/2}^{-1} \sum_{1 \leq j_1 < \dots < j_{k/2} \leq n} h(X_{j_1}, \dots, X_{j_{k/2}}) + o(1) \end{aligned}$$

for  $h(u_1, \dots, u_{k/2}) = \prod_{i=1}^{k/2} \mathbb{E} (f^2 + (2fg + g^2)\pi | X = u_i)$ . Strong law of large numbers for U-statistics (Aaronson et al., 1996) ensures that  $\binom{n}{k/2}^{-1} \sum_{1 \leq j_1 < \dots < j_{k/2} \leq n} h(X_{j_1}, \dots, X_{j_{k/2}})$  converges almost surely to  $\mathbb{E}(h(X_1, \dots, X_{k/2})) = (V_0)^{k/2}$  and  $\lim_n n^{-k/2} \binom{n}{k/2} = \frac{1}{(k/2)!}$ . Then,  $\lim_n M_{n,k} = 0$  for  $k$  odd, and  $\lim_n M_{n,k} = \frac{k!}{2^{k/2}(k/2)!} V_0^{k/2}$  for  $k$  even. By the adapted form of the theorem in Takacs (1991), if  $f_i$  and  $g_i$ , are bounded, we have that, conditional on  $(X_i)_{i \geq 1}$   $\frac{1}{\sqrt{n}} \sum_{i=1}^n f_i + g_i D_i$  converges almost surely to a Gaussian of variance  $V_0$ .



**Second step:**  $\mathbb{E}(Y(0)^2 + Y(1)^2 + \|X\|^2) < \infty$  **implies** (C.1)

Assumption 1 ensures only that  $f_i$  and  $g_i$  admit moments of order 2. Then, for  $M > 0$ , let  $f_{\leq M,i}$ ,  $f_{>M,i}$ ,  $g_{\leq M,i}$  and  $g_{>M,i}$  the truncated variables  $f_{\leq M,i} = f_i \mathbb{1}\{|f_i| \leq M\}$ ,  $f_{>M,i} = f_i \mathbb{1}\{|f_i| > M\}$ ,  $g_{\leq M,i} = g_i \mathbb{1}\{|g_i| \leq M\}$  and  $g_{>M,i} = g_i \mathbb{1}\{|g_i| > M\}$ . We define  $\tilde{f}_{\leq M,i} = f_{\leq M,i} - \mathbb{E}[f_{\leq M,i}|X_i]$ ,  $\tilde{f}_{>M,i} = f_{>M,i} - \mathbb{E}[f_{>M,i}|X_i]$ ,  $\tilde{g}_{\leq M,i} = g_{\leq M,i} - \mathbb{E}[g_{\leq M,i}|X_i]$  and  $\tilde{g}_{>M,i} = g_{>M,i} - \mathbb{E}[g_{>M,i}|X_i]$ . We have:

$$\begin{aligned}
& \mathbb{E} \left[ \left| \frac{1}{\sqrt{n}} \sum_{i=1}^n (\tilde{f}_{>M,i} + \tilde{g}_{>M,i} D_i) \right|^2 \middle| (X_\ell)_{\ell \geq 1} \right] \\
&= \frac{1}{n} \sum_{i=1}^n \mathbb{E} [(\tilde{f}_{>M,i} + \tilde{g}_{>M,i} D_i)^2 | (X_\ell)_{\ell \geq 1}] + \frac{1}{n} \sum_{\substack{1 \leq i, j \leq n \\ i \neq j}} \mathbb{E} [(\tilde{f}_{>M,i} + \tilde{g}_{>M,i} D_i) (\tilde{f}_{>M,j} + \tilde{g}_{>M,j} D_j) | (X_\ell)_{\ell \geq 1}] \\
&= \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\tilde{f}_{>M,i}^2 | (X_\ell)_{\ell \geq 1}] + \mathbb{E} [(2\tilde{f}_{>M,i} \tilde{g}_{>M,i} + \tilde{g}_{>M,i}^2) | (X_\ell)_{\ell \geq 1}] \mathbb{E} [D_i | (X_\ell)_{\ell \geq 1}] \\
&+ \frac{1}{n} \sum_{\substack{1 \leq i, j \leq n \\ i \neq j}} \left( \mathbb{E} [\tilde{f}_{>M,i} \tilde{f}_{>M,j} | (X_\ell)_{\ell \geq 1}] + \mathbb{E} [\tilde{f}_{>M,i} \tilde{g}_{>M,j} | (X_\ell)_{\ell \geq 1}] \mathbb{E} [D_j | (X_\ell)_{\ell \geq 1}] \right. \\
&\quad \left. + \mathbb{E} [\tilde{f}_{>M,j} \tilde{g}_{>M,i} | (X_\ell)_{\ell \geq 1}] \mathbb{E} [D_i | (X_\ell)_{\ell \geq 1}] + \mathbb{E} [\tilde{g}_{>M,i} \tilde{g}_{>M,j} | (X_\ell)_{\ell \geq 1}] \mathbb{E} [D_i D_j | (X_\ell)_{\ell \geq 1}] \right) \\
&= \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\tilde{f}_{>M,i}^2 | X_i] + \mathbb{E} [(2\tilde{f}_{>M,i} \tilde{g}_{>M,i} + \tilde{g}_{>M,i}^2) | X_i] \pi_i \\
&+ \frac{1}{n} \sum_{\substack{1 \leq i, j \leq n \\ i \neq j}} \left( \mathbb{E} [\tilde{f}_{>M,i} | X_i] \mathbb{E} [\tilde{f}_{>M,j} | X_j] + \mathbb{E} [\tilde{f}_{>M,i} | X_i] \mathbb{E} [\tilde{g}_{>M,j} | X_j] \pi_j \right. \\
&\quad \left. + \mathbb{E} [\tilde{f}_{>M,j} | X_j] \mathbb{E} [\tilde{g}_{>M,i} | X_i] \pi_i + \mathbb{E} [\tilde{g}_{>M,i} | X_i] \mathbb{E} [\tilde{g}_{>M,j} | X_j] \mathbb{E} [D_i D_j | (X_\ell)_{\ell \geq 1}] \right) \\
&= \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\tilde{f}_{>M,i}^2 + (2\tilde{f}_{>M,i} \tilde{g}_{>M,i} + \tilde{g}_{>M,i}^2) \pi_i | X_i]
\end{aligned}$$

The second equality holds because  $(f_1, \dots, f_n, g_1, \dots, g_n) \perp\!\!\!\perp (D_1, \dots, D_n) | X_1, \dots, X_n$  by Assumption 2. The third equality holds because  $(f_i, g_i, X_i)_{i \geq 1}$  are independent across  $i$  by Assumption 1 and  $\mathbb{E}[D_i | (X_\ell)_{\ell \geq 1}] = \pi_i$  by Assumption 2. The fourth equality holds because  $\mathbb{E}[\tilde{f}_{>M,\ell} | X_\ell] = \mathbb{E}[\tilde{g}_{>M,\ell} | X_\ell] = 0$ .

The SLLN ensures that  $\frac{1}{n} \sum_{i=1}^n \mathbb{E} [\tilde{f}_{>M,i}^2 + (2\tilde{f}_{>M,i} \tilde{g}_{>M,i} + \tilde{g}_{>M,i}^2) \pi_i | X_i]$  converges almost-surely to  $\mathbb{E} [\tilde{f}_{>M,1}^2 + (2\tilde{f}_{>M,1} \tilde{g}_{>M,1} + \tilde{g}_{>M,1}^2) \pi_1]$ . It follows that by Cauchy-Schwarz inequality:

$$\limsup_n \mathbb{E} \left[ \left| \frac{1}{\sqrt{n}} \sum_{i=1}^n (\tilde{f}_{>M,i} + \tilde{g}_{>M,i} D_i) \right|^2 \middle| (X_\ell)_{\ell \geq 1} \right]$$

$$\begin{aligned}
&\leq \limsup_n \mathbb{E} \left[ \left| \frac{1}{\sqrt{n}} \sum_{i=1}^n (\tilde{f}_{>M,i} + \tilde{g}_{>M,i} D_i) \right|^2 \middle| (X_\ell)_{\ell \geq 1} \right]^{1/2} \\
&= \mathbb{E} [\tilde{f}_{>M,1}^2 + (2\tilde{f}_{>M,1}\tilde{g}_{>M,1} + \tilde{g}_{>M,1}^2) \pi_1]^{1/2} \tag{C.2}
\end{aligned}$$

which, by dominated convergence, is arbitrarily small for a sufficiently large  $M$ .

Let  $h$  a bounded Lipschitz function of constant  $c_h$ ,  $V(M) = \mathbb{E} [\tilde{f}_{\leq M,1}^2 + (2\tilde{f}_{>M,1}\tilde{g}_{\leq M,1} + \tilde{g}_{\leq M,1}^2) \pi_1]$ , and  $N \sim \mathcal{N}(0, 1)$ . We have by triangle and Lipschitz inequities, and the fact that  $f_i + g_i D_i = \tilde{f}_{\leq M,i} + \tilde{g}_{\leq M,i} D_i + \tilde{f}_{>M,i} + \tilde{g}_{>M,i} D_i$ :

$$\begin{aligned}
&\left| \mathbb{E} \left[ h \left( \frac{1}{\sqrt{n}} \sum_{i=1}^n f_i + g_i D_i \right) \middle| (X_\ell)_{\ell \geq 1} \right] - \mathbb{E} [h(V_0^{1/2} N)] \right| \\
&\leq \left| \mathbb{E} \left[ h \left( \frac{1}{\sqrt{n}} \sum_{i=1}^n f_i + g_i D_i \right) \middle| (X_\ell)_{\ell \geq 1} \right] - \mathbb{E} \left[ h \left( \frac{1}{\sqrt{n}} \sum_{i=1}^n \tilde{f}_{\leq M,i} + \tilde{g}_{\leq M,i} D_i \right) \middle| (X_\ell)_{\ell \geq 1} \right] \right| \\
&+ \left| \mathbb{E} \left[ h \left( \frac{1}{\sqrt{n}} \sum_{i=1}^n \tilde{f}_{\leq M,i} + \tilde{g}_{\leq M,i} D_i \right) \middle| (X_\ell)_{\ell \geq 1} \right] - \mathbb{E} [h(V(M)^{1/2} N)] \right| \\
&+ \left| \mathbb{E} [h(V(M)^{1/2} N)] - \mathbb{E} [h(V_0^{1/2} N)] \right| \\
&\leq c_h \mathbb{E} \left[ \left| \frac{1}{\sqrt{n}} \sum_{i=1}^n \tilde{f}_{>M,i} + \tilde{g}_{>M,i} D_i \right|^2 \middle| (X_\ell)_{\ell \geq 1} \right] \\
&+ \left| \mathbb{E} \left[ h \left( \frac{1}{\sqrt{n}} \sum_{i=1}^n \tilde{f}_{\leq M,i} + \tilde{g}_{\leq M,i} D_i \right) \middle| (X_\ell)_{\ell \geq 1} \right] - \mathbb{E} [h(V(M)^{1/2} N)] \right| \\
&+ c_h |V(M)^{1/2} - V_0^{1/2}| \mathbb{E}(|N|).
\end{aligned}$$

The first step of the proof and (C.2) ensure that for any value of  $M > 0$ :

$$\begin{aligned}
&\limsup_n \left| \mathbb{E} \left[ h \left( \frac{1}{\sqrt{n}} \sum_{i=1}^n f_i + g_i D_i \right) \middle| (X_\ell)_{\ell \geq 1} \right] - \mathbb{E} [h(V_0^{1/2} N)] \right| \\
&\leq c_h \left( \mathbb{E} [\tilde{f}_{>M,1}^2 + (2\tilde{f}_{>M,1}\tilde{g}_{>M,1} + \tilde{g}_{>M,1}^2) \pi_1]^{1/2} + |V(M)^{1/2} - V_0^{1/2}| \right).
\end{aligned}$$

By dominated convergence,  $\lim_M V(M) = V_0$  and  $\lim_M \mathbb{E} [\tilde{f}_{>M,1}^2 + (2\tilde{f}_{>M,1}\tilde{g}_{>M,1} + \tilde{g}_{>M,1}^2) \pi_1] = 0$ . Next, considering  $M$  tending to  $\infty$ , dominated convergence ensures

$$\limsup_n \left| \mathbb{E} \left[ h \left( \frac{1}{\sqrt{n}} \sum_{i=1}^n f_i + g_i D_i \right) \middle| (X_\ell)_{\ell \geq 1} \right] - \mathbb{E} [h(V_0^{1/2} N)] \right| = 0.$$

This achieves the proof.

## D Additional Tables

Table D.1: Empirical application: Standard deviation of ATE estimators

	Number of covariates	Complete Randomization (1)	Stratified Randomization (2)	Matched Pairs (3)	Cube Method (4)	Double Lasso (5)
$n = 100$	1	2.958	2.662	2.363	2.382	2.371
	2	–	2.673	2.357	2.438	–
	3	–	2.658	2.390	2.448	–
	5	–	2.816	2.427	2.428	–
	7	–	3.153	2.459	2.434	–
	9	–	3.745	2.530	2.389	–
	12	–	6.770	2.563	2.391	–
$n = 256$	1	1.825	1.675	1.455	1.478	1.457
	2	–	1.664	1.456	1.491	–
	3	–	1.673	1.459	1.485	–
	5	–	1.673	1.506	1.483	–
	7	–	1.771	1.532	1.492	–
	9	–	2.004	1.545	1.468	–
	12	–	2.818	1.593	1.505	–
$n = 500$	1	1.299	1.189	1.052	1.055	1.031
	2	–	1.174	1.052	1.048	–
	3	–	1.173	1.040	1.060	–
	5	–	1.188	1.050	1.057	–
	7	–	1.227	1.055	1.059	–
	9	–	1.317	1.085	1.063	–
	12	–	1.691	1.107	1.055	–
$n = 1000$	1	0.918	0.836	0.738	0.751	0.736
	2	–	0.835	0.735	0.739	–
	3	–	0.837	0.736	0.734	–
	5	–	0.832	0.741	0.750	–
	7	–	0.844	0.745	0.748	–
	9	–	0.898	0.761	0.748	–
	12	–	1.059	0.778	0.744	–

This table shows the standard deviation of PATE estimators for different allocation designs and experimental sample sizes. For each allocation design, the standard deviation estimates are computed over 10,000 simulations. For columns 1, 3, and 4, the estimator used is the Horvitz-Thompson algorithm. In column 1 the design used is complete randomization. For column 3, we assign treatment using a matched pairs design, pairing individuals to the closest unit and using the Mahalanobis distance whenever more than one covariate is balanced. Column 4 shows the results for the cube method with two moments for each variable. For column 2, we run stratified randomization. We use median values for continuous variables and we estimate the PATE using OLS with strata fixed-effects in accordance with Bugni et al. (2018). Finally, column 5 uses a double-lasso selection procedure as described in Belloni et al. (2014).

Table D.2: Empirical application: Average bias of ATE estimators

	Number of covariates	Complete Randomization (1)	Stratified Randomization (2)	Matched Pairs (3)	Cube Method (4)	Double Lasso (5)
$n = 100$	1	0.046	-0.065	-0.007	0.032	0.256
	2	-	-0.025	-0.028	0.023	-
	3	-	-0.047	-0.058	-0.031	-
	5	-	-0.040	-0.064	-0.051	-
	7	-	-0.010	-0.011	-0.003	-
	9	-	0.115	-0.036	0.008	-
	12	-	0.359	-0.008	-0.043	-
$n = 256$	1	-0.016	-0.028	-0.006	-0.038	0.112
	2	-	-0.010	-0.019	-0.027	-
	3	-	-0.014	-0.019	-0.019	-
	5	-	-0.032	-0.022	-0.004	-
	7	-	-0.015	-0.010	-0.025	-
	9	-	0.048	-0.033	-0.030	-
	12	-	0.222	-0.030	-0.018	-
$n = 500$	1	-0.010	-0.005	0.009	0.007	0.088
	2	-	0.011	-0.027	0.000	-
	3	-	-0.013	0.005	-0.011	-
	5	-	0.028	-0.002	0.006	-
	7	-	0.003	-0.007	0.002	-
	9	-	0.074	-0.019	-0.005	-
	12	-	0.146	-0.024	-0.015	-
$n = 1000$	1	-0.015	-0.008	-0.013	-0.002	0.056
	2	-	-0.000	0.002	-0.006	-
	3	-	-0.002	-0.008	-0.002	-
	5	-	0.003	-0.006	0.002	-
	7	-	0.004	-0.009	0.004	-
	9	-	0.021	-0.010	-0.003	-
	12	-	0.064	-0.004	-0.008	-

This table shows the average bias of PATE estimators for different allocation designs and experimental sample sizes. For each allocation design, the average bias estimates are computed over 10,000 simulations. For columns 1, 3, and 4, the estimator used is the Horvitz-Thompson algorithm. In column 1, the design used is complete randomization. For column 3, we assign treatment using a matched pairs design, pairing individuals to the closest unit and using the Mahalanobis distance whenever more than one covariate is balanced. Column 4 shows the results for the cube method with two moments for each variable. For column 2, we run stratified randomization. We use median values for continuous variables and we estimate the PATE using OLS with strata fixed-effects in accordance with Bugni et al. (2018). Finally, column 5 uses a double-lasso selection procedure as described in Belloni et al. (2014).

Table D.3: Empirical application: Coverage Rates

	Number of covariates	Complete Randomization (1)	Stratified Randomization (2)	Matched Pairs (3)	Cube Method (4)	Double Lasso (5)
$n = 100$	1	0.945	0.943	0.948	0.940	0.941
	2	–	0.937	0.949	0.932	–
	3	–	0.937	0.950	0.931	–
	5	–	0.910	0.953	0.927	–
	7	–	0.884	0.956	0.924	–
	9	–	0.865	0.956	0.928	–
	12	–	0.776	0.955	0.927	–
$n = 256$	1	0.951	0.946	0.951	0.944	0.946
	2	–	0.947	0.953	0.939	–
	3	–	0.943	0.954	0.943	–
	5	–	0.938	0.952	0.943	–
	7	–	0.918	0.952	0.939	–
	9	–	0.902	0.955	0.941	–
	12	–	0.858	0.955	0.930	–
$n = 500$	1	0.952	0.951	0.949	0.950	0.952
	2	–	0.950	0.948	0.949	–
	3	–	0.951	0.954	0.950	–
	5	–	0.943	0.955	0.945	–
	7	–	0.931	0.957	0.946	–
	9	–	0.919	0.957	0.941	–
	12	–	0.883	0.958	0.944	–
$n = 1000$	1	0.955	0.952	0.953	0.945	0.949
	2	–	0.953	0.953	0.950	–
	3	–	0.952	0.954	0.947	–
	5	–	0.947	0.954	0.947	–
	7	–	0.939	0.957	0.948	–
	9	–	0.927	0.959	0.946	–
	12	–	0.902	0.955	0.946	–

This table shows the coverage rate of 95%-confidence intervals of PATE estimators for different allocation designs and experimental sample sizes. For each allocation design, the coverage rate estimates are computed over 10,000 simulations. For columns 1, 3, and 4, the estimator used is the Horvitz-Thompson algorithm. In column 1, the design used is complete randomization and the confidence intervals are constructed using White standard errors. For column 3, we assign treatment using a matched pairs design, pairing individuals to the closest unit and using the Mahalanobis distance whenever more than one covariate is balanced. Confidence intervals are constructed as described in (Bai, 2022). Column 4 shows the results for the cube method with two moments for each variable. Confidence intervals follow the asymptotic-based procedure described in Section 4.3. For column 2, we run stratified randomization. We use median values for continuous variables and we estimate the PATE using OLS with strata fixed-effects in accordance with Bugni et al. (2018). Finally, column 5 uses a double-lasso selection procedure as described in Belloni et al. (2014) with White standard errors for the post-lasso estimation.

Table D.4: Empirical application: Test Power

	Number of covariates	Complete Randomization (1)	Stratified Randomization (2)	Matched Pairs (3)	Cube Method (4)	Double Lasso (5)
$n = 100$	1	0.151	0.191	0.223	0.226	0.192
	2	–	0.194	0.215	0.234	–
	3	–	0.207	0.209	0.249	–
	5	–	0.233	0.197	0.256	–
	7	–	0.235	0.181	0.257	–
	9	–	0.220	0.180	0.255	–
	12	–	0.251	0.171	0.262	–
$n = 256$	1	0.325	0.385	0.465	0.477	0.440
	2	–	0.389	0.464	0.475	–
	3	–	0.393	0.457	0.476	–
	5	–	0.416	0.448	0.478	–
	7	–	0.429	0.423	0.486	–
	9	–	0.382	0.403	0.496	–
	12	–	0.288	0.394	0.499	–
$n = 500$	1	0.552	0.638	0.740	0.739	0.719
	2	–	0.637	0.749	0.743	–
	3	–	0.642	0.735	0.7491	–
	5	–	0.640	0.725	0.743	–
	7	–	0.663	0.714	0.750	–
	9	–	0.605	0.692	0.752	–
	12	–	0.483	0.681	0.754	–
$n = 1000$	1	0.847	0.907	0.961	0.957	0.953
	2	–	0.904	0.960	0.958	–
	3	–	0.904	0.959	0.959	–
	5	–	0.912	0.957	0.959	–
	7	–	0.914	0.951	0.956	–
	9	–	0.893	0.945	0.961	–
	12	–	0.810	0.932	0.959	–

This table shows the rejection power of 95%-confidence intervals of PATE estimators for different allocation designs and experimental sample sizes. For each allocation design, the power estimates are computed over 10,000 simulations. For columns 1, 3, and 4, the estimator used is the Horvitz-Thompson algorithm. In column 1, the design used is complete randomization and the confidence intervals are constructed using White standard errors. For column 3, we assign treatment using a matched pairs design, pairing individuals to the closest unit and using the Mahalanobis distance whenever more than one covariate is balanced. Confidence intervals are constructed as described in (Bai, 2022). Column 4 shows the results for the cube method with two moments for each variable. Confidence intervals follow the asymptotic-based procedure described in Section 4.3. For column 2, we run stratified randomization. We use median values for continuous variables and we estimate the PATE using OLS with strata fixed-effects in accordance with Bugni et al. (2018). Finally, column 5 uses a double-lasso selection procedure as described in Belloni et al. (2014) with White standard errors for the post-lasso estimation.