

# Sequential rerandomization

BY QUAN ZHOU, PHILIP A. ERNST

*Department of Statistics, Rice University, 6100 Main St., Houston, Texas 77005, U.S.A.*  
quan.zhou@rice.edu philip.ernst@rice.edu

KARI LOCK MORGAN

*Department of Statistics, Pennsylvania State University, 323 Thomas Building, University Park, Pennsylvania 16801, U.S.A.*  
klm47@psu.edu

DONALD B. RUBIN

*Department of Statistics, Harvard University, One Oxford Street, Cambridge, Massachusetts 02138, U.S.A.*  
rubin@stat.harvard.edu

AND ANRU ZHANG

*Department of Statistics, University of Wisconsin-Madison, 1300 University Ave., Madison, Wisconsin 53706, U.S.A.*  
anruzhang@stat.wisc.edu

## SUMMARY

The seminal work of Morgan & Rubin (2012) considers rerandomization for all the units at one time. In practice, however, experimenters may have to rerandomize units sequentially. For example, a clinician studying a rare disease may be unable to wait to perform an experiment until all the experimental units are recruited. Our work offers a mathematical framework for sequential rerandomization designs, where the experimental units are enrolled in groups. We formulate an adaptive rerandomization procedure for balancing treatment/control assignments over some continuous or binary covariates, using Mahalanobis distance as the imbalance measure. We prove in our key result that given the same number of rerandomizations, in expected value, under certain mild assumptions, sequential rerandomization achieves better covariate balance than rerandomization at one time.

*Some key words:* Experimental design; Mahalanobis distance; Noncentral chi-squared distribution; Sequential enrolment.

## 1. INTRODUCTION

Rerandomization is a method for achieving balanced distributions of covariates across treatment groups before conducting an experiment (Holschuh, 1980; Urbach, 1985; Imai et al., 2008; Morgan & Rubin, 2012). Despite advocacy for rerandomization dating back to Sir Ronald Fisher (Savage, 1962, p. 88), a solid mathematical foundation for rerandomization was only recently developed in the seminal work of Morgan & Rubin (2012), which advises rerandomization only if ‘the decision to rerandomize or not is based on a pre-specified criterion’ (Morgan & Rubin, 2012, p. 1265). This work has catalysed a surge of research in rerandomization, both theoretical and applied in nature. For theoretical contributions, see Li & Ding (2017) and Morgan & Rubin (2015). For more applied contributions, see Athey & Imbens (2017), Delavande et al. (2016) and Xu & Kalbfleisch (2013).

The main objective of this work is to balance treatment/control assignments over some continuous, or binary, covariates by rerandomization. The majority of the traditional randomization procedures were developed for discrete covariates only, and continuous covariates are simply discretized by binning. However, both the number and the boundaries of such bins are very difficult to choose, as discussed in [Hu & Hu \(2012\)](#). [Morgan & Rubin \(2012\)](#) considered rerandomization for a finite sample where all units were recruited at one time, using Mahalanobis distance as the imbalance measure; henceforth we refer to this as Morgan–Rubin complete rerandomization. The theoretical advantages of using Mahalanobis distance for continuous covariates are discussed in [Rubin \(1979\)](#) and [Greevy et al. \(2004\)](#). When the data contain categorical covariates, as advocated by [Morgan & Rubin \(2012\)](#), one may combine blocking with rerandomization by applying a stratified randomization procedure to the most important categorical covariates. In practice, however, a researcher may be unable to wait to perform an experiment until all experimental units can be recruited, and thus covariate-adaptive minimization methods (see [Lin et al., 2015](#)) might be preferred. To solve this problem, in the present work we consider rerandomization for sequential enrolment designs where participants arrive in groups, which we henceforth term sequential rerandomization. To the best of our knowledge, a mathematical framework for sequential rerandomization has not been developed previously. A unique advantage of sequential rerandomization is that while it is adaptive, it still allows for rerandomization, and thus is much less liable to selection bias than are minimization procedures ([Berger, 2010](#)). For more discussion on the relationship between rerandomization and other methods, such as the finite selection model, see [Morgan & Rubin \(2012, § 5\)](#).

Given the same number of rerandomizations, in expected value, a seemingly natural conjecture is that the balance created by employing Morgan–Rubin complete rerandomization would, in expectation, be superior to that created using sequential rerandomization, since Morgan–Rubin complete rerandomization allows for all possible allocations of units. Under only mild asymptotic conditions, our main result in this paper, Theorem 3, shows the opposite to be true; see § 4. The key mathematical implications for sequential rerandomization and the results needed to prove Theorem 3 are provided in § 2 and § 3. In § 5 we extend our results to more general settings and conclude with a discussion on optimal randomization procedures. All proofs and the results of simulation studies are given in the [Supplementary Material](#).

## 2. SEQUENTIAL RERANDOMIZATION

Consider a sequential trial in which  $2N$  units are to be divided into  $K$  sequential groups, each group containing  $2n_1, \dots, 2n_K$  experimental units, where  $n_1 + \dots + n_K = N$ . Let the matrix  $X = (X_1, \dots, X_K) \in \mathbb{R}^{p \times (2N)}$  represent the  $p$  covariates for these  $2N$  units, where  $X_1, \dots, X_K$  are block matrices with corresponding dimensions  $p \times 2n_1, \dots, p \times 2n_K$ ; assume that  $X_1, \dots, X_K$  are observed sequentially. The matrix  $X$  will be treated as fixed, and the sample covariance matrix of the  $k$ th group, which is denoted by  $\text{cov}(X_k)$  and has dimension  $p \times p$ , is assumed to have rank equal to  $p$ .

Consider the following rerandomization procedure. For the first group of  $2n_1$  units, we randomly assign  $n_1$  patients to the treatment group and the other  $n_1$  to the control group. We denote this randomization by  $W_1^* = (W_{1,1}^*, \dots, W_{1,2n_1}^*)^T$ , a vector of dimension  $2n_1$ , where  $W_{1,i}^* = 1$  if the  $i$ th patient of the first group is assigned to treatment and  $W_{1,i}^* = 0$  otherwise. Throughout this paper, a superscript  $*$  will denote results from a tentative allocation, subject to being accepted or rerandomized based on a specific criterion, whereas results without a superscript  $*$  correspond to the actual treatment administered. The Mahalanobis distance between the treatment and control groups corresponding to  $W_1^*$  is

$$M_1^* = \frac{n_1}{2} (\bar{X}_{T,1}^* - \bar{X}_{C,1}^*)^T \text{cov}(X_1)^{-1} (\bar{X}_{T,1}^* - \bar{X}_{C,1}^*),$$

where  $\bar{X}_{T,1}^* = n_1^{-1} X_1 W_1^*$  and  $\bar{X}_{C,1}^* = n_1^{-1} X_1 (1 - W_1^*)$  are the  $p$ -dimensional mean vectors of the treatment, T, and control, C, groups respectively. This expression is based on the observation that  $\text{cov}(\bar{X}_{T,1}^* - \bar{X}_{C,1}^* | X_1) = 2 \text{cov}(X_1)/n_1$ ; see the [Supplementary Material](#) for details. As in [Morgan & Rubin \(2012\)](#), we let  $(\varphi_1, a_1)$  represent a prespecified rerandomization criterion such that  $\varphi_1(X_1, W_1^*) = 1$  if  $M_1^* < a_1$  and 0 otherwise, where  $\varphi_1 = 1$  indicates an acceptable rerandomization. If  $\varphi_1 = 0$ , then  $W_1^*$  is not acceptable

and the randomization is repeated; otherwise we set  $W_1 = W_1^*$ ,  $M_1 = M_1^*$ ,  $\bar{X}_{T,1} = \bar{X}_{T,1}^*$  and  $\bar{X}_{C,1} = \bar{X}_{C,1}^*$  and proceed to consider the second group of  $2n_2$  units. If  $K = 1$ , we simply stop and sequential rerandomization reduces to Morgan–Rubin complete rerandomization.

The above procedure continues as follows. For the  $k$ th group of units, we randomize  $n_k$  units to treatment and  $n_k$  units to control and denote the tentative assignment by  $W_k^*$ . It should be emphasized that sequential rerandomization takes into account all the data and fixed assignments from the first  $k - 1$  groups, namely  $X_{1:(k-1)} = (X_1, \dots, X_{k-1})$  and  $W_{1:(k-1)} = (W_1^T, \dots, W_{k-1}^T)^T$ , in addition to the data from the  $k$ th group. The total number of subjects used to assess the acceptability of  $W_k^*$  for the  $k$ th group is  $2n_{1:k}$ , where  $n_{1:k} = \sum_{j=1}^k n_j$ . The assignment of the first  $k$  groups using  $W_k^*$  is denoted by

$$W_{1:k}^* = (W_1^T, \dots, W_{k-1}^T, W_k^{*T})^T,$$

which is a vector with  $2n_{1:k}$  components. The superscript  $*$  on the right-hand side occurs only at the  $k$ th term because the assignment vectors of the first  $k - 1$  groups are already fixed. The mean vectors of the first  $k$  treatment and control groups are written as

$$\bar{X}_{T,1:k}^* = \frac{1}{n_{1:k}} X_{1:k} W_{1:k}^*, \quad \bar{X}_{C,1:k}^* = \frac{1}{n_{1:k}} X_{1:k} (1 - W_{1:k}^*),$$

with corresponding Mahalanobis distance for the first  $k$  groups

$$M_k^* = \frac{n_{1:k}}{2} (\bar{X}_{T,1:k}^* - \bar{X}_{C,1:k}^*)^T \text{cov}(X_{1:k})^{-1} (\bar{X}_{T,1:k}^* - \bar{X}_{C,1:k}^*), \quad (1)$$

where  $\text{cov}(X_{1:k})$  is the sample covariance matrix of  $X_{1:k}$ , which is assumed to be of full rank. Given  $a_k$ , we decide whether  $W_k^*$  is acceptable by evaluating the prespecified rerandomization criterion

$$\varphi_k(X_{1:k}, W_{1:k}^*) = \begin{cases} 1, & M_k^* < a_k, \\ 0, & \text{otherwise,} \end{cases} \quad (2)$$

where  $n_1, \dots, n_K$  must be large enough to ensure that an acceptable randomization can be realized. The threshold  $a_k$  can be chosen as a function of  $M_1, \dots, M_{k-1}$ , but as we will see shortly in § 3.2,  $M_{k-1}$  alone is sufficient for choosing  $a_k$ . After the experimenter has concluded the sequential allocations, the Mahalanobis distance is calculated on the complete dataset  $X = X_{1:K}$  using the appropriate version of (1).

### 3. PROPERTIES OF SEQUENTIAL RERANDOMIZATION

#### 3.1. Average treatment effect estimation

Now we present the key mathematical consequences of the sequential rerandomization framework outlined in § 2. We begin with the estimation of the true average treatment effect for the entire sample. Suppose that the potential outcome for unit  $i$  after treatment or control is  $y_i(1)$  or  $y_i(0)$ , respectively, according to the Rubin causal model (Rubin, 1974). Let the observed response be given by  $Y_i = y_i(1)$  if  $W_i = 1$  and  $Y_i = y_i(0)$  otherwise. The average treatment effect is

$$\tau = \frac{\sum_{i=1}^{2N} y_i(1) - \sum_{i=1}^{2N} y_i(0)}{2N}.$$

The usual estimate of  $\tau$  is the difference between the treatment group and control group sample means:

$$\hat{\tau} = \bar{Y}_T - \bar{Y}_C = \frac{1}{N} \sum_{i=1}^{2N} Y_i W_i - \frac{1}{N} \sum_{i=1}^{2N} Y_i (1 - W_i) = \frac{1}{N} Y^T (2W - 1), \quad (3)$$

where  $Y$  is the vector of the outcomes. As expected,  $\hat{\tau}$  is an unbiased estimator for  $\tau$ .

PROPOSITION 1. *For our sequential rerandomization,  $E(\hat{\tau} | X, \varphi_1 = \dots = \varphi_K = 1) = \tau$ .*

*Proof.* See the [Supplementary Material](#). In fact, for  $\hat{\tau}$  to be unbiased, we only require that the rerandomization criterion satisfy  $\varphi_k(X_{1:k}, W_{1:k}^*) = \varphi_k(X_{1:k}, 1 - W_{1:k}^*)$  for each  $k$  and that each group contain the same number of treatment and control units.  $\square$

Next consider the sampling variance of  $\hat{\tau}$ . Following the argument of [Morgan & Rubin \(2012\)](#), when the treatment effect is an additive constant for all units, we decompose  $Y_i$  as

$$Y_i = \hat{\beta}_0 + \hat{\beta}^T X_i + \tau W_i + \hat{e}_i \quad (i = 1, \dots, 2N), \quad (4)$$

where  $\hat{\beta}_0 + \hat{\beta}^T X_i$  is the projection of  $y_i(0)$  onto the space spanned by  $(1, X^T)$ , and the error  $\hat{e}_i$  is the projection of  $y_i(0)$  onto the orthogonal complement of that space. Letting  $\bar{e}_T$  and  $\bar{e}_C$  be the error means for the treatment and control groups, respectively, by (3) and (4) we have

$$\text{var}(\hat{\tau}) = \text{var}\{\hat{\beta}^T(\bar{X}_T - \bar{X}_C) + \bar{e}_T - \bar{e}_C\} = \hat{\beta}^T \text{cov}(\bar{X}_T - \bar{X}_C) \hat{\beta} + \text{var}(\bar{e}_T - \bar{e}_C), \quad (5)$$

where  $\bar{X}_T$  and  $\bar{X}_C$  are the covariate mean vectors of the treatment and control groups. A natural line of enquiry is to find the reduction in  $\text{cov}(\bar{X}_T - \bar{X}_C | X, \varphi_1 = \dots = \varphi_K = 1)$  under sequential rerandomization relative to  $\text{cov}(\bar{X}_T - \bar{X}_C | X)$  under complete randomization, which could be used to derive the reduction in the variance of the estimation for  $\tau$ . Recall that  $M_K$  is the Mahalanobis distance of the entire dataset after all sequential randomized allocations have been conducted.

THEOREM 1. *Let  $\nu = E(M_K | X, \varphi_1 = \dots = \varphi_K = 1)/p$ . We have*

$$\text{cov}(\bar{X}_T - \bar{X}_C | X, \varphi_1 = \dots = \varphi_K = 1) = \nu \text{cov}(\bar{X}_T - \bar{X}_C | X).$$

Proofs of Theorem 1 and all the following theorems, lemmas and propositions can be found in the [Supplementary Material](#).

THEOREM 2. *Let  $\tilde{\tau}$  be the estimator for  $\tau$  for complete randomization, and let  $\hat{\tau}$  be the estimator for  $\tau$  for sequential rerandomization. Assuming that the treatment effect is additive, we have*

$$\frac{\text{var}(\tilde{\tau}) - \text{var}(\hat{\tau})}{\text{var}(\tilde{\tau})} = (1 - \nu)R^2,$$

where  $R^2$  is the squared multiple correlation between  $Y$  and  $X$  in either the treatment or the control group.

These results can be seen as extensions of those presented in [Morgan & Rubin \(2012\)](#). When  $K = 1$ , the expression for  $\nu$  reduces to equation (9) in [Morgan & Rubin \(2012\)](#). Henceforth we shall write simply  $E(M_k | X)$  instead of  $E(M_k | X, \varphi_1 = \dots = \varphi_k = 1)$ , since the notation  $M_k$  clearly implies that sequential rerandomization has been conducted.

### 3.2. Asymptotic minimization of the expected Mahalanobis distance

Theorem 1 tells us that, under the additive treatment effect model,  $\text{var}(\hat{\tau})$  is minimized when  $E(M_K | X)$  is minimized. In this section we propose an asymptotically optimal strategy that minimizes  $E(M_K | X)$  and thus makes the estimation of average treatment effect most precise. To this end, we first seek the distribution of  $M_k$ , which is a truncated version of the distribution of  $M_k^*$ . Recall that  $X_1, \dots, X_K$  are treated as fixed and the randomness comes only from the treatment/control assignment. We further assume that the data are homogeneous so that  $\text{cov}(X_k) \approx \text{cov}(X_{1:k})$ ; the heterogeneous case will be discussed in § 5.1. By (1), the distribution of  $M_k^*$  depends on the  $p$ -dimensional random variable

$$D_k^* = \bar{X}_{T,k}^* - \bar{X}_{C,k}^* = \frac{1}{n_k} X_k W_k^* - \frac{1}{n_k} X_k (1 - W_k^*) = 2\bar{X}_{T,k}^* - 2\bar{X}_k. \quad (6)$$

As shown below in Lemma 1, when  $D_k^*$  is normally distributed, the distribution of  $M_k^*$  is fully determined by the value of  $M_{k-1}$ , which is a noncentral chi-squared distribution with noncentrality parameter proportional to  $M_{k-1}$ . Consequently, when choosing the threshold  $a_k$  in (2), we only need to use  $M_{k-1}$ , since conditional on  $M_{k-1}$ ,  $M_k^*$  is independent of  $M_1, \dots, M_{k-2}$ .

LEMMA 1. *Assume that  $D_k^* | X_k \sim \mathcal{N}\{0, 2n_k^{-1} \text{cov}(X_k)\}$  and  $\text{cov}(X_k) \approx \text{cov}(X_{1:k})$ . Let  $M_{k-1}$  be the Mahalanobis distance for the first  $k-1$  treatment and control groups after rerandomization with  $M_0 = 0$ ; then*

$$M_k^* | X_k, M_{k-1} \sim \frac{n_k}{n_{1:k}} \chi_p^2 \left( \frac{n_{1:k} - n_k}{n_k} M_{k-1} \right), \quad (7)$$

where  $\chi_p^2(\lambda)$  denotes a noncentral chi-squared distribution with  $p$  degrees of freedom and noncentrality parameter  $\lambda$ .

Remark 1. For sufficiently large  $n_1, \dots, n_K$ , the assumption that  $D_k^* | X_k \sim \mathcal{N}\{0, 2n_k^{-1} \text{cov}(X_k)\}$  holds under very general settings (Li & Ding, 2017). According to our sequential rerandomization procedure, the covariate mean of the  $k$ th treatment group, i.e., the term  $\bar{X}_{T,k}^*$  in (6), can be viewed as the mean of samples from a finite population without replacement. Under certain regularity conditions, the latter is known to follow a normal distribution asymptotically (Wald & Wolfowitz, 1944). By Hoeffding (1951) and Hájek (1961), a sufficient condition is as follows: the column vectors  $X_1, \dots, X_{2N}$  are independent and identically distributed  $p$ -dimensional random vectors from a distribution with finite third absolute moments and a positive-definite covariance matrix. Then, as  $n_k \rightarrow \infty$ ,  $\sqrt{n_k} D_k^*$  converges in distribution to  $\mathcal{N}\{0, 2 \text{cov}(X_k)\}$ . This result will be used to compute  $E(M_k | X)$  and derive the optimal strategy for sequential rerandomization.

Recall the sequential rerandomization criteria  $\varphi_1, \dots, \varphi_K$  defined in (2). We use the distribution given in (7) to choose  $a_k$  so that  $F_{M_k^*}(a_k) = \alpha_k$ , where  $F_{M_k^*}$  is the conditional distribution function of  $M_k^*$  given  $M_{k-1}$  and  $\alpha_k$  is the acceptance probability of each rerandomization. The number of randomizations required for  $\varphi_k$  to evaluate to 1 is distributed as a geometric random variable with expectation  $s_k = 1/\alpha_k$ . Hence, if we know how to choose  $s_k$ , we can choose  $a_k$  accordingly by the distribution of  $M_k^*$  given in Lemma 1, and we denote this by writing  $a_k = a_k(M_{k-1}, s_k)$ . It is reasonable to suppose that the experimenter, equipped with modern computational resources, may perform rerandomization a very large number of times. We may therefore assume that  $s_1, \dots, s_K$  are sufficiently large and  $M_1, \dots, M_K$  are correspondingly small. Using an asymptotic result for truncated noncentral chi-squared distributions, Lemma 2 below, we proceed to find an asymptotic expression, in Lemma 3, for the expected value of  $M_k$  conditional on  $M_{k-1}$ .

LEMMA 2. *Let  $M$  be a random variable that follows  $\chi_p^2(\lambda)$ , and let  $F_M$  be its cumulative distribution function. As  $a \downarrow 0$ ,*

$$F_M(a) \sim \frac{a^{p/2} \exp(-\lambda/2)}{2^{p/2} \Gamma(p/2 + 1)}, \quad E(M | M < a) \sim \frac{pa}{p + 2},$$

where  $\sim$  denotes asymptotic equivalence, i.e., for two positive functions  $f(x)$  and  $g(x)$  we write  $f \sim g$  as  $x \rightarrow x_0$  if and only if  $\lim_{x \rightarrow x_0} f(x)/g(x) = 1$ .

LEMMA 3. *Suppose that  $M_k^*$  ( $k = 1, \dots, K$ ) follows the distribution given in Lemma 1 and  $\text{pr}(M_k^* < a_k | X_k, M_{k-1}) = 1/s_k$ . Then, as  $s_k \uparrow \infty$  and  $M_{k-1} \downarrow 0$ ,*

$$E(M_k | X_k, M_{k-1}) \sim \frac{n_k}{n_{1:k}} C_p s_k^{-2/p} \left( 1 + \frac{n_{1:k} - n_k}{pn_k} M_{k-1} \right),$$

where  $C_p = 2p\{\Gamma(p/2 + 1)\}^{2/p}/(p + 2)$ .

Let the expected total number of rerandomizations  $S = s_1 + \dots + s_K$  be sufficiently large. Proposition 2 details the asymptotically optimal strategy for choosing  $s_1, \dots, s_K$ , in which optimality is achieved by asymptotically minimizing  $E(M_K | X)$  for fixed  $S$ .

**PROPOSITION 2.** *Suppose that  $M_k^*$  ( $k = 1, \dots, K$ ) follows the distribution given in Lemma 1. As  $S \uparrow \infty$ , in order to minimize  $E(M_K | X)$ , one should choose  $s_1, \dots, s_K$  so that*

$$s_{k-1} \approx \left( \frac{C_p n_{k-1}}{p n_k} s_k \right)^{p/(p+2)}, \quad C_p = \frac{2p}{p+2} \Gamma(p/2 + 1)^{2/p}. \quad (8)$$

#### 4. COMPARING SEQUENTIAL RERANDOMIZATION WITH MORGAN–RUBIN COMPLETE RERANDOMIZATION

In this section, we compare sequential rerandomization with Morgan–Rubin complete rerandomization. We begin by recalling the Morgan–Rubin complete rerandomization algorithm;  $2N$  units are assumed to be enrolled when the rerandomization starts and randomizations are conducted until the Mahalanobis distance  $M^*$  is smaller than some prespecified threshold  $a$ , where

$$M^* = \frac{N}{2} (\bar{X}_T^* - \bar{X}_C^*)^\top \text{cov}(X)^{-1} (\bar{X}_T^* - \bar{X}_C^*).$$

When the rerandomization stops, let  $M = M^*$ . Asymptotically, the distribution of  $M$  is a truncated chi-squared distribution with support  $(0, a)$ . This statistic  $M$  represents the same quantity as the statistic  $M_K$  in sequential rerandomization: namely, it is the Mahalanobis distance calculated on the entire sample after all units have received treatment assignment. If, in expectation, the same number of rerandomizations are conducted in Morgan–Rubin complete rerandomization and sequential rerandomization, it is tempting to conjecture that  $E(M | X)$ , which we define as the expected Mahalanobis distance from Morgan–Rubin complete rerandomization, is smaller than  $E(M_K | X)$ , since Morgan–Rubin complete rerandomization considers all  $(2N)!/(N!N!)$  possible allocations, whereas sequential rerandomization selects from a subset of those that are allowed by the sequential design. Surprisingly, as we will now show in Theorem 3, under certain asymptotic conditions the opposite is true.

**THEOREM 3.** *Let  $n_1, \dots, n_K$  be given and let  $S \in \mathbb{N}$  be the expected total number of rerandomizations. For Morgan–Rubin complete rerandomization, choose the threshold  $a$  such that  $\text{pr}(M^* < a | X) = 1/S$ ; for sequential rerandomization, choose  $s_1, \dots, s_K$  according to Proposition 2 under the constraint  $\sum_{i=1}^K s_i = S$ , and then choose thresholds  $a_k$  such that  $\text{pr}(M_k^* < a_k | X_k, M_{k-1}) = 1/s_k$ . Then, assuming that  $M_k^*$  given  $M_{k-1}$  ( $k = 1, \dots, K$ ) follows the distribution given in Lemma 1, as  $S \uparrow \infty$ ,*

$$E(M_K | X) \sim \frac{n_K}{N} E(M | X).$$

**COROLLARY 1.** *Under the assumptions of Theorem 3 and assuming  $n_1 = \dots = n_K$ , as  $S$  grows to infinity,  $E(M_K | X) \sim E(M | X)/K$ .*

**Remark 2.** We pause to offer some intuition for Theorem 3. The rerandomization of the last group is the most important step, because any imbalance between the first  $K - 1$  treatment and control groups may be cancelled out, making the entire dataset balanced once again. Heuristically, an efficient sequential rerandomization strategy need only ensure that the imbalance accumulated in the first  $K - 1$  groups is sufficiently small and then perform most rerandomizations for the last group. In fact, any strategy that satisfies the following two conditions would make Theorem 3 hold: (i) as  $S \uparrow \infty$ , every  $s_k$  does so too; (ii)  $S \sim s_K$ . The first condition ensures that every  $M_k$  will decrease to zero and so, by Lemma 1,  $Nn_K^{-1}M_K^*$  will converge to a  $\chi_p^2$  random variable. The second condition guarantees that, asymptotically,  $Nn_K^{-1}M_K$  and  $M$  are equivalent, in expectation, because they are truncated at the same threshold.

The consequences of the results in this section can be significant for clinical trials research. If a large number of individuals are enrolled simultaneously, Theorem 3 says that it is advantageous to use sequential rerandomization in lieu of Morgan–Rubin complete rerandomization. In the [Supplementary Material](#),

we conduct multiple simulation studies using both simulated and real datasets to show that sequential rerandomization achieves a smaller Mahalanobis distance in almost every practical setting.

## 5. DISCUSSION

### 5.1. Generalizations of our results

In practice, experimenters may prefer unequal allocation schemes where the numbers of treatment and control units are not the same (Hey & Kimmelman, 2014). Let  $\omega$  be the proportion of treatment assignments. If  $\omega$  is constant across all the groups, then as long as we use the correct version of Mahalanobis distance, see the [Supplementary Material](#), all our results still hold. More generally, our results can be extended to a heterogeneous dataset where  $\text{cov}(X_k)$  is very different across the groups. This happens when the clinical trial has a large time span or the groups of samples are collected at different places. The key is to find an appropriate form of Mahalanobis distance. In the [Supplementary Material](#), we propose to standardize the data separately for each group and then compute the Mahalanobis distance using the standardized variables with proper weights. This can be viewed as a generalization of the Mahalanobis distance defined in (1), and a corresponding generalized version of Lemma 1 is also proved; see Lemma A1 in the [Supplementary Material](#). Since Lemma 1 characterizes the conditional distribution of  $M_k^*$  and is the foundation of all the subsequent results, our main results, Proposition 2 and Theorem 3, follow by the same argument.

When the data are heterogeneous, minimizing the overall imbalance may not be sufficient and thus the strategy given in Proposition 2 becomes undesirable. For example, the effects of the covariates may change between the groups, and each group may have a unique systematic effect on the outcome. In such cases, one may want to achieve good balance within each group and choose a more uniform value for  $(s_1, \dots, s_K)$ . We point out that, in terms of within-group balance, sequential rerandomization is still superior to Morgan–Rubin complete rerandomization; see the [Supplementary Material](#) for details. Heuristically, this is because even if we let the threshold  $a$  of Morgan–Rubin complete rerandomization go to zero, we are only enforcing cancellation of the within-group imbalances, but their absolute values can be arbitrarily large.

### 5.2. Towards an optimal procedure

For classical dynamic randomization procedures, it is often assumed that the assignment of a unit must be determined as soon as the individual is enrolled. In the seminal work of Atkinson (1982), a type of Efron’s biased coin procedure (Efron, 1971) was proposed that achieves optimum performance when the underlying model is linear (see also Smith, 1984). A natural avenue of enquiry would be to find the optimal procedure for the case where the participants arrive in groups and the rerandomization technique is employed. Such questions have to be formulated very carefully. Even if all the participants arrive at the same time, the deterministic construction that minimizes the Mahalanobis distance is usually undesirable for the following two reasons. Firstly, for large sample sizes, the construction is not practical since finding the deterministic optimum is a nonconvex optimization problem. Secondly, we want the procedure to possess a certain degree of randomness to avoid selection bias (Antognini & Zagoraiou, 2017). Qin et al. (2016) introduced a procedure which can be applied when the participants arrive in pairs. For the  $k$ th pair, they considered the two 1 : 1 assignment schemes and chose the one that gives a smaller Mahalanobis distance of the first  $2k$  units with probability  $q \in (1/2, 1)$ . In the [Supplementary Material](#) we apply this procedure to a real dataset. When  $q = 0.75$ , the value suggested in Qin et al. (2016), the Mahalanobis distance of the entire dataset is only slightly smaller than that of Morgan–Rubin complete rerandomization, but it is greater than those of our sequential designs. For comparison, when  $q = 1$ , which makes the whole procedure deterministic, the Mahalanobis distance reduces dramatically. We believe that when the group size of a sequential design is small, the covariate imbalance can only be efficiently minimized at the cost of selection bias, i.e., the procedure being more deterministic. Kapelner & Krieger (2014) offered a more complicated dynamic procedure which also uses Mahalanobis distance, but some participants may wait a long time before being assigned. We believe that to find an optimal procedure, one needs to strike a balance between the following factors: the covariate imbalance as measured by Mahalanobis distance, the group size of the sequential enrolment design, and the randomness and computational cost of the procedure.

## ACKNOWLEDGEMENT

We thank the referees for comments which helped to improve the quality of the paper. Ernst was funded by the U.S. Office of Naval Research. Rubin was funded by the U.S. National Science Foundation and National Institutes of Health and by a Google Faculty Fellowship.

## SUPPLEMENTARY MATERIAL

Supplementary material available at *Biometrika* online includes proofs of all the theoretical results, some theoretical generalizations of our results, and simulation studies with both simulated and real datasets.

## REFERENCES

ANTOGNINI, A. B. & ZAGORAIOU, M. (2017). Estimation accuracy under covariate-adaptive randomization procedures. *Electron. J. Statist.* **11**, 1180–206.

ATHEY, S. & IMBENS, G. W. (2017). The econometrics of randomized experiments. In *Handbook of Economic Field Experiments*, vol. 1, E. Duflo & A. Banerjee, eds. Amsterdam: Elsevier, pp. 73–140.

ATKINSON, A. C. (1982). Optimum biased coin designs for sequential clinical trials with prognostic factors. *Biometrika* **69**, 61–7.

BERGER, V. W. (2010). Minimization, by its nature, precludes allocation concealment, and invites selection bias. *Contemp. Clin. Trials* **31**, 406.

DELAVANDE, A., WAGNER, Z. & SOOD, N. (2016). The impact of repeat HIV testing on risky sexual behavior: Evidence from a randomized controlled trial in Malawi. *J. AIDS Clin. Res.* **7**, DOI: 10.4172/2155-6113.1000549.

EFRON, B. (1971). Forcing a sequential experiment to be balanced. *Biometrika* **58**, 403–17.

GREEVY, R., LU, B., SILBER, J. H. & ROSENBAUM, P. (2004). Optimal multivariate matching before randomization. *Biostatistics* **5**, 263–75.

HÁJEK, J. (1961). Some extensions of the Wald-Wolfowitz-Noether theorem. *Ann. Math. Statist.* **32**, 506–23.

HEY, S. P. & KIMMELMAN, J. (2014). The questionable use of unequal allocation in confirmatory trials. *Neurology* **82**, 77–9.

HOEFFDING, W. (1951). A combinatorial central limit theorem. *Ann. Math. Statist.* **22**, 558–66.

HOLSCHUH, N. (1980). Randomization and design: I. In *R. A. Fisher: An Appreciation*. Berlin: Springer, pp. 35–45.

HU, Y. & HU, F. (2012). Balancing treatment allocation over continuous covariates: A new imbalance measure for minimization. *J. Prob. Statist.* **2012**, DOI: 10.1155/2012/842369.

IMAI, K., KING, G. & STUART, E. A. (2008). Misunderstandings between experimentalists and observationalists about causal inference. *J. R. Statist. Soc. A* **171**, 481–502.

KAPELNER, A. & KRIEGER, A. (2014). Matching on-the-fly: Sequential allocation with higher power and efficiency. *Biometrics* **70**, 378–88.

LI, X. & DING, P. (2017). General forms of finite population central limit theorems with applications to causal inference. *J. Am. Statist. Assoc.* **112**, 1759–69.

LIN, Y., ZHU, M. & SU, Z. (2015). The pursuit of balance: An overview of covariate-adaptive randomization techniques in clinical trials. *Contemp. Clin. Trials* **45**, 21–5.

MORGAN, K. L. & RUBIN, D. B. (2012). Rerandomization to improve covariate balance in experiments. *Ann. Statist.* **40**, 1263–82.

MORGAN, K. L. & RUBIN, D. B. (2015). Rerandomization to balance tiers of covariates. *J. Am. Statist. Assoc.* **110**, 1412–21.

QIN, Y., LI, Y. & HU, F. (2016). An optimal method for covariate balancing and its properties. *arXiv*: 1611.02802.

RUBIN, D. B. (1974). Estimating causal effects of treatments in randomized and nonrandomized studies. *J. Educ. Psychol.* **66**, 688–701.

RUBIN, D. B. (1979). Using multivariate matched sampling and regression adjustment to control bias in observational studies. *J. Am. Statist. Assoc.* **74**, 318–28.

SAVAGE, L. J. (1962). *The Foundations of Statistical Inference*. London: Methuen.

SMITH, R. L. (1984). Properties of biased coin designs in sequential clinical trials. *Ann. Statist.* **12**, 1018–34.

URBACH, P. (1985). Randomization and the design of experiments. *Phil. Sci.* **52**, 256–73.

WALD, A. & WOLFOWITZ, J. (1944). Statistical tests based on permutations of the observations. *Ann. Math. Statist.* **15**, 358–72.

XU, Z. & KALBFLEISCH, J. D. (2013). Repeated randomization and matching in multi-arm trials. *Biometrics* **69**, 949–59.

[Received on 9 October 2017. Editorial decision on 5 April 2018]