



## Rerandomization Strategies for Balancing Covariates Using Pre-Experimental Longitudinal Data

Per Johansson & Mårten Schultzberg

To cite this article: Per Johansson & Mårten Schultzberg (2020) Rerandomization Strategies for Balancing Covariates Using Pre-Experimental Longitudinal Data, Journal of Computational and Graphical Statistics, 29:4, 798-813, DOI: [10.1080/10618600.2020.1753531](https://doi.org/10.1080/10618600.2020.1753531)

To link to this article: <https://doi.org/10.1080/10618600.2020.1753531>



© 2020 The Author(s). Published with license by Taylor & Francis Group, LLC.



Published online: 22 May 2020.



Submit your article to this journal [↗](#)



Article views: 1022



View related articles [↗](#)



View Crossmark data [↗](#)



Citing articles: 1 View citing articles [↗](#)

# Rerandomization Strategies for Balancing Covariates Using Pre-Experimental Longitudinal Data

Per Johansson<sup>a,b</sup> and Mårten Schultzberg<sup>c</sup>

<sup>a</sup>Uppsala University and IFAU, Uppsala, Sweden; <sup>b</sup>Tsinghua University, Beijing, China; <sup>c</sup>Uppsala University, Uppsala, Sweden

## ABSTRACT

This article considers experimental design based on the strategy of rerandomization to increase the efficiency in experiments. Two aspects of rerandomization are addressed. First, we propose a two-stage allocation sample scheme for randomization inference to the units in experiments that guarantees that the difference-in-mean estimator is an unbiased estimator of the sample average treatment effect for any experiment, conserves the exactness of randomization inference, and halves the time consumption of the rerandomization design. Second, we propose a rank-based covariate-balance measure which can take into account the estimated relative weight of each covariate. Several strategies for estimating these weights using pre-experimental data are proposed. Using Monte Carlo simulations, the proposed strategies are compared to complete randomization and Mahalanobis-based rerandomization. An empirical example is given where the power of a mean difference test of electricity consumption of 54 households is increased by 99%, in comparison to complete randomization, using one of the proposed designs based on high frequency longitudinal electricity consumption data. Supplementary materials for this article are available online.

## ARTICLE HISTORY

Received January 2019  
Revised March 2020

## KEYWORDS

Experimental design; Fisher exact test; High frequency longitudinal data; Mahalanobis distance criterion; Mirror allocation sampling; Rerandomization

## 1. Introduction

There is today a substantial literature on how to design randomized experiments to increase the efficiency as compared to complete randomization. All designs try to improve the similarity in the potential outcome of the groups of comparison by, in different ways, making the groups balanced in covariates that are observed before the experimental design is decided. The most common design used to improve balance is stratified or blocked randomization. The idea of stratified randomization is to divide units into strata (i.e., groups or blocks) based on similarity on covariates and then perform complete randomization within each strata. In this way, units from all strata will be represented in both the treatment and control groups<sup>1</sup> and thereby imbalance in any of these covariates are avoided (see, e.g., Imbens and Rubin 2015 for a recent overview).

An alternative design is rerandomization which was originally suggested by Fisher in the early twentieth century but was first written up and implemented in Morgan and Rubin (2012). As the name suggests, rerandomization consists of redoing the randomization until some prespecified balance criterion on the observed covariates is met. That is, the randomization is restricted to a subset of admissible allocations that fulfill the rerandomization covariate balance criterion. Based on the affinely invariant Mahalanobis distance covariate balance criterion, Morgan and Rubin (2012) showed that rerandomization can decrease the variance in the effect estimate substantially


as compared to complete randomization. Morgan and Rubin (2015) extended the results in Morgan and Rubin (2012) to deal with the case when large numbers of covariates are available and when some covariates can be a priori defined as more important than others.

The strategy of rerandomization is especially useful when continuous covariates are available, as in principle, even a single continuous covariate implies infinitely many strata and therefore must be discretized with information loss as a consequence. Compared to stratification, rerandomization is computationally demanding. However, with today's computers it provides a very interesting and powerful alternative design. As also Morgan and Rubin (2012) pointed out, rerandomization is not a design strategy that replaces stratification, rather a researcher should block on what covariates are possible and then use rerandomization on remaining covariates within these strata. Like with blocked designs, it is possible to apply rerandomization in sequential randomization designs (Zhou et al. 2018), as is often used in, for example, clinical trials.

One potential caveat with rerandomization is that common test-statistics are no longer asymptotic normally distributed. Li, Ding, and Rubin (2018) showed that the specific asymptotic distribution under rerandomization depends on which covariate balance measure is used and derive the asymptotic sampling distribution of the difference-in-mean estimator under rerandomization based on the Mahalanobis distance measure.

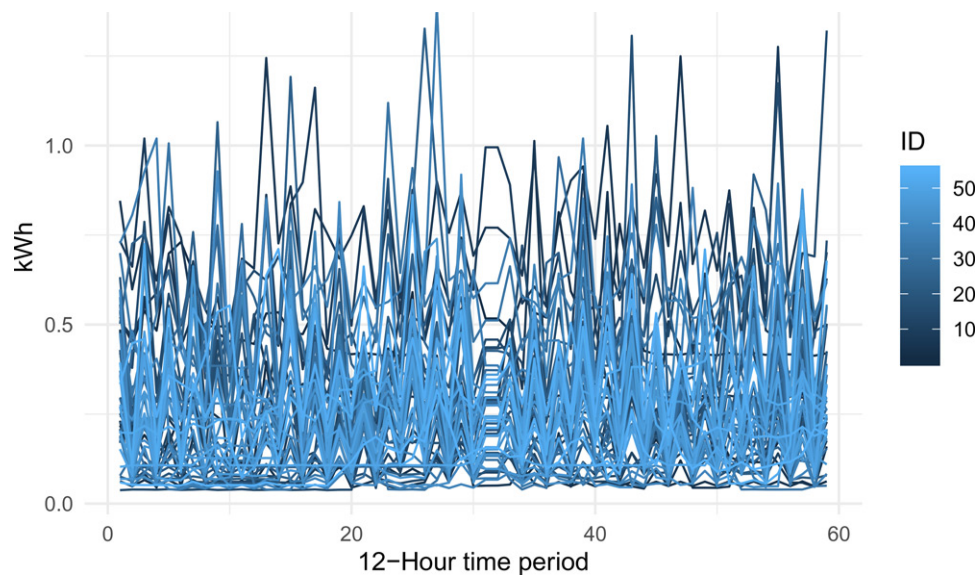
**CONTACT** Per Johansson  [per.johansson@statistics.uu.se](mailto:per.johansson@statistics.uu.se)  Uppsala University and IFAU, Uppsala, Sweden.

<sup>1</sup>The concepts in this article extend to any number of treatment groups. To simplify discussions, the number of treatment groups is restricted throughout this article to two, referred to as treatment and control.

 Supplementary materials for this article are available online. Please go to [www.tandfonline.com/r/JCGS](http://www.tandfonline.com/r/JCGS).

© 2020 The Author(s). Published with license by Taylor and Francis Group, LLC

This is an Open Access article distributed under the terms of the Creative Commons Attribution-NonCommercial-NoDerivatives License (<http://creativecommons.org/licenses/by-nc-nd/4.0>), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the original work is properly cited, and is not altered, transformed, or built upon in any way.



**Figure 1.** The electricity consumption (kWh) during the pretreatment month for the 54 households included in the experiment.

One appealing alternative to asymptotic inference was given in Morgan and Rubin (2012), namely, to restrict the inference to the units in the experiments and to base the analysis on exact (randomization) inference (Fisher 1935). Due to the assumption-free nature of the randomization inference, this strategy is valid for all well-behaved (discussed in Section 3) balance criteria. It is advantageous if the subset of allocations from all admissible allocations is of a moderate size when conducting the exact inference. A formal strategy of choosing the “best” subset from the admissible is, to our knowledge, not available in the literature.

The contribution of this article is 2-fold. First, an allocation sampling scheme for choosing the approximate “best” subset of admissible allocations is proposed. In addition to providing the exact level for the exact test, the sampling scheme: (i) guarantees that the difference-in-mean estimator is unbiased for the sample average treatment effect (SATE) and (ii) reduces the computational time for the design by half. Second, we develop a rerandomization covariate balance measure that is easy to use when pre-experimental outcome data (possibly high frequency longitudinal) are available, a situation that has not previously been addressed in the literature.

The article should be of broad interest as the situation where the pretreatment outcome is observed for many time periods is becoming more common. The last few decades’ technological development of personal electronic devices like smart phones, smart watches, fitness trackers, and the “Internet of Things,” has made the collection of high frequency longitudinal data substantially simplified and cheaper. This development has also led to an increased interest in what kinds of research questions these data might help us answer (see, e.g., Hamaker and Wichers 2017; Hamaker et al. 2018). The present article points out yet another possibility that these data brings, namely that of improving designs, enabling informative causal analysis also in relatively small experiments. In addition, we present practical guidelines for any rerandomization design that should be useful for any practitioner that want to use this strategy.

As a motivating example, data from an electricity consumption experiment are considered. In the design stage, repeated measurements of the outcome (kWh), displayed in Figure 1, are available for a sample of 54 households for the month before the treatment assignment. The average consumption for every 12 hr period, that is, 60 measurements per month is observed. Clearly, there are large variations in several aspects (level, variation, etc.) of the consumption behavior the months before the treatment assignment. The interest in this particular study is to see how user’s electricity consumptions behavior is affected by information campaigns about money saving consumption behaviors. Using one of the strategies proposed in this article, the power of a mean difference test with equal sized groups at the first time period after the pretreatment period, is increased by 99% as compared to complete randomization. All the details of this study are presented in Section 6.

The rest of the article is structured as follows. Section 2 introduces the concept of rerandomization and the Mahalanobis distance rerandomization criterion specifically. Section 3 introduces the formalized allocation sample scheme to select the subset of admissible allocations for exact inference. Section 4 introduces the new rerandomization covariate balance measure based on the ranks of the mean differences in the covariates. Section 5 provides a Monte Carlo simulation where the performance of the new balance measure is compared to Mahalanobis-based rerandomization and complete randomization. Section 6 provides the empirical analysis of electricity consumption data illustrating the proposed procedure and balance measure. Section 7 contains a discussion and concludes the article.

## 2. Mahalanobis-Based Rerandomization

To provide a better understanding of the underpinnings of the rerandomization framework proposed by Morgan and Rubin (2012), it is here compared to complete randomization and classical stratified design.

The technical difference between complete and stratified randomization is that allocations that are possible in complete randomization are excluded in the stratified randomization. More precisely, the allocations from complete randomization associated with imbalances in the stratification covariates are excluded. For example, consider a study where two equal sized treatment and control groups will be compared. A sample of 10 males and 10 females are randomly sampled from the population. Under complete randomization there are  $\binom{20}{10} = 184,756$  possible treatment allocations. If instead, randomization is stratified on sex, that is, 5 males and 5 females are allocated to treatment the number of possible treatment allocations are reduced to  $\binom{10}{5}\binom{10}{5} = 63,504$ , that is, the 121,252 ( $= 184,756 - 63,504$ ) allocations that are unbalanced on sex are excluded. To stratify, the covariate space needs to be partitioned into finite sets. With a set of a few categorical covariates this strategy is easy to implement, at least if there is more than one individual within each strata. However, as mentioned above, if one would like to include continuous covariates (e.g., pre-experiment outcomes) in the design, simple stratification methods run into problems as continuous covariates must be discretized (Hu and Hu 2012).

Rerandomization is similar to stratified randomization in the sense that certain allocations are excluded, the main difference is the exclusion criterion. To formally introduce the concept of rerandomization, we start by describing the basic idea as it is outlined in Morgan and Rubin (2012) which also forms the basis for the extensions in Morgan and Rubin (2015), Zhou et al. (2018), and Li, Ding, and Rubin (2018).

Consider a trial with  $N$  individuals of which  $N/2 = N_1 = N_0$  are assigned to treatment and control, respectively. Let  $\mathbf{z}$  be a fixed  $N \times K$  matrix containing variables observed prior to the treatment assignment.  $\mathbf{z}$  may contain both covariates and pretreatment outcomes. When the separation is important,  $\mathbf{x}$  denotes the  $N \times K_x$  covariate matrix, and  $\mathbf{y}_{pre}$  denotes the  $N \times K_y$  pretreatment outcomes matrix, where  $K = K_x + K_y$ .

Let  $W_i = 1$  if individual  $i$  is treated and  $W_i = 0$  if not. Let  $\mathbf{W}$  be the matrix of all  $\binom{N}{N/2} = N_A$  possible random assignments (i.e., before treatment groups are assigned). For a given allocation  $j, j = 1, \dots, N_A$  the Mahalanobis distance between the covariate mean vectors of those assigned to treatment ( $T$ ) and control ( $C$ ), respectively, is defined as

$$M(\mathbf{z}, \mathbf{w}^j) = M_j = \frac{N}{4} (\bar{\mathbf{z}}_T^j - \bar{\mathbf{z}}_C^j)' \text{cov}(\mathbf{z})^{-1} (\bar{\mathbf{z}}_T^j - \bar{\mathbf{z}}_C^j),$$

$$j = 1, \dots, N_A, \quad (1)$$

where  $\mathbf{w}^j$  is the  $j$ th column vector in  $\mathbf{W}$ ,  $\text{cov}(\mathbf{z})$  is the sample covariance matrix and  $\bar{\mathbf{z}}_T^j - \bar{\mathbf{z}}_C^j$  is the difference in mean vectors which is a  $K \times 1$  stochastic vector as it depends on the allocation  $j$ . Morgan and Rubin (2012) suggested randomizing within the set

$$\{\mathbf{W} | M(\mathbf{z}, \mathbf{w}^j) - a \leq 0\}, \quad (2)$$

where  $a$  is a constant. This means that instead of randomly choosing one of the  $N_A$  possible allocations, a smaller set of allocations with small Mahalanobis distances is considered. If the means are normally distributed then  $M^j \sim \chi^2(K)$ . This means that  $a$  can be indirectly determined by setting  $p_a$  in

$$p_a = \Pr(\chi^2(K) \leq a).$$

As the number of rerandomizations is geometrically distributed with expected value  $1/p_a$ , the expected number of rerandomizations before drawing a randomization fulfilling the criterion with, for example,  $p_a = 0.001$ , is 1000. As this holds for any  $N$ , the time it takes to find the allocations from which to finally make the randomization is independent of  $N$  for a fixed  $p_a$ .

If  $\mathbf{z}$  is ellipsoidally symmetric then, as a consequence of Mahalanobis metric being multivariate affinely invariant, the variation reduction is equally large for each covariate and, given  $M^j \sim \chi^2(K)$ , equal to

$$\text{cov}(\bar{\mathbf{z}}_T^j - \bar{\mathbf{z}}_C^j | \mathbf{z}, M^j \leq a) = v_a \text{cov}(\bar{\mathbf{z}}_T^j - \bar{\mathbf{z}}_C^j | \mathbf{z}), \quad (3)$$

where

$$v_a = \frac{\Pr(\chi^2(K+2) \leq a)}{\Pr(\chi^2(K) \leq a)}; \quad 0 < v_a < 1.$$

This implies that the variance of the covariates in the subset of allocations is reduced in comparison to the complete randomization. The equal percent variance reduction (EPVR) of the included covariates is equal to

$$100(1 - v_a). \quad (4)$$

Furthermore,  $v_a$  can be written as

$$\frac{2}{K} \times \frac{\gamma(K/2 + 1, a/2)}{\gamma(K/2, a/2)}, \quad (5)$$

where  $\gamma(b, c) = \int_0^c y^{b-1} e^{-y} dy$ . Equation (5) shows that the variance reduction is increasing in  $a$  and decreasing in  $K$ .

Let  $Y_i(0)$  be the potential outcome if not treated and let  $R^2$  be the coefficient of determination of a regression of  $\mathbf{Y}(0)$  on  $\mathbf{z}$ , where  $\mathbf{Y}(0) = (Y_1(0), \dots, Y_N(0))'$ . Under the assumptions of conditionally normally distributed outcomes and additive homogeneous treatment effects, Morgan and Rubin (2012) showed that the percent reduction in variance of the differences in mean estimator is equal to

$$100R^2(1 - v_a). \quad (6)$$

In the two following sections, the two main contributions of this article are presented, respectively. First, the practical strategy for implementing rerandomization for a class of rerandomization balance measures is proposed, after which the new covariate balance measure is presented.

### 3. Sampling Scheme for Exact Inference Under Rerandomization

This section presents a formal strategy for how to implement rerandomization with exact inference for empirical applications. This strategy applies to all balance measures with symmetric criterion functions, defined in detail below.

The exact  $p$ -value (Fisher 1935) of a test-statistic for a randomly selected allocation is obtained from the percentile of the histogram of the observed values of the test-statistic for all possible allocations. In complete randomization there are  $N_A$  possible allocations. To calculate the exact  $p$ -value, the test statistic is calculated for all  $j = 1, \dots, N_A$  allocations to create the histogram. The implication of this procedure is that the smallest possible  $p$ -values under the null hypothesis will be restricted



by the size  $N$  and whether there are ties in the test statistic. If the number of possible allocations is restricted further, as when using rerandomization, the set of allocations fulfilling the rerandomization criteria must be kept large enough to calculate the desired percentiles of the randomization distribution. This means that, randomization within a set of “best” allocations is valid only as long as the set of allocations is large enough to obtain the desired percentile, what we define as the *resolution* of the exact  $p$ -value.

That is, the upper bound of the *resolution* of the exact  $p$ -value from a two-sided hypothesis test is  $2/S_{\text{unique}}$ , where  $S_{\text{unique}}$  is the number of unique test-statistic values in the set of admissible allocations.<sup>2</sup> For a guaranteed resolution,  $r$ , on the hundredth in the two-sided  $p$ -value, 200 allocations must be selected, given no ties in the statistic. For continuous outcomes and a statistic that is a smooth function of the outcome, for example, the difference-in-means, all allocations are likely to be unique.

For most sample sizes, the number of allocations fulfilling even very strict rerandomization criteria is usually large why too low resolution of the  $p$ -value is not usually an issue. However, too large number of allocations means that it becomes intractable to calculate the exact  $p$ -value based on *all* admissible allocations. One alternative is to randomly draw an allocation from the admissible allocations, and then Monte Carlo approximate the exact  $p$ -value. Another alternative is to first strategically choose a subset of  $H$  allocations from all admissible allocations and then randomly draw one of these  $H$  allocations. The exact  $p$ -value can then be calculated over only the subset of the  $H$  allocations. This procedure will in general require sampling a much smaller number of allocations than with Monte Carlo approximation and provides, by definition, the correct level of the hypothesis test for an effect. In the following section we discuss a strategy of finding the “best”  $H$  allocations during a fixed computational time.

### 3.1. Criterion Function for Admissible Allocations

Define the general rerandomization criterion function

$$\varphi(\mathbf{z}, \mathbf{w}^j, B, c) = \begin{cases} 1, & \text{if } B(\mathbf{z}, \mathbf{w}^j) \leq c, \\ 0, & \text{otherwise,} \end{cases} \quad (7)$$

where  $B(\mathbf{z}, \mathbf{w}^j)$  is the scalar covariate balance measurement of allocation  $j$  for the set of covariates  $\mathbf{z}$ , and  $c$  is a criterion. Let  $\mathbf{W}^\varphi$  be the subset of admissible allocations, that is, the set of allocations  $j$  for which  $\varphi(\mathbf{z}, \mathbf{w}^j, B, c) \equiv 1$ . As an example, for the Mahalanobis distance balance measure with inclusion criterion  $a$ ,

$$\varphi(\mathbf{z}, \mathbf{w}^j, M, a) = \begin{cases} 1, & \text{if } M(\mathbf{z}, \mathbf{w}^j) \leq a, \\ 0, & \text{otherwise.} \end{cases} \quad (8)$$

<sup>2</sup>The uniqueness is required to be able to calculate the unique rank on which the  $p$ -value is based. Given no ties in the test statistic the resolution of the  $p$ -values from a two-sided hypothesis test in a balanced complete randomized experiment is equal to  $2/\binom{N}{N/2}$ .

### 3.2. Mirror Allocations and Unbiased Estimators of SATE Under Rerandomization

In addition to ensuring unbiasedness of the difference-in-means estimator, the inclusion of mirror allocations (pairs of mirror allocations) has the advantage of reducing the time it takes to find a fixed number of allocations by a factor of two. This follows since for any symmetric balance measure, if an allocation is admissible, so is the corresponding mirror allocation and thus only the covariate balance of one allocation in each pair of mirror allocations has to be evaluated.<sup>3</sup>

For the symmetric balance measure  $B$ , let  $B_j$  be the balance measurement of the  $j$ th allocation for  $j = 1, \dots, N_A$ , where the allocations are ordered lexicographic such that the first  $N_A/2$  contains no pair of mirror allocations. Based on this measure, the best allocation among the  $N_A/2$  first allocations is given by

$$\mathbf{w}_{(1)} = \underset{\mathbf{w}^j}{\operatorname{argmin}} B_j, \quad j = 1, \dots, N_A/2.$$

That is, the “best” allocation is the allocation with the smallest imbalance. If, for example, the balance measure is the Mahalanobis distance, the smallest possible Mahalanobis distance for a given sample is  $M_{(1)} = M(\mathbf{z}, \mathbf{w}_{(1)})$ .

A reasonable goal in a rerandomization design is that of finding the set of the  $H$  allocations with the globally smallest imbalances  $\mathbf{W}_{(1:H)}$ , defined as

$$\mathbf{W}_{(1:H)} = \{\mathbf{w}_{(1)}, \dots, \mathbf{w}_{(H/2)}, \mathbf{w}_{(1)}^M, \dots, \mathbf{w}_{(H/2)}^M\}, \quad (9)$$

where  $\mathbf{w}_{(o)}$ ,  $o = 1, \dots, N_A/2$  is the allocations with the  $N_A/2$  smallest imbalances and  $\mathbf{w}_{(o)}^M$  is the mirror allocation to  $\mathbf{w}_{(o)}$ .

Finding  $\mathbf{W}_{(1:H)}$  requires that we compute all balance measurement in the set  $\mathbf{W}$ . However,  $N_A$  is often so large that computational time is a restriction in finding  $\mathbf{W}_{(1:H)}$ . For this reason we suggest an algorithm that finds the *approximate best set* of allocations for large  $N$  within, what can be defined as, a reasonable time limit.

Define the approximate best set

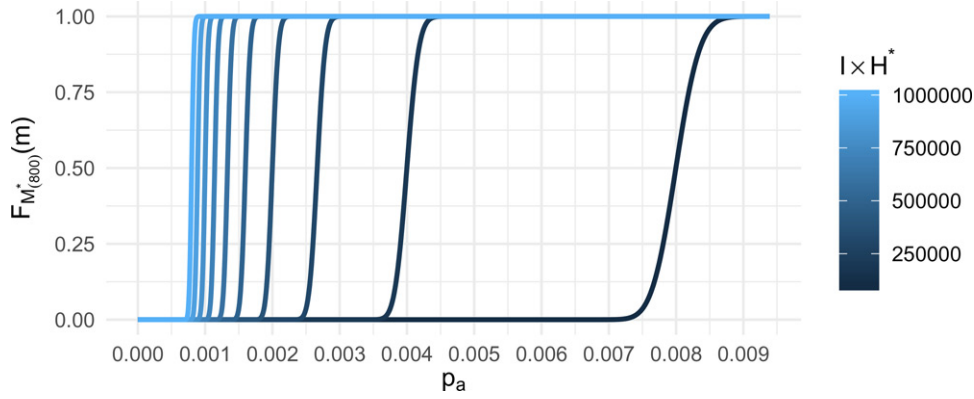
$$\mathbf{W}_{(1:H)}^* = \{\mathbf{w}_{(1)}^*, \dots, \mathbf{w}_{(H/2)}^*, \mathbf{w}_{(1)}^{*M}, \dots, \mathbf{w}_{(H/2)}^{*M}\}, \quad (10)$$

where  $\mathbf{w}_{(o)}^*$  for  $o = 1, \dots, \operatorname{card}(\mathbf{W}^*)/2$  is the allocation with  $o$ th smallest  $B_j$  for  $j \in \mathbf{W}^* \subset \mathbf{W}$ ,  $\mathbf{W}^*$  is the set of considered allocations, and  $\operatorname{card}(\mathbf{W}^*) = N_s$  is the cardinality of this set.

To illustrate the idea of finding the set  $\mathbf{W}_{(1:H)}^*$ , let  $N_s = I \times H^*$  where  $I$  is the number of times we are reading  $H^*$  allocations from the disk. Let  $\mathbf{C}_M$  denote the  $N_A/2$  allocations where all mirror allocations are removed and assume for simplicity that the columns in  $\mathbf{C}_M$  are ordered according to the order they are being randomly selected, that is  $\mathbf{c}_M^1$  is the first randomly selected allocation vector, and so on. Let  $i = 1, \dots, I$ , and let  $\mathbf{A}_1$  be the first set of  $H^*$  randomly selected allocations from  $\mathbf{C}_M$ .

For concretization, assume that we use the Mahalanobis balance measure to find the  $H$  best allocation in  $\mathbf{W}_{(1:H)}^*$ . This means that we need to sort  $M(\mathbf{z}, \mathbf{c}_M^j), j = 1, \dots, H^*$  and then keep the  $H/2$  with the best balance in  $\mathbf{A}_1$ . Denote this set

<sup>3</sup>Note that in a balanced factorial experiment with  $G$  treatments, the time it takes to compare the balance measure of any number of allocations is reduced by a factor of  $1/G!$  when using this mirror allocation design.



**Figure 2.** The cumulative distribution function for the Mahalanobis distance ( $df = 3$ ) of the 800th order statistic in the random sample of allocations for numbers of allocations between 100,000 and 1,000,000.

$\mathbf{M}_{(1)}^1 = (M(\mathbf{z}, \mathbf{a}_{1(1)}), \dots, M(\mathbf{z}, \mathbf{a}_{1(H/2)}))$ , where  $\mathbf{a}_{1(o)}$  is a vector of allocations with the  $o$ th smallest imbalance in  $\mathbf{A}_1$ . For  $i = 2$ , we select  $H^*$  new allocations in  $\mathbf{A}_2$  and calculate  $M(\mathbf{z}, \mathbf{c}_M^j)$ ,  $j = H^* + 1, \dots, 2H^*$ . The balance measure of these  $H^* + H/2$  allocations are then sorted and the  $H/2$  allocations with best balance in  $\mathbf{A}_1 \cup \mathbf{A}_2$  are stored in  $\mathbf{M}_{(1)}^2$ . For  $i = 3$ , we select  $H^*$  new allocations in the set  $\mathbf{A}_3$ , and so on. For the final set  $\mathbf{W}^*$  we add the mirror allocations.

The proposed allocation sample scheme, using the mirror-allocation strategy and for the Mahalanobis distance criterion, is presented in detailed in Allocation sample scheme 1.

**Allocation sample scheme 1.** Order all possible allocations lexicographically<sup>4</sup> such that the first  $N_A/2$  allocations contains no pairs of mirror allocations. Call the set containing the first  $N_A/2$  allocations  $\mathbf{C}_M$ . Choose a desired level of resolution  $r$ . Assuming no ties, this gives the number of allocation  $H = 2/r$ .

1. Randomly sample, without replacement, a set  $\mathbf{A}_1$  from  $\mathbf{C}_M$  containing  $H^*$  allocations, where  $H^* \geq H/2$ .
2. Calculate the balance measure for all allocations in  $\mathbf{A}_1$ , store the  $H/2$  allocation with the smallest imbalances from  $\mathbf{A}_1$  in the set  $\mathbf{W}_s^*$ .
3. For  $i = 2, \dots, I$ 
  - (a) Sample a new set of allocations,  $\mathbf{A}_i$ , of size  $H^*$  from the  $N_A/2 - (i - 1) \times H^*$  remaining allocations in  $\mathbf{C}_M$ .
  - (b) Calculate the Mahalanobis distance for  $\mathbf{A}_i$ . Replace the set  $\mathbf{W}_s^*$  with the  $H/2$  allocations with smallest imbalances in the set  $\{\mathbf{W}_s^*, \mathbf{A}_i\}$ .
4. Save  $\mathbf{W}_{(1:H)}^* = \{\mathbf{W}_s^*, \mathbf{W}^{M*}\}$  as the final subset of admissible allocations, where  $\mathbf{W}^{M*}$  contains all mirror allocations for the allocations in  $\mathbf{W}_s^*$ .

With this algorithm,  $\mathbf{W}^* = \mathbf{A}_1 \cup \mathbf{A}_2 \cup \dots \cup \mathbf{A}_{I-1} \cup \mathbf{A}_I$ . As the sampling of allocations are random,  $\mathbf{W}_{(1:H)}^*$  is an approximation of the set containing the globally best allocations  $\mathbf{W}_{(1:H)}$ . This approximation will be better the longer the sampling of allocations is allowed to continue since  $\mathbf{W}^* \rightarrow \mathbf{W}$  as  $I \times H^* \rightarrow N_A$ .

If the balance measure is the Mahalanobis distance, the set containing the  $H$  allocations with the globally smallest

imbalances follows from Equation (9). For  $M_{(j)} = M(\mathbf{z}, \mathbf{w}_{(j)})$ , the set of the globally smallest Mahalanobis distances is  $\{M_{(1)}, \dots, M_{(H/2)}, M_{(1)}^M, \dots, M_{(H/2)}^M\}$ . That is, the allocations with the  $H$  first-order statistics of the Mahalanobis distance in the set of all  $N_A$  allocations. This corresponds to  $a = M_{(H/2)}$  in  $p_a = \Pr(\chi^2(K) \leq a)$ . Let  $M_{(H/2)}^*$  be the Mahalanobis distance of the pair of mirror allocation with the  $H/2$  smallest Mahalanobis distance in  $\mathbf{W}^*$ . The sample scheme implies that  $p_{a^*} = \Pr(\chi^2(K) < M_{(H/2)}^*)$  is unknown and stochastically dependent on  $I \times H^*$ . This is in contrast to Morgan and Rubin (2012) where the inclusion criterion  $p_a$  is fixed and the number of included allocations is random for a fixed number of rerandomizations. If the number of considered allocations is small, that is,  $I \times H^* \ll \binom{N}{N_1}$ , then  $M_{(H/2)}^*$ , and thus  $p_{a^*}$ , may be large. However, as  $I \times H^* \rightarrow N_A$  it follows that  $M_{(H/2)}^* \rightarrow M_{(H/2)}$  and  $p_{a^*} \rightarrow \Pr(\chi^2(K) \leq M_{(H/2)})$ .

For chi-square distributed Mahalanobis distances, that is, under normal covariates or a large sample (Morgan and Rubin 2012), the distribution of the  $H$ th order statistic is known for any  $I \times H^*$ , given by

$$F_{M_{(H)}^*}(m) = F_{\beta(H, I \times H^* - H)}\left(F_{\chi_{(K)}^2}(m)\right), \quad (11)$$

where  $F_{M_{(H)}^*}$ ,  $F_{\beta(H, I \times H^* - H)}$ , and  $F_{\chi_{(K)}^2}$  denotes cumulative distribution functions (CDFs) of the  $H$ th order statistics of the Mahalanobis distances of a random sample of allocations, the  $\beta(H, I \times H^* - H)$  distribution, and the  $\chi^2(K)$  distribution, respectively (Gut 2009, chap. 4). This means we can calculate (numerically) the probability of  $p_{a^*} < \epsilon$  implied by  $M_{(H/2)}^*$  for any given  $I \times H^*$ , for an arbitrary  $\epsilon \in (0, 1)$ . As an example, Figure 2 displays the CDF of the 800th ( $H/2 = 800$ ) order statistic for a  $\chi^2(3)$  distribution as a function of  $p_a$  for different  $I \times H^*$ . It is clear from the figure that, with probability close to one,  $\Pr(\chi^2(K) \leq M_{(800)}^*) = 0.009$  when  $I \times H^* = 100,000$ . Thus, with probability close to one, the set  $\mathbf{W}_{(1:800)}^*$  only contains allocations from the 1% of the globally best allocations, or better. With  $I \times H^* = 1,000,000$  the corresponding probability equals 0.001, that is, the set  $\mathbf{W}_{(1:800)}^*$  contains allocations from the 0.1% of the globally best allocations, or better. This implies that the cardinality of  $\mathbf{W}^*$  can be chosen such that the allocations in the obtained set  $\mathbf{W}_{(1:H)}^*$ , with a probability arbitrarily close to one, all will have Mahalanobis distances smaller than some desired value.

<sup>4</sup>Most programming language have functions for generating combinations in lexicographic order, for example, the `RcppAlgos`-package in R.

The size of  $H^*$  will affect the computational time as it depends on the speed of reading from the disk and the time it takes to sort the  $H^* + H/2$  allocations in each iteration. For each iteration we would like to find as many allocations as possible with better balance than as the currently best  $H/2$  allocations. This means that  $H^*$  should be as large as possible, making  $I$  as small as possible for a fixed  $N_s$ . However, a too large set  $A_i$  in memory slows the sorting time down substantially. This means that if we want to have as large  $N_s$  as possible for a given computation time,  $H^*$  should be tuned to the memory capacity of the computer and with the speed of reading from the disk. In summary, there are two parameters that governs the proposed sample scheme. The number of included allocations  $H$  (implied by the desired resolution of the exact  $p$ -value) and the number of considered allocations ( $I \times H^*$ ). The number of iterations in step 3 can as consequence be decided based on computational time given an “optimal” choice of  $H^*$ .<sup>5</sup>

#### 4. A Rank-Based Balance Measure

In this section, we present the second contribution of this article, an alternative rerandomization balance measure based on the rank of the covariate mean difference. The suggested balance measure is marginal affinely invariant and avoids the problem with inverting potentially large and singular covariance matrices that might occur using the Mahalanobis distance, especially with highly correlated covariates. In addition, the suggested balance measure provides a convenient alternative to Morgan and Rubin (2015) if the researcher has a priori information on the relative importance of the observed covariates.

Define the mean of the covariate  $k$  for a given allocation  $j$  for those assigned treatment and control  $\bar{z}_{kT}^j = \frac{1}{N/2} \sum_{i:W_i^j=1} z_{ik}$  and  $\bar{z}_{kC}^j = \frac{1}{N/2} \sum_{i:W_i^j=0} z_{ik}$ , respectively. Consider the covariate balance measure

$$\sum_{k=1}^K \text{Rank}(|\bar{z}_{kT}^j - \bar{z}_{kC}^j|), \tag{12}$$

where the rank is calculated over  $j = 1, \dots, \text{card}(\mathbf{W})$ . Again, let  $\mathbf{w}_{(j)}$  be the allocation with the  $j$ th smallest value of (12) over the first half of the lexicographic order of allocations, the set of the  $H$ th best allocation as  $\mathbf{W}_{(1:H)} = \{\mathbf{w}_{(1)}, \dots, \mathbf{w}_{(H/2)}, \mathbf{w}_{(1)}^M, \dots, \mathbf{w}_{(H/2)}^M\}$ . As, for any covariate, the exact  $p$ -value from a two-sided mean difference test within a set of allocations is a monotonous transformation of the rank of the absolute differences within the same set the  $p$ -values can be used instead of the ranks to simplify notation and to avoid large numeric values in large sets of allocations. Let  $p_{jk}$  denote the exact  $p$ -value of the two-sided mean difference test for allocation  $j$  and covariate  $k$ , and define the balance measure

$$R(\mathbf{z}, \mathbf{w}^j) = - \left( \sum_{k=1}^K \omega_k p_{jk} \right), \tag{13}$$

where  $\omega_k$  is the weight given to covariate  $k$ , with  $\sum_{k=1}^K \omega_k = 1$ . Here, the largest  $p$ -value implies the smallest absolute mean difference. Therefore, to fit the general criterion function, the negative sign is added so that the  $H$  allocations with the smallest  $R(\mathbf{z}, \mathbf{w}^j)$ ,  $j = 1, \dots, N_A$  gives the set  $\mathbf{W}_{(1:H)}$ .  $R$  is a symmetric balance measure for any choices of  $\omega_k$ 's. Using Allocation sample scheme 1, the criterion function for  $R$  has a criterion that itself is a function of  $\mathbf{W}^*$ , that is,

$$\varphi(\mathbf{z}, \mathbf{w}^j, R, R(\mathbf{z}, \mathbf{W}_{(H)}^*)) = \begin{cases} 1, & \text{if } R(\mathbf{z}, \mathbf{w}^j) \leq R(\mathbf{z}, \mathbf{W}_{(H)}^*), \\ 0, & \text{if } R(\mathbf{z}, \mathbf{w}^j) > R(\mathbf{z}, \mathbf{W}_{(H)}^*), \end{cases} \tag{14}$$

where  $\mathbf{w}^j \in \mathbf{W}^*$  for all  $j$ .

With one single covariate, that is,  $K = 1$ , the  $H$  allocations with the smallest imbalances, that is, the allocations with the largest  $R(\mathbf{z}, \mathbf{w}^j)$  among  $j = 1, \dots, N_A$ , are the same allocations that have the  $H$  smallest Mahalanobis distances. This implies that with one covariate the rerandomization based on this new balance measure yields identical variance reductions as in Morgan and Rubin (2012) and, hence, that the percent reduction in variance of the estimator under normality of the outcome is given in Equation (6). However, for  $K > 1$  the two balance measures may give different sets of allocations. The proposed measure does not explicitly use the covariance structure of the covariate matrix, but instead make use of marginal weights for each covariate. However, in situations where the weights can be estimated from data, the covariate covariance structure can be taken into account in the weights themselves. The following section discuss this further.

The balance measure  $R(\mathbf{z}, \mathbf{w}^j)$  is robust against outliers in single covariates, as the impact of the imbalance in one covariate is bounded by the exact  $p$ -value, that is,  $[0,1]$ . The Mahalanobis balance measure has no upper bound for the influence of a single covariates. This difference could be important for small samples where single observations are more influential. A drawback with the measure is an increase in computational time in contrast to the Mahalanobis measure. The reason is that the  $p$ -values in the criterion is fixed in the set  $\mathbf{W}$ , but not across subsets. The implication is that the  $p$ -values needs to be calculated in the set  $\{\mathbf{W}_s^*, \mathbf{A}_i\}$ , for all  $i = 1, \dots, I$ . This means that the 3(b) in Algorithm 1 with the rank-based measure is 3(b'): Calculate the balance measure for the set  $\{\mathbf{W}_s^*, \mathbf{A}_i\}$ . Replace the set  $\mathbf{W}_s^*$  with the  $H/2$  allocations with the smallest imbalances.<sup>6</sup>

Throughout the remainder of this article, we denote the balance measure using the procedure given in Equation (13)  $R$ , while the procedure in Equation (1) using the Mahalanobis distance balance measure is denoted  $M$ .  $R$  with uniform weights is simply denoted  $R$  while if the  $R$  balance measure is based on nonuniform weights it is denoted  $R(\omega, G)$ , where  $G$  is a generic term defining the method being used to set these weights.

Morgan and Rubin (2015) addressed the case when covariates vary in importance and suggest rerandomization based on Mahalanobis distance within tiers of covariates, grouped by a priori importance. The  $R$  balance measure is a complement to Morgan and Rubin (2015) which is simple to implement and

<sup>5</sup>Making use of this procedure with the rank-based balance measure (see Section 4) in the empirical example (see Section 6), the allocation sampling scheme was left to work for 11 hr after which about one billion allocations had been considered and the  $H = 800$  approximately best had been retrieved.

<sup>6</sup>An example on the computational time of different choices of  $H^*$  is included in the supplementary R-file.

allows the weights  $\omega_k, k = 1, \dots, K$ , to be estimated. This ability is especially useful when  $K$  is large and some weights might be zero or close to zero. In the following section we discuss strategies for estimating the weights in the situation where (at least) one premeasured outcome is observed at the design stage.

#### 4.1. Estimating the Weights From Pretreatment Data

Assume data on the outcome is observed for  $T$  time periods at the time of the design and experiment. For this time period, we also observe a set of covariates. We first discuss the case with  $T = 1$ . This means that  $\mathbf{z} = \{\mathbf{Y}_1, \mathbf{x}_1, \dots, \mathbf{x}_{K_x}\}$  where  $\mathbf{Y}_1 = (Y_{11}, \dots, Y_{n1})'$ . Then, we turn to the case with  $T > 1$ . Focusing on the situation with  $K_x = 0$ .

##### 4.1.1. One Pretreatment Outcome

If the pretreatment outcome,  $\mathbf{Y}_1$ , can be assumed to be correlated with the post-treatment outcome in the absence of treatment,  $\mathbf{Y}_2(0)$ , an obvious strategy is to estimate the weights in Equation (13) using the partial correlation of the individual covariates and the pretreatment outcomes.

Denote the weights for  $\{\mathbf{Y}_1, x_1, \dots, x_{K_x}\}$  as  $\{\omega_0, \omega_1, \dots, \omega_{K_x}\}$ . The weight for the covariates are estimated using ordinary least square (OLS) on

$$\tilde{y}_{i1} = \sum_{k=1}^{K_x} \beta_k \tilde{x}_{ik} + \varepsilon_{iT},$$

where  $\sim$  indicate a standardized variables, that is,  $\tilde{z}_i = (Z_i - \bar{Z})/\sqrt{\text{var}(Z_i)}$ . With  $K_x$  covariates and one pretreatment outcome there are  $K = K_x + 1$  variables to base rerandomization on.<sup>7</sup> By normalizing the weights to sum to one we get  $\hat{\omega}' = (\frac{1}{\delta}, \frac{|\hat{\beta}_1|}{\delta}, \dots, \frac{|\hat{\beta}_{K_x}|}{\delta})$ , where  $\delta = \sum_{k=1}^{K_x} |\hat{\beta}_k| + 1$ . As all variables are standardized to have unit variance the coefficients are therefore bounded in theory, that is,  $0 \leq |\hat{\beta}_j| \leq 1 \forall j = 1, \dots, K_x$ , and it follows that  $\frac{1}{\delta} \geq \frac{|\hat{\beta}_j|}{\delta}$ . This means that the pretreatment outcome will always have the largest weight, which is natural given that the weights are estimated under the assumption that the pretreatment outcome is associated with the outcome at the time period of the experiment. This strategy is denoted as  $R(\omega, O)$  throughout the rest of this article.

If there are many covariates and/or the sample size is small in comparison to  $K_x$ , it might be useful to estimate the partial correlations with some regularization estimator, for example, LASSO (Tibshirani 1996). Using LASSO tuned with cross-validation, some covariates can be given zero weights, which might substantially reduce the noise in the weight estimation with many highly correlated covariates. This strategy, using leave-one-out cross-validation, is denoted as  $R(\omega, L)$ .

There is a close relation between Mahalanobis-based rerandomization and the proposed measure with regression-estimated weights. As can be seen from Equation (6), the variance reduction in the outcome in Mahalanobis-based rerandomization is governed by  $R^2$  and  $\nu_a$ . As  $\nu_a$  is decreasing in  $K$  (Equation (5)) this means that including unnecessary

covariates (i.e., not increasing  $R^2$ ) will reduce the potential efficiency gain from the rerandomization. This property stems from the fact that the covariates are treated symmetrically, that is, the percentage reduction in the imbalance is the same for all covariates (Equation (4)). Our measure with estimated weights will not reduce the imbalance equally for all covariates. If  $\mathbf{Y}_1$  is a good proxy for  $\mathbf{Y}_2(0)$ , better balance is obtained for the covariates that are most important for the potential outcome. Using LASSO, some estimated weights may be exactly zero, this corresponds to removing unimportant covariates from  $\mathbf{x}$ , which in Mahalanobis rerandomization would imply obtaining a smaller  $\nu_a$  but unchanged  $R^2$ . This would give larger balance improvements in the remaining covariates leading to larger variance reduction in the difference-in-means estimator. This indicates that the proposed balance measure with estimated weights is a useful complement to Mahalanobis rerandomization when there are many covariates and their importance is unknown a priori.

In the situation where the outcome is observed at several time periods pretreatment, that is,  $T > 1$ , LASSO is still useful to estimate the weights for the  $R$  balance measure, however, for larger numbers of time periods other strategies for performing the rerandomization might be preferable as is discussed in the next section.

##### 4.1.2. Several Pretreatment Outcomes

With  $\mathbf{z} = (\mathbf{Y}_1, \dots, \mathbf{Y}_T, \mathbf{x}_1, \dots, \mathbf{x}_{K_x})$  the strategies in the previous section applies directly. However, because the exact  $p$ -values used in the measure must be calculated for each pretreatment time period with nonzero estimated weight, the computational time estimating  $R(\omega, O)$  or  $R(\omega, L)$  increases drastically with the number of pretreatment time periods why alternative strategies may be preferable. We limit the discussion to the case with  $K_x = 0$  as the extension to the situation with  $K_x > 0$  is straightforward applying the results of the previous section. Thus, in the following  $\mathbf{z} = (\mathbf{Y}_1, \dots, \mathbf{Y}_{iT})$ .

An alternative in this situation is to predict the outcome value at the time period of the experiment by fitting a time series model to the pretreatment outcomes of each individual and then use the one step forecast to base the rerandomization on. In this setting the time-series prediction model is denoted  $R(\omega, F)$ . Since there is only one forecast value for each individual,  $M$  based on the forecast will give the same allocations as with  $R(\omega, F)$ . Note that the  $R(\omega, F)$  strategy is not only time saving, it also allows for heterogeneity across experimental units.

The two strategies  $R(\omega, L)$  and  $R(\omega, O)$  are likely to work well if the processes across individuals are homogeneous. If the processes are heterogeneous across the individuals, for example, with differences in the “memory” or time-dependency in the outcomes, the  $R(\omega, F)$  may be more efficient in reducing the variance as it can incorporate heterogeneity, although, at the cost of estimating more parameters. Thus, if  $T$  is large and heterogeneity is present, which is possible to detect using the pretreatment data,  $R(\omega, F)$  should likely be preferred over  $R(\omega, L)$  and  $R(\omega, O)$ .

The following section provides a Monte Carlo simulation study where the different strategies are evaluated and compared.

<sup>7</sup>The set of covariates can, of course, be extended to include transformations of the originally observed covariates.



### 5. Monte Carlo Simulation Studies

Throughout the Monte Carlo studies, complete randomization will serve as a benchmark for the reduction of variance of the covariates (including pretreatment outcomes) and the estimated treatment effect under the null for the different rerandomization strategies.

In traditional Neyman–Pearson asymptotic inference, a test with a small variance of the estimator will be asymptotically more powerful than a test based on an estimator with large variance given that both tests are based on consistent estimators of the effect and the variance. As the power of the Fisher randomized test (FRT) is based on shifts in the rank due to the shift under the alternative it is not possible to evaluate the power of two strategies based only on the variance of the estimators under null.<sup>8</sup> For this reason, the relative power under the alternative is also presented.

All evaluations of variance are made on the data for the period after the allocation is made, that is if the design and treatment assignment is performed at period  $T$ , the power and variance is calculated at  $T + 1$ . Denote complete randomization with  $c$ , and the different rerandomization strategies with  $d = M, R, R(\omega, O), R(\omega, L)$  and  $R(\omega, F)$ . For each arm of the study, 4000 replications ( $N_{\text{rep}}$ ) are considered.

Let  $W_{ir} = 1$  or  $0$  if unit  $i$  is treated or control in replicate  $r$  and let  $z_{rki}$  be value of covariate  $k$  for unit  $i$  in replication  $r$  and define

$$\begin{aligned} \bar{z}_{rqT} &= \frac{1}{N_1} \sum_{i=1}^N W_{ir} z_{rqi}, \\ \bar{z}_{rqC} &= \frac{1}{N_0} \sum_{i=1}^N (1 - W_{ir}) z_{rqi}, \quad r = 1, \dots, N_{\text{rep}}, \\ \bar{\mathbf{z}}_{kT} &= (\bar{z}_{1kT}, \dots, \bar{z}_{N_{\text{rep}}kT})' \quad \text{and} \quad \bar{\mathbf{z}}_{kC} = ((\bar{z}_{1kC}, \dots, \bar{z}_{N_{\text{rep}}kC})'. \end{aligned}$$

The relative (compared to complete randomization) change in variance in the mean difference between the treated and controls in covariate/pretreatment outcome  $k$ , at time period  $T$ , using design  $d$  is then defined as

$$VC(k|d) \equiv \frac{\text{var}(\bar{\mathbf{z}}_{kT} - \bar{\mathbf{z}}_{kT}|d) - \text{var}(\bar{\mathbf{z}}_{kT} - \bar{\mathbf{z}}_{kT}|c)}{\text{var}(\bar{\mathbf{z}}_{kT} - \bar{\mathbf{z}}_{kT}|c)}, \quad (15)$$

for  $k = 1, \dots, K$ . The treatment effect estimate for each replicate,  $r$ , is defined

$$\begin{aligned} \hat{\tau}_{rc} &= \frac{1}{N_1} \sum_{i=1}^N W_{ir} Y_{iT+1} - \frac{1}{N_0} \sum_{i=1}^N (1 - W_{ir}) Y_{iT+1}, \\ \hat{\tau}_{dr} &= \frac{1}{N_1} \sum_{i=1}^N W_{ir} Y_{iT+1} - \frac{1}{N_0} \sum_{i=1}^N (1 - W_{ir}) Y_{iT+1} \Big| W \in \mathbf{W}_d^*. \end{aligned}$$

<sup>8</sup>The distribution of the FRT test is only known empirically (i.e., the histogram). Under the Fisher null (i.e., homogeneous treatment effects and the same variance of treated and controls and no treatment effect) the asymptotic variance is equal to  $\hat{v} = Ns^2/(N_1N_0)$  where  $s^2$  is the sample variance and

$$\frac{\hat{\tau} - 0}{\sqrt{\hat{v}}} \xrightarrow{d} N(0, 1).$$

Under these assumptions the  $t$ -test (i.e., Neyman–Pearson inference) and FRT have the same size asymptotically.

The empirical variance under rerandomization and complete randomization is then defined

$$\begin{aligned} \text{var}(\hat{\tau}_d) &= \frac{1}{N_{\text{rep}}} \sum_{r=1}^{N_{\text{rep}}} (\hat{\tau}_{rd} - \bar{\tau}_d)^2 \quad \text{and} \\ \text{var}(\hat{\tau}_c) &= \frac{1}{N_{\text{rep}}} \sum_{r=1}^{N_{\text{rep}}} (\hat{\tau}_{rc} - \bar{\tau}_c)^2, \end{aligned}$$

where  $\bar{\tau}_d$  and  $\bar{\tau}_c$  are the average estimated treatment effects across the replications. The relative change in variance of the treatment effect of strategy  $d$  is defined

$$VC_{\tau}(d) \equiv \frac{\text{var}(\hat{\tau}_d) - \text{var}(\hat{\tau}_c)}{\text{var}(\hat{\tau}_c)}. \quad (16)$$

The exact  $p$ -value for the complete randomization and the rerandomization strategies are defined as

$$\pi_{rc} = \Pr(|\hat{\tau}_{rc}(\mathbf{W}, \mathbf{Y})| \geq |\bar{\tau}_c|) \quad \text{and} \quad \pi_{rd} = \Pr(|\hat{\tau}_{rd}(\mathbf{W}_d^*, \mathbf{Y})| \geq |\bar{\tau}_d|).$$

Here,  $\hat{\tau}_{rc}(\mathbf{W}, \mathbf{Y})$  is the distribution of estimates over all allocations in replication  $r$  under complete randomization and  $\hat{\tau}_{rd}(\mathbf{W}_d^*, \mathbf{Y})$  is the corresponding distribution under rerandomization. The relative power is evaluated with  $\tau$  being varied from  $0.4\sigma_Y$  to  $2\sigma_Y$  in steps of  $0.4\sigma_Y$ , and estimated as

$$\text{Power}(\tau, d) \equiv \frac{p_{\tau d} - p_{\tau c}}{p_{\tau c}}, \quad \tau = 0.4, 0.8, \dots, 2.00, \quad (17)$$

where

$$\begin{aligned} p_{\tau c} &= \frac{1}{N_{\text{rep}}} \sum_{r=1}^{N_{\text{rep}}} \mathbb{1}(\pi_{rc} \leq 0.05) \quad \text{and} \\ p_{\tau d} &= \frac{1}{N_{\text{rep}}} \sum_{r=1}^{N_{\text{rep}}} \mathbb{1}(\pi_{rd} \leq 0.05). \end{aligned}$$

#### 5.1. Cross-Section Data and One Pretreatment Outcome

Consider the data-generating process (DGP)

$$Y_{it} = \mathbf{x}'_i \boldsymbol{\beta} + \epsilon_{it}, \quad i = 1, \dots, N, \quad t = 0, 1, \quad (18)$$

where  $\mathbf{x}_i$  is a  $K_x \times 1$  vector of normal distributed variables with mean 2 and covariance matrix  $\boldsymbol{\Sigma}$ , that is,  $\mathbf{x}_i \sim N(\mathbf{2}, \boldsymbol{\Sigma})$  and  $\epsilon_{it} = 0.3 \times \epsilon_{it-1} + \zeta_{it}$ , where  $\zeta_{it}$  is independent and identical distributed (iid) and normal, that is,  $\zeta_{it} \sim N(0, \sigma^2)$ .<sup>9</sup> Due to independence the marginal variance of the outcome is equal to

$$\text{var}(Y) = \boldsymbol{\beta} \boldsymbol{\Sigma} \boldsymbol{\beta}' + \frac{\sigma^2}{1 - 0.3^2}. \quad (19)$$

We compare the proposed rerandomization balance measure,  $R$ , with  $M$  for  $T = 0$  ( $K_y = 1, K_x = 3$ ) and  $T = 1$  ( $K_y = 1, K_x = 3$  and  $K_y = 1, K_x = 10$ ) under different specifications of  $\boldsymbol{\beta}$  and  $\boldsymbol{\Sigma}$ . The residual variance,  $\sigma^2$ , is chosen to obtain  $R^2 = 0.25$  in expectation.

<sup>9</sup>The results from the Monte Carlo simulations with regards to power are not sensitive with respect to the choice of distributions of the covariates or the error term. The covariates and the error terms are chosen to be normally distributed only to compare the results from the Monte Carlo to the theoretical expected variance reductions given in Equations (4) and (6) for small  $N$ .

**Table 1.** Relative variance change in the covariates and the effect estimate as compared to complete randomization (Equations (15) and (16)) for Mahalanobis distance ( $M$ ) and ranked  $p$ -values ( $R$ ) rerandomization designs.

Balance measure	A		B	
	$R$	$M$	$R$	$M$
$VC(X_1 .)$	-0.71	-0.74	-0.83	-0.74
$VC(X_2 .)$	-0.71	-0.75	-0.83	-0.75
$VC(X_3 .)$	-0.70	-0.75	-0.73	-0.75
$VC_{\tau}(\cdot)$	-0.20	-0.19	-0.22	-0.20

NOTE: The left and right panels display the results for DGPs A and B, respectively.

**5.1.1. Cross-Sectional Data**

This section serves to compare the proposed balance measure with the balance measure presented in Morgan and Rubin (2012) in a setting when the weights cannot be estimated from data. That is, the weights of the covariates are assumed uniform, that is, not estimated, which implies that we expect no improvement using our balance measure as compared to the Mahalanobis distance. As this setting is not primary focus, the results are restricted to  $N = 14$ , for which the 800 allocations with the (globally) smallest value of the balance measure are used. This implies that  $p_a = 800/\binom{14}{7} = 0.233$ .

Two DGPs are considered. In the first, (A), we let  $\beta = (1, 1, 1)$  and  $\Sigma = \text{diag}(1, 1, 1)$  and in the second, (B), we let  $\beta = (1, -1, 1)$  and

$$\Sigma = \begin{bmatrix} 1 & -0.8 & 0 \\ -0.8 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

Table 1 displays the relative variance reduction in the covariates and the estimated treatment effect. The variance reduction of the covariates for the  $M$  balance measure is, as expected, of the same magnitude for all covariates and around 75%. Given that  $v_a = 0.213 (= P(\chi_5^2 \leq 0.233)/P(\chi_3^2 \leq 0.233))$  this variance reduction is close to the one theoretically expected of 78.7% (cf. Equation (4)). The variance reductions of the covariates using the  $R$  balance measure is around 71% under DGP A. Under DGP B the variance reduction is 83% for the two correlated covariates but only 73% for the independent covariates. The variance reduction of the treatment effect under the null is around 20% for all strategies. Given that  $v_a = 0.213$

and  $R^2 = 0.25$  this is in line with the theoretically expected variance reduction using the  $M$  balance measure of 19.7% (see Equation (6)). Figure 3 displays the relative power gain of the rerandomization as compared to complete randomization for the two DGPs A and B. From the left panel, displaying the result under DGP A, one can see that both criteria increase the power by 20% for small treatment effects. Given that the covariates are uncorrelated, the similarity of the results with  $M$  and  $R$  with uniform weights is expected. The results displayed in panel B show that the criteria have similar power gains for small effect sizes. However, for larger effect sizes,  $R$  gives substantially larger power gains.

**5.1.2. Results With One Pretreatment Outcome**

Turning to the case with  $T = 1$ , two new DGPs are considered. The first, (A), with  $K_x = 3$  we let  $\beta = (0, 0.25, 0.75)$  and  $\Sigma = \text{diag}(1, 1, 1)$ . In the second, (B), with  $K_x = 10$  we let  $\beta = (0, 0, 0, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7)$  and  $\Sigma = \text{diag}(1, 1, 1, 1, 1, 1, 1, 1, 1, 1)$ . For both DGP's the importance of the covariates is increasing in  $k$ . In DGP A,  $x_1$  does not contribute to the variation of the outcome, and in DGP B,  $x_1, x_2,$  and  $x_3$  do not contribute to the variation of the outcome. These two DGPs illustrate how the criteria  $M, R, R(\omega, O),$  and  $R(\omega, L)$  are affected by the number of included covariates and how well the criteria incorporates the relative importance of the covariates. The considered sample sizes are  $N = 14, 50,$  and  $100$ . Including one pretreatment outcome observation implies  $K = 4$  and  $K = 11$  covariates. In the  $K = 11, N = 14$  setting, the  $M$  balance measure is expected to perform poorly, as the covariance matrix is very large in comparison to the sample size.

When  $N = 14$ , the 800 globally best allocations for a given sample are used for each criteria. As  $p_a = 0.233$  for the Mahalanobis balance measure this implies  $v_a = 0.29 (= P(\chi_6^2 \leq 0.233)/P(\chi_4^2 \leq 0.233))$  with  $Q = 4$ . For  $Q = 11$ , we get  $v_a = 0.51 (= P(\chi_{13}^2 \leq 0.233)/P(\chi_{11}^2 \leq 0.233))$ . This means that, using the  $M$  balance measure, the expected variance reductions in the covariates are 71% and 49% for  $K = 4$  and  $K = 11$ , respectively. The corresponding variance reduction in the variance of the effect estimate are 17.99% and 12.25%.



**Figure 3.** Relative power as compared to complete randomization for Mahalanobis distance ( $M$ ) and ranked  $p$ -values ( $R$ ) rerandomization designs. The left and right figures display the results for DGPs A and B, respectively.

**Table 2.** Relative change in variance of the covariates and in the effect estimate as compared to complete randomization (Equations (15) and (16)) for Mahalanobis distance ( $M$ ) and ranked  $p$ -values ( $R$ ) rerandomization designs.

N	$R(\omega, L)$			$R(\omega, O)$			$R$			$M$		
	14	50	100	14	50	100	14	50	100	14	50	100
$K = 4$												
$VC(Y_{\tau=T} .)$	-0.94	-0.96	-0.96	-0.92	-0.95	-0.95	-0.66	-0.71	-0.73	-0.64	-0.73	-0.74
$VC(X_1 .)$	-0.07	-0.09	-0.05	-0.25	-0.15	-0.09	-0.62	-0.66	-0.67	-0.65	-0.73	-0.74
$VC(X_2 .)$	-0.17	-0.16	-0.12	-0.33	-0.23	-0.15	-0.63	-0.69	-0.66	-0.67	-0.74	-0.73
$VC(X_3 .)$	-0.31	-0.55	-0.63	-0.51	-0.61	-0.65	-0.67	-0.71	-0.72	-0.63	-0.73	-0.73
$VC_{\tau}(.)$	-0.18	-0.20	-0.22	-0.20	-0.22	-0.24	-0.17	-0.20	-0.21	-0.16	-0.20	-0.20
$K = 11$												
$VC(Y_{\tau=T} .)$	-0.93	-0.96	-0.96	-0.68	-0.93	-0.94	-0.39	-0.50	-0.46	-0.21	-0.48	-0.47
$VC(X_1 .)$	-0.04	-0.08	-0.05	-0.24	-0.16	-0.09	-0.37	-0.43	-0.41	-0.21	-0.47	-0.49
$VC(X_2 .)$	-0.04	-0.07	-0.03	-0.31	-0.15	-0.13	-0.43	-0.42	-0.44	-0.25	-0.47	-0.49
$VC(X_3 .)$	-0.05	-0.04	-0.08	-0.25	-0.18	-0.12	-0.36	-0.43	-0.42	-0.21	-0.48	-0.50
$VC(X_4 .)$	0.03	-0.01	-0.07	-0.24	-0.11	-0.14	-0.36	-0.43	-0.45	-0.17	-0.46	-0.49
$VC(X_5 .)$	-0.04	-0.04	-0.08	-0.29	-0.17	-0.14	-0.37	-0.41	-0.43	-0.23	-0.46	-0.50
$VC(X_6 .)$	-0.02	-0.10	-0.07	-0.26	-0.18	-0.17	-0.40	-0.45	-0.42	-0.21	-0.47	-0.48
$VC(X_7 .)$	-0.08	-0.10	-0.16	-0.24	-0.24	-0.23	-0.36	-0.45	-0.43	-0.17	-0.50	-0.48
$VC(X_8 .)$	-0.13	-0.11	-0.18	-0.30	-0.28	-0.28	-0.41	-0.46	-0.42	-0.23	-0.46	-0.52
$VC(X_9 .)$	-0.11	-0.19	-0.31	-0.30	-0.35	-0.34	-0.39	-0.44	-0.44	-0.20	-0.49	-0.51
$VC(X_{10} .)$	-0.15	-0.22	-0.32	-0.33	-0.34	-0.40	-0.42	-0.44	-0.45	-0.23	-0.47	-0.52
$VC_{\tau}(.)$	-0.17	-0.16	-0.17	-0.15	-0.21	-0.22	-0.12	-0.12	-0.13	-0.05	-0.16	-0.11

NOTE: The top and bottom panels display the results including four ( $K = 4$ ) and eleven ( $K = 11$ ) covariates, respectively.

To limit the computational time for  $N = 50$  and  $N = 100$ , we randomly sampled 4000 from all possible allocations. To get the same resolution of the exact  $p$ -value as in the  $N = 14$  case, the rerandomization is performed by selecting the 800 best allocations according to the different criteria, within this set of 4000 allocations. This implies that the 20% best allocations among these 4000 are selected. Since  $p_a$  is a little bit smaller than with  $N = 14$  the expected variance reduction is a few percent larger. However, as a consequence of the initial random sampling of the 4000 allocations the theoretical results of the variance reductions for the  $M$  balance measure may be less exact for  $N = 50$  and 100.

Table 2 displays the change in variance in the covariates and the effect estimates for the four criteria compared to complete randomization for both DGP's. Figure 4 displays the results with respect to relative power of the four criteria.

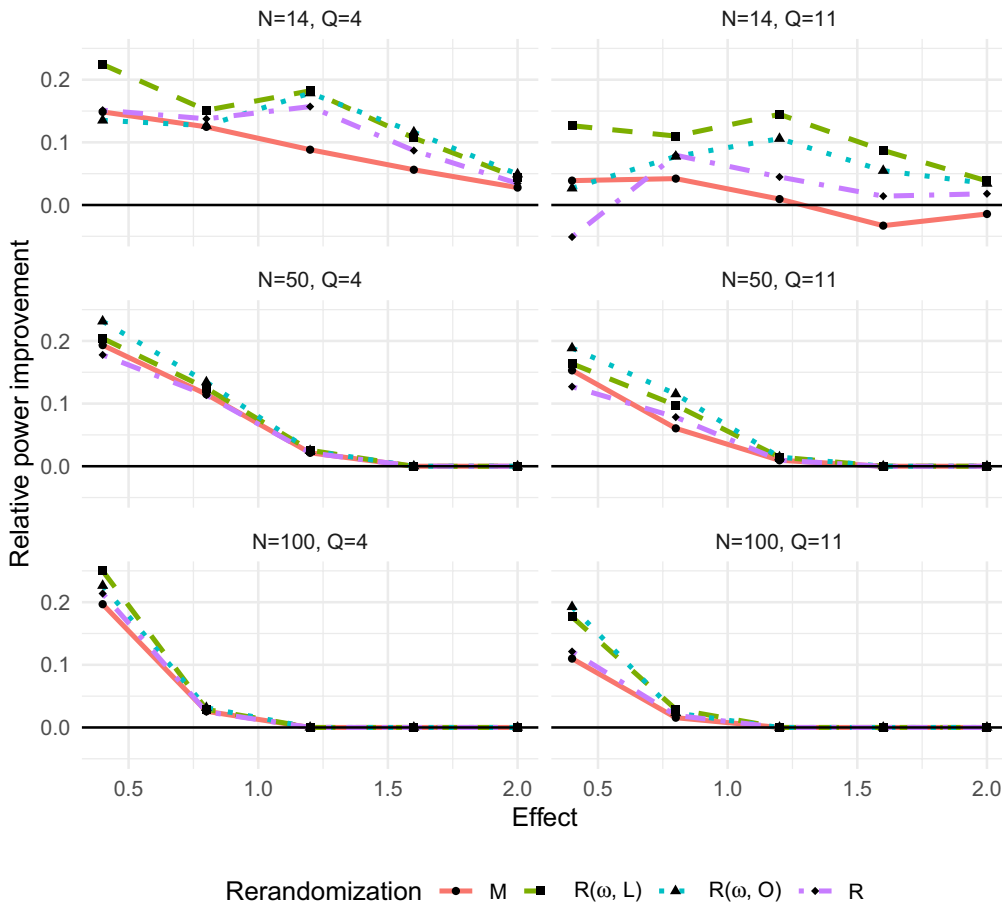
From the top panel, displaying the results with  $K = 4$ , we can see that the variance reduction of the  $R$  and  $M$  balance measure are very similar both in the covariates and the effect estimates. For  $N = 14$ , the variance reductions are a little bit lower than what is theoretically expected (71% in the covariates and 18% in the effect estimate). With  $N = 50$  and 100 the variance reduction is a little bit higher than what is theoretically expected. For the  $R(\omega, O)$  and  $R(\omega, L)$  criteria, the variance reductions in the covariates are the highest for the pretreatment outcome and monotonically increasing for the covariates which is in accordance with the increasing importance of the covariates in the DGPs. With regard to variance reduction of the estimated treatment effect, the  $R(\omega, O)$  balance measure is performing best overall.

From the bottom panel, displaying the results with  $K = 11$ , it is clear that the variance reductions of the  $M$  balance measure and the  $R$  balance measure are performing similar but only when  $N \geq 50$ . The variance reductions are close to the theoretically expected of 49% and 12% for the  $M$  balance measure. However

with  $N = 14$ , the variance reduction is much smaller than the theoretical value as expected given the small number of observations to estimate the  $10 \times 10$  covariance matrix. It is interesting to note that the variance reduction for the  $R$  balance measure with  $N = 14$  is on the same levels as with  $N = 50$  and 100. Turning to the two criteria using the estimated weights  $R(\omega, L)$  and  $R(\omega, O)$ , the pattern as with  $K = 4$  is repeated. Both criteria are reducing the variance of the lagged outcome the most and are both successfully picking up the relative importance of the covariates in the DGP, even though the  $R(\omega, L)$  balance measure are better. Except for  $N = 14$ , the  $R(\omega, O)$  balance measure is giving the largest variance reduction in the estimated treatment effect. With  $N = 14$ , the  $R(\omega, L)$  balance measure performs better than the  $R(\omega, O)$  balance measure.

The left panels of Figure 4 displays the results for the relative power when  $K = 4$ . From this figure, it is clear that all the rerandomization criteria increase power as compared to complete randomization. The power gain for the smallest effect sizes is around 15%–25% and increasing with  $N$ . Both the  $R(\omega, O)$  and  $R(\omega, L)$  criteria have higher power than the  $M$  and  $R$  criteria. With  $N = 14$ , the power gain using the  $R(\omega, L)$  balance measure is substantial also compared to the  $R(\omega, O)$  balance measure. Turning to the right panels displaying the results with  $K = 11$ , we find that the power is increasing in  $N$  and that the  $R(\omega, O)$  and  $R(\omega, L)$  criteria are more efficient than the  $M$  and  $R$  criteria. Once again, with  $N = 14$ , the power gain using the  $R(\omega, L)$  balance measure is substantial in comparison with the other criteria. This results displays the advantage of using penalizing in the situation of many covariates in combination with small samples. With  $N = 100$  the relative power increase is almost twice as good with the  $R(\omega, O)$  and  $R(\omega, L)$  criteria as compared to the  $M$  and  $R$  criteria.

To summarize, for both DGPs the power was substantially improved by using estimated weights. This result most likely stems from giving more weight to the covariates most correlated



**Figure 4.** Relative power as compared to complete randomization for Mahalanobis distance ( $M$ ) and ranked  $p$ -values ( $R$ ) rerandomization designs. The left and right panels display the results including four ( $K = 4$ ) and eleven ( $K = 11$ ) covariates, respectively.

with the outcome. Furthermore, with many covariates and small  $N$  it is important to penalize the number of covariates to be included in the criteria to obtain real power gains.

**5.2. Longitudinal Data**

In this section, two different time series DGP are considered with  $T = 10, 100$  and  $N = 1, 450, 100$ . In the first DGP the times series process, denoted homogeneous in the following, is set to be the same for all units

$$Y_{it} = 0.5 \times Y_{it-1} + \zeta_{it}, \quad i = 1, \dots, N, \quad t = 1, \dots, T + 1, \quad (20)$$

where  $\zeta_{it}$  iid  $N(0, 1)$ . In the second DGP, denoted heterogeneous in the following, the time series processes differ across four strata according to

$$Y_{it} = \phi_j Y_{i,\text{lag}} + \epsilon_{it}, \quad i = 1, \dots, N, \quad t = 1, \dots, T + 1, \quad j = 1, \dots, 4, \quad (21)$$

where  $Y_{i,\text{lag}} = (Y_{it-1}, Y_{it-2}, Y_{it-3})'$ ,  $\phi_1 = (0.5, 0, 0)$ ,  $\phi_2 = (0, 0.5, 0)$ ,  $\phi_3 = (0, 0, 0.5)$ ,  $\phi_4 = (0.39, 0.32, 0)$ , and  $\epsilon_{it}$  iid  $N(0, 1)$ . For both DGPs, the parameters are chosen such that  $R^2 = 0.25$ .

When  $T = 10$ , the  $R(\omega, O)$ ,  $R(\omega, L)$ , and  $R(\omega, F)$  criteria are used. In this context, the Mahalanobis-based rerandomization is difficult to apply due to singular covariance matrices and is therefore excluded. Given the large number of correlated

pretreatment outcomes with  $T = 100$ , also the  $R(\omega, O)$  balance measure runs into singularity problems in the estimation in that setting and is therefore excluded in that case.

In this setting with only outcome data, the OLS estimated weights are estimated as

$$\tilde{y}_{iT} = \beta_0 + \sum_{t=1}^{T-1} \beta_t \tilde{y}_{it} + \epsilon_i$$

using OLS. Denote  $\omega_\tau$  the weight for  $\tilde{y}_{i\tau}$  then  $\omega = (\omega_1, \dots, \omega_T)'$  and the estimated  $\hat{\omega}' = \left( \frac{|\hat{\beta}_1|}{\delta}, \dots, \frac{|\hat{\beta}_{T-1}|}{\delta}, \frac{1}{\delta} \right)$  where  $\delta = \sum_{t=1}^{T-1} |\hat{\beta}_t| + 1$ . This implies  $\omega_T \geq \omega_t \forall t = 1, \dots, T - 1$  in expectation.

**5.2.1. Results With  $T = 10$**

The top panel of Table 3 displays the results under the homogeneous DGP. With  $N \geq 50$ , all criteria are successful in giving the latter time periods larger weights which is in line with the DGP. The variance reduction for  $R(\omega, O)$  varies more across different  $N$ , than the two other criteria. The advantage of using penalized regression for small  $N$  is confirmed also in this setting. It is clear that the  $R(\omega, L)$  and  $R(\omega, O)$  criteria give larger reductions in the variance of the lagged outcomes and the estimated treatment effect under the null than  $R(\omega, F)$ . The variance reduction in the estimated effects is significantly better for the  $R(\omega, L)$  and  $R(\omega, O)$  criteria than for  $R(\omega, F)$  balance measure. The bottom



**Table 3.** Variance reduction of the covariates and in the effect estimate as compared to complete randomization (Equations (15) and (16)) for the ranked  $p$ -values ( $R$ ) rerandomization designs.

Rerandomization	$R(\omega, F)$			$R(\omega, L)$			$R(\omega, O)$		
	14	50	100	14	50	100	14	50	100
DPG = Homogeneous									
VC(Y1 <sub> .</sub> )	-0.01	0.04	0.01	-0.07	-0.02	-0.01	-0.33	-0.17	-0.12
VC(Y2 <sub> .</sub> )	-0.01	-0.02	-0.00	-0.09	-0.04	-0.05	-0.37	-0.25	-0.18
VC(Y3 <sub> .</sub> )	-0.00	0.01	-0.04	-0.08	-0.03	-0.07	-0.39	-0.24	-0.25
VC(Y4 <sub> .</sub> )	-0.01	-0.02	-0.06	-0.05	-0.03	-0.07	-0.42	-0.26	-0.20
VC(Y5 <sub> .</sub> )	-0.08	-0.09	-0.08	-0.06	-0.11	-0.07	-0.42	-0.30	-0.16
VC(Y6 <sub> .</sub> )	-0.12	-0.08	-0.09	-0.10	-0.10	-0.11	-0.41	-0.29	-0.22
VC(Y7 <sub> .</sub> )	-0.17	-0.14	-0.16	-0.08	-0.08	-0.10	-0.39	-0.26	-0.25
VC(Y8 <sub> .</sub> )	-0.28	-0.23	-0.25	-0.20	-0.24	-0.21	-0.45	-0.34	-0.32
VC(Y9 <sub> .</sub> )	-0.37	-0.31	-0.32	-0.42	-0.62	-0.69	-0.51	-0.68	-0.73
VC(Y10 <sub> .</sub> )	-0.39	-0.42	-0.39	-0.92	-0.95	-0.95	-0.69	-0.92	-0.93
VC $_{\tau}$ (.)	-0.10	-0.10	-0.13	-0.19	-0.27	-0.21	-0.14	-0.25	-0.24
DPG = Heterogeneous									
VC(Y1 <sub> .</sub> )	0.01	0.00	-0.03	-0.01	-0.03	-0.08	-0.29	-0.18	-0.17
VC(Y2 <sub> .</sub> )	-0.03	-0.03	0.02	-0.10	-0.07	0.03	-0.33	-0.23	-0.10
VC(Y3 <sub> .</sub> )	-0.05	-0.01	-0.11	-0.05	-0.01	-0.06	-0.30	-0.18	-0.17
VC(Y4 <sub> .</sub> )	-0.08	-0.06	-0.05	-0.10	-0.05	-0.06	-0.35	-0.22	-0.18
VC(Y5 <sub> .</sub> )	-0.10	-0.11	-0.08	-0.06	-0.04	-0.05	-0.35	-0.19	-0.17
VC(Y6 <sub> .</sub> )	-0.14	-0.10	-0.15	-0.06	-0.05	-0.11	-0.36	-0.22	-0.20
VC(Y7 <sub> .</sub> )	-0.19	-0.19	-0.21	-0.08	-0.15	-0.24	-0.36	-0.28	-0.28
VC(Y8 <sub> .</sub> )	-0.19	-0.21	-0.16	-0.14	-0.25	-0.29	-0.38	-0.39	-0.34
VC(Y9 <sub> .</sub> )	-0.29	-0.27	-0.24	-0.14	-0.09	-0.13	-0.35	-0.24	-0.20
VC(Y10 <sub> .</sub> )	-0.31	-0.27	-0.30	-0.93	-0.96	-0.97	-0.71	-0.94	-0.95
VC $_{\tau}$ (.)	-0.06	-0.08	-0.10	-0.04	-0.07	-0.05	-0.06	-0.09	-0.06

NOTE: The top and bottom panels display the results for the homogeneous and heterogeneous DGPs, respectively.

panel of Table 3 displays the results with the heterogeneous DGP. Also in this case, the variance reduction increases with  $t$  when  $N \geq 50$  for all criteria. Once again we see that the variance reduction is of similar magnitudes for the  $R(\omega, F)$  and  $R(\omega, L)$  criteria but that the variance reduction of  $R(\omega, O)$  with  $N = 14$  differs to a large extent from the variance reduction with larger  $N$ . The variance reduction in the effect estimate is of similar magnitude for  $N = 14$  and 50. With  $N = 100$ , the  $R(\omega, F)$  balance measure gives the largest variance reduction.

Figure 5 displays the relative power of the different rerandomization strategies as opposed to complete randomization under the homogeneous (left panel) and heterogeneous (right panel) DGPs. From the left panel one can see that for small effects the power gain is around 20% for the  $R(\omega, O)$  and  $R(\omega, L)$  criteria for all  $N$ . Furthermore, these criteria are superior to the  $R(\omega, F)$  balance measure. From the right hand panels one can see that the power gains for  $N = 50$  and 100 is around 10% for the  $R(\omega, F)$  balance measure. For these sample sizes, this balance measure gives almost 100% larger relative improvement than the two other criteria. With  $N = 14$  it is hard to get any improvements for any of these criteria in comparison to the complete randomization, at least with  $R^2 \leq 0.25$ .

### 5.2.2. Results With $T = 100$

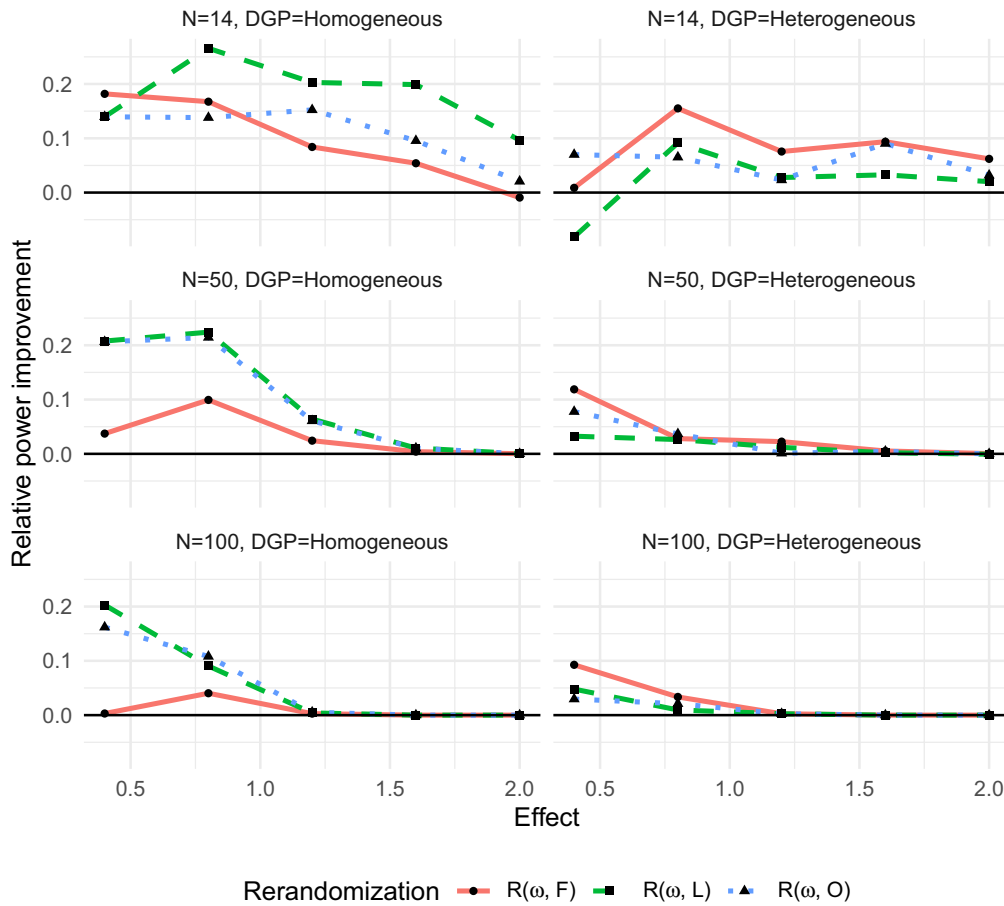
Table 4 displays the change in variance in the effect estimate for the two criteria. With the homogeneous DGP, the  $R(\omega, F)$  and  $R(\omega, L)$  criteria give similar variance reductions in the range 19%–25%. With the heterogeneous DGP the variance reduction obtained with the  $R(\omega, F)$  balance measure is of the same magnitude as with the homogeneous DGP. The variance reduction using the  $R(\omega, L)$  balance measure is now only around 10%. The variance reduction in the 100 lags shows a pattern (not

displayed) very similar to the pattern presented in Table 3. The first 90 time periods have close to zero weights for both strategies and the last 10 have increasing weights. Figure 6 displays the relative power of the  $R(\omega, F)$  and  $R(\omega, L)$  criteria under the two DGPs. In the left panels, displaying the results from the homogeneous DGP we can see that with the smallest effects size the increase in power is around 10% with  $N = 14$  and around 20% with  $N = 50$  and 100. For  $N = 50$  and 100, the  $R(\omega, L)$  balance measure performs better than the  $R(\omega, F)$  balance measure. From the right hand panels, displaying the results with the heterogeneous DGP, we can see that the  $R(\omega, F)$  balance measure with  $N = 50$  and 100 is in the range 20%–15% in comparison to the complete randomization. With  $N = 14$  the increase in power is only 5%. The  $R(\omega, L)$  balance measure gives hardly any improvement in power in comparison to the complete randomization.

In summary of Section 5.2, when the number of time periods and sample size is small, and the DGPs are heterogeneous, it is difficult to improve the power by the rerandomization strategies presented in this section. However, if either the sample size or the number of time periods increase, there are gains to be made using the strategies presented here. With long time series of pretreatment outcomes there are large gains for both considered DGPs. Of course, with large  $T$ , the level of heterogeneity can be evaluated by preanalysis to guide the choice of strategy.

## 6. Empirical Example—An Information Experiment on Electricity Consumption

This section illustrates how the proposed design strategies can be applied in practice in a small sample experiment. The



**Figure 5.** Relative power as compared to complete randomization for the ranked  $p$ -values ( $R$ ) rerandomization designs for  $T = 10$ . The left panel and right panel display the results from the Homogenous and the heterogeneous DGPs, respectively.

**Table 4.** Relative change in variance in the estimated treatment effect as compared to complete randomization (Equations (15) and (16)) for the ranked  $p$ -values ( $R$ ) rerandomization designs with forecast-based and LASSO estimated weights, respectively.

Rerandomization	$R(\omega, F)$			$R(\omega, L)$		
	14	50	100	14	50	100
Homogeneous						
$VC_{\tau}(\cdot)$	-0.19	-0.25	-0.21	-0.22	-0.25	-0.24
Heterogeneous						
$VC_{\tau}(\cdot)$	-0.23	-0.26	-0.22	-0.07	-0.11	-0.06

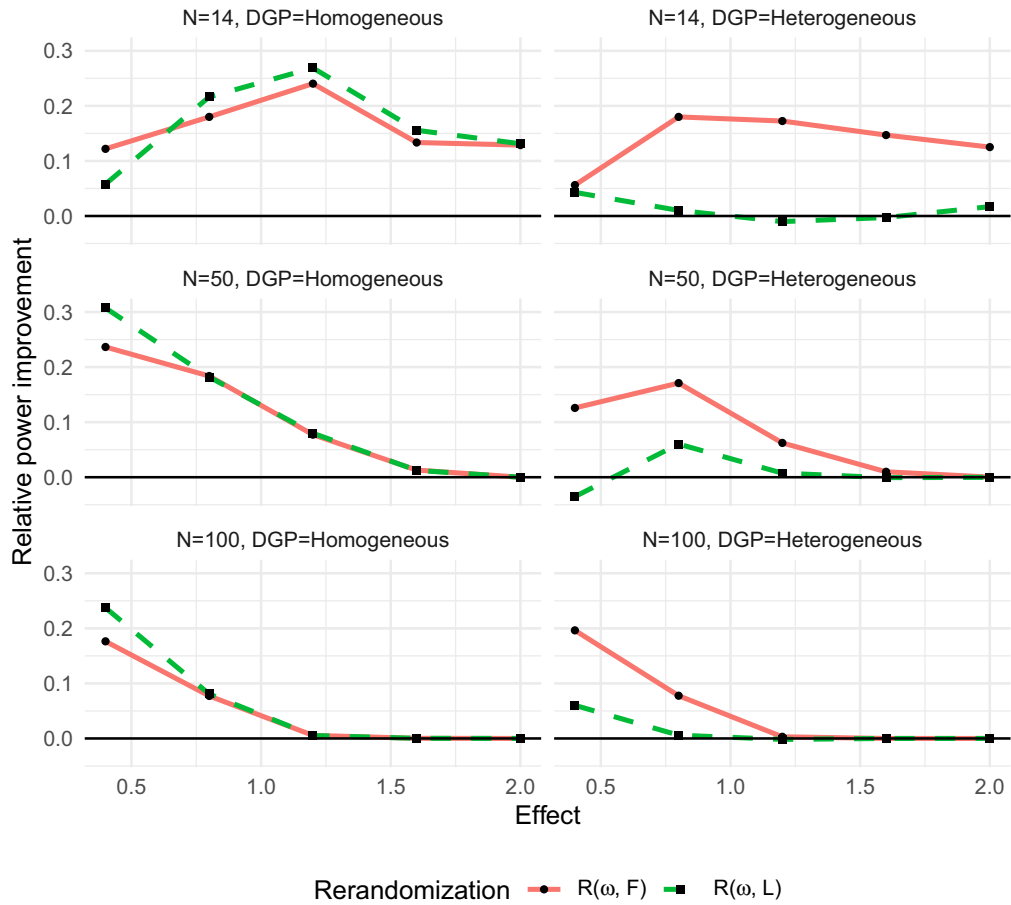
NOTE: The top and bottom panels display the results for the homogeneous and heterogeneous DGPs, respectively.

pre-experimental data used in this example originates from Öhrlund, Stikvoort, and Schultzberg (2018), where a small randomized experiment is conducted. The interest is in how electricity consumption behavior can be affected toward a behavior more suitable for solar energy by providing information on energy consumption. Several outcomes are of interest in the original study, here we focus on one of them, the mean consumption difference between the group that received the information and the group that did not. The sample consists of 54 households for which the electricity consumption is observed for each hour for the 4 months before treatment assignment. To increase the efficiency, a complex stratified randomization

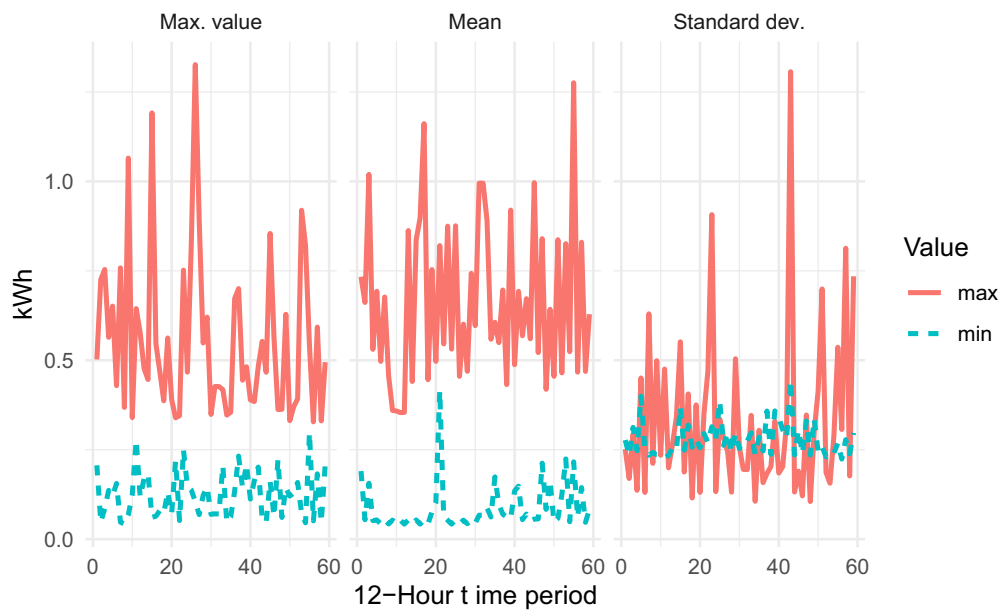
was used to assign treatment in the original study using all pretreatment data. Here, as an illustration, only the first 30 days are used, where consumption data are aggregated to 60 distinct 12-hr periods. The last period ( $T + 1 = 60$ ) is left out from the design stage and is used to evaluate the design.

All the pretreatment outcome time series are presented in Figure 1. Figure 7 displays the heterogeneity in the pretreatment outcome by showing the households with the smallest and largest maximum consumption value, the smallest and largest household mean, and the smallest and largest standard deviation over the pretreatment time periods, in the panels from left to right, respectively. It is clear that there are quite large differences in several aspects of the electricity consumption between the households during the pretreatment period and it is clearly not trivial to find a balanced design.

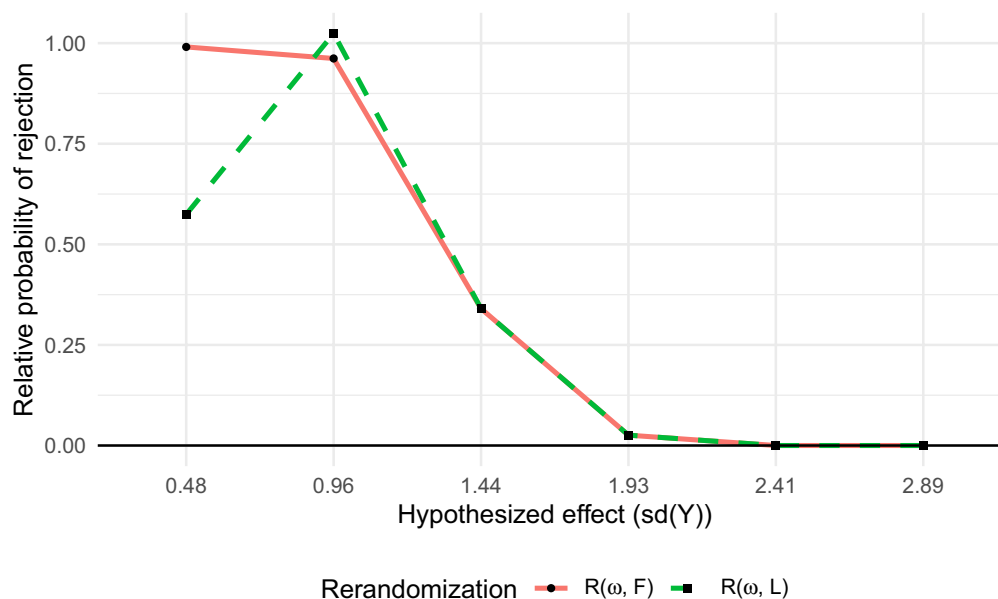
Since the pretreatment data are measured with high frequency and no other covariates are available, the two strategies presented in the latter part of Section 5.2 are used, that is select the best allocations according the  $R(\omega, L)$  and  $R(\omega, F)$  criteria. Since the number of possible allocations equals  $\binom{54}{27} = 1.95e15$ , the globally best allocations cannot be found and instead the procedure presented in Section 3 is applied. We chose here a resolution of  $1/400$  implying  $H = 800$ , that is, allocations were sampled randomly without replacement and the best 800 was kept. The procedure was left working for 11 hr which in this case meant that a random sample of one billion allocations were



**Figure 6.** Relative power as compared to complete randomization the ranked  $p$ -values ( $R$ ) rerandomization designs with forecast-based and LASSO estimated weights, respectively, for  $T = 100$ . The left panel and right panel display the results from the homogenous and the heterogeneous DGP's, respectively.



**Figure 7.** The electricity consumption (kWh) for the households with the largest (max) and smallest (min) maximum consumption, mean consumption, and standard deviation in their consumption, respectively.



**Figure 8.** The relative probability as compared to complete randomization of randomly selecting an allocation that gives a significant result for two different rerandomization strategies given different hypothesized treatment effects.

considered. As a benchmark to the rerandomization strategies, complete randomization was conducted. The exact  $p$ -value for the complete randomization was Monte Carlo approximated by a random draw of 40,000 allocations from the considered billion.

To evaluate the potential power gains under the proposed designs as compared to complete randomization, hypothetical homogeneous treatment effects were added to the treated group. That is, for each of the 800 selected allocations, the hypothetical treatment effect was added to  $Y_{60}$  for the treatment group and the exact  $p$ -value calculated. The same procedure was applied to the 40,000 randomly selected allocations (complete randomization), and the relative number of allocations that had exact  $p$ -values less than  $\alpha = 0.05$  were compared by calculating the relative probability of drawing an allocation that rejects the null under the alternative when using  $R(\omega, L)$  or  $R(\omega, F)$  as compared to complete randomization.

Figure 8 displays the relative probability of rejecting the null under the alternative when random assignment is restricted to the set of allocation defined by the two criteria as compared to random assignment in the set of 40,000 randomly selected allocations.<sup>10</sup> From the figure we can see that for a half standard deviation (of the outcome) effect, the probability of obtaining a statistically significant effect after making a random draw from the 800 allocations determined by the  $R(\omega, F)$  and  $R(\omega, L)$  measures is around 99% and 57% higher, respectively, than when making a random draw from the 40,000 randomly chosen allocations. It is worth noting that the forecast procedure, that is to use  $R(\omega, F)$  balance measure is less computationally demanding than the  $R(\omega, L)$  balance measure. The  $R(\omega, L)$  takes around  $T$  time longer than the  $R(\omega, F)$  balance measure to calculate.

An alternative way of displaying the difference between the complete randomization and the rerandomization strategies is

**Table 5.** Relative variance change in the effect estimates across the possible allocation using  $R(\omega, F)$  and  $R(\omega, L)$  as compared to complete randomization (Equation (16)).

Rerandomization	$R(\omega, F)$	$R(\omega, L)$
$VC_{\tau}(\cdot)$	-0.56	-0.61

NOTE: Complete randomization is here Monte Carlo approximated by 40,000 random allocations.

to look at the empirical variance of the effect estimate under these designs. The variance of the effect estimate is thus obtained by estimating the effect for the restricted set of 800 allocations under the two criteria and for all 40,000 allocations for the complete randomization. Table 5 displays the percentage change in variance in the effect estimate compared to complete randomization. It is clear that, in comparison to complete randomization the variance using the  $R(\omega, F)$  and  $R(\omega, L)$  criteria is reduced by 56% and 61%, respectively.

## 7. Discussion

Based on the results in Morgan and Rubin (2012), this article develops strategies for rerandomization as a means to increase efficiency in randomized experiments. Morgan and Rubin (2012) suggested randomization based on the Mahalanobis distance balance measure of the covariate mean-difference vector between potential treated and controls. This article has two main contributions. First, a strategy to sample from set of admissible allocations fulfilling an implicit rerandomization balance measure to find the best possible design, given a balance measure, within a certain time limit. With the proposed sampling strategy and a symmetric balance measure, the difference-in-means estimator is unbiased and the Fisher test is exact by design. Second, a new covariate balance measure is proposed as an alternative to the Mahalanobis distance. The balance measure differs from the Mahalanobis distance in several ways;

<sup>10</sup>Note that there is no difference in rejection rate under the null as all tests are based on exact inference.



it has computational advantages over the Mahalanobis balance measure in the situation of a large set of highly correlated covariates, and importantly, as the proposed balance measure expresses the weights of each covariate explicitly, it enables for various strategies of estimating the weights from data. For given a priori weights, the strategy can be considered an alternative to the strategy of Morgan and Rubin (2015) whom suggested rerandomization within tiers of importance. The proposed criterion balance measure is especially useful with one pretreatment outcome or longitudinal pretreatment outcome data. In this situation, the correlation structure of the data can be estimated using data to give different weights to the covariates and the pretreatment outcomes accordingly.

The Monte Carlo simulations show: (i) that with traditional cross-section data (i.e., only covariates) the suggested criterion has similar performance as the Mahalanobis criterion, (ii) an advantage with the new strategy to the Mahalanobis criterion when one or several pretreatment outcomes are available. Finally (iii), two Monte Carlo simulations with only pretreatments observations (as in the empirical illustration) shows advantages with the new strategies in comparison to complete randomization.

Taking use of a sample of 54 households with electricity consumption over 60 time periods, it is shown that the power of a mean difference test in a balanced randomized experiment can be increased by up to 100% using one of the proposed rerandomization strategies as compared to complete randomization.

## Supplementary Materials

The file “R-functions.R” contains a function for performing the designs proposed in the article. Moreover, examples of several designs are given with simulated data. Code for a small Monte Carlo simulation study is given for comparison of different designs. The code for the design of the Empirical example is given (Note that smaller number of considered allocations is specified so that the whole code will run in reasonable time). Finally, a function for calculating the CDF for the  $H$ :th order statistic of the Mahalanobis distance in the set of considered allocations is given together with the example from the paper.

The file “EXACT\_INFERENCE\_RCPP\_ARMADILLO.cpp” contains C++ functions for calculations of the exact  $p$ -values needed in the designs. This file is sourced from the “R-functions.R” and is necessary for running the function therein.

The file “EMPIRICAL\_DATA\_ELECTRICITY.Rdata” contains the data used in the empirical example. This file is loaded in “R-functions.R”.

## Acknowledgments

The authors thank Rauf Ahmad, Bengt Muthén, and Mattias Nordin, seminar participants at the statistics department Uppsala University, the audience at the The First Beijing Symposium in Biostatistics and Data Science, November 16–17, 2018, two anonymous reviewers, and the associate editor for helpful suggestions.

## References

- Fisher, R. A. (1935), *The Design of Experiments*, Edinburgh: Oliver and Boyd. [799,800]
- Gut, A. (2009), *An Intermediate Course in Probability*, Springer Texts in Statistics, New York: Springer. [802]
- Hamaker, E. L., Asparouhov, T., Brose, A., Schmiedek, F., and Muthén, B. (2018), “At the Frontiers of Modeling Intensive Longitudinal Data: Dynamic Structural Equation Models for the Affective Measurements From the COGITO Study,” *Multivariate Behavioral Research*, 53, 820–841. [799]
- Hamaker, E. L., and Wichers, M. (2017), “No Time Like the Present,” *Current Directions in Psychological Science*, 26, 10–15. [799]
- Hu, Y., and Hu, F. (2012), “Balancing Treatment Allocation Over Continuous Covariates: A New Imbalance Measure for Minimization,” *Journal of Probability and Statistics*, 2012, 842369. [800]
- Imbens, G. W., and Rubin, D. B. (2015), *Causal Inference in Statistics, Social, and Biomedical Sciences: An Introduction*, Cambridge: Cambridge University Press. [798]
- Li, X., Ding, P., and Rubin, D. B. (2018), “Asymptotic Theory of Rerandomization in Treatment–Control Experiments,” *Proceedings of the National Academy of Sciences of the United States of America*, 115, 9157–9162. [798,800]
- Morgan, K. L., and Rubin, D. B. (2012), “Rerandomization to Improve Covariate Balance in Experiments,” *Annals of Statistics*, 40, 1263–1282. [798,799,800,802,803,806,812]
- Morgan, K. L., and Rubin, D. B. (2015), “Rerandomization to Balance Tiers of Covariates,” *Journal of the American Statistical Association*, 110, 1412–1421. [798,800,803,813]
- Öhrlund, I., Stikvoort, B., and Schultzberg, M. (2018), “Taking the Sun in Our on Hands: Implications of Shared Residential Solar Installations on Pre-Environmental Behaviours,” Mimeo. [810]
- Tibshirani, R. (1996), “Regression Selection and Shrinkage via the Lasso,” *Journal of the Royal Statistical Society, Series B*, 58, 267–288. [804]
- Zhou, Q., Ernst, P. A., Morgan, K. L., Rubin, D. B., and Zhang, A. (2018), “Sequential Rerandomization,” *Biometrika*, 105, 745–752. [798,800]