

Fast computation of exact confidence intervals for randomized experiments with binary outcomes

P. M. Aronow Haoge Chang Patrick Lopatto*

Abstract

Consider a randomized experiment with a binary outcome, such that half of units are randomly assigned to receive a treatment and the other half are assigned to a control group. In this setting, exact confidence intervals for the average causal effect of the treatment can be computed through a series of permutation tests. This approach requires minimal assumptions and is valid for all sample sizes, as it does not rely on large-sample approximations such as the central limit theorem. We show that these confidence intervals can be found in $O(n \log n)$ permutation tests. Prior to this work, the most efficient known construction was given by Li and Ding, who showed $O(n^2)$ permutation tests suffice. We also demonstrate how to construct confidence intervals using Monte Carlo permutation tests, instead of more computationally demanding exact tests, while retaining the guarantee that the intervals contain the true parameter at no less than the given rate. Taken together, the number of computations required by our method is $O(n^2 \epsilon^{-2} (\log n) \log(n \epsilon^{-1} \log n))$, with a precision loss proportional to ϵ , improved from the previous best of $O(2^n n^{5/2})$. Our results thus facilitate exact inference as a viable option for randomized experiments far larger than those accessible by previous methods.

*P.L. is partially supported by NSF postdoctoral fellowship DMS-220289.

Contents

1	Introduction	1
1.1	Statistical Background	1
1.2	Algorithmic Perspective and Main Results	3
1.3	Open Problems	5
1.4	Outline	6
2	Preliminaries	6
2.1	Notation	6
2.2	Previous Work	7
3	Basic Properties	9
4	Fast Computation of Confidence Intervals	10
5	Monte Carlo Intervals	11
6	Missing Data	13
7	Proofs for Section 3	14
7.1	Proof of Lemma 3.1	15
7.2	Proof of Proposition 3.2	15
7.3	Proof of Proposition 3.3	16
8	Proofs for Section 4	17
8.1	Preliminary Results	17
8.2	Proof of Theorem 4.1	19
8.3	Proof of Lemma 4.2	20
8.4	Proof of Lemma 8.10	26
9	Proof of Proposition 5.3	28
10	Proof of Proposition 6.2	30

1 Introduction

There is a vast literature on causal inference in randomized experiments with binary outcomes (see, for example, [IR15] and the references therein). This setting includes a large class of important examples, such as experiments probing the causal effect of an treatment for a dangerous disease, a design change in an online advertisement, or a reminder about the date of an upcoming election. In each case, the outcome for each subject – whether they are hospitalized during the study period, click a banner, or cast a ballot, respectively – is binary.

The study of causal inference for binary outcomes naturally leads to a variety of algorithmic questions, which are the focus of this article. To contextualize our results, we first review the relevant statistical background. We then describe previous work on algorithms for inference on randomized experiments with binary outcomes, and present our contributions. We conclude this introduction by discussing some open problems.

1.1 Statistical Background

A fundamental problem in causal inference is the construction of *confidence intervals* for an unknown causal parameter $\theta \in \mathbb{R}$. We observe data \mathbf{D} generated by some random process, such as a randomized experiment, whose distribution is assumed to be function of the unknown θ . We would like to use the observed data \mathbf{D} to provide an *interval estimate* for θ . Specifically, we fix some $\alpha \in (0, 1)$ and aim to construct an interval $\mathcal{J}_\alpha(\mathbf{D}) = [L_\alpha(\mathbf{D}), U_\alpha(\mathbf{D})]$ such that

$$\mathbb{P}_\theta(\theta \in \mathcal{J}_\alpha(\mathbf{D})) \geq 1 - \alpha \tag{1.1}$$

holds for all possible θ , where \mathbb{P}_θ denotes the probability with respect to the random process arising from θ . A conventional choice is $\alpha = 0.05$, leading to the well-known 95% confidence interval. Among intervals that satisfying (1.1), shorter intervals \mathcal{J}_α more precisely estimate θ than longer ones, and are therefore more desirable. Crucially, \mathcal{J}_α contains θ with high probability, regardless of the true value of θ . This uniform coverage property guarantees that if many intervals \mathcal{J}_α are constructed from independent data sets (and α is small), then with high probability most of these intervals will cover their corresponding true parameters.¹ This is the standard justification for the use of confidence intervals in modern science, where they are enormously popular. Further details may be found in [Was04].

Many traditional approaches to constructing confidence intervals for experiments with binary outcomes are based on a binomial model for the outcome distribution. Examples include Wald intervals [Was04, Chapter 10], which are based on a normal approximation valid in large samples, as well as other interval estimators [BB19, SS80]. However, the assumptions underlying the binomial model are highly problematic in typical experimental designs. This was noted by Robins in [Rob88], whose discussion we now briefly summarize.

Consider an experiment with n subjects, where m subjects are assigned to treatment and the remaining $n - m$ are assigned to control, and the random assignment is such that all possible configurations are equally likely. Robins identifies two modeling assumptions under which binomial confidence intervals will cover the true average causal effect at the nominal rate.

1. If a subject is exposed to treatment, their outcome may be modeled as a Bernoulli random variable, and the treatment outcomes across subjects are independent with common mean

¹Precise quantitative statements of this kind can easily be deduced from concentration of measure bounds. Specifically, one may apply Hoeffding’s inequality to the indicator functions for the events $\{\theta \in \mathcal{J}_\alpha\}$.

p_1 . Similarly, the control outcomes are independent Bernoulli random variables with common mean p_2 . Define the average causal effect of treatment as $p_1 - p_2$.

2. The subjects in the study are drawn uniformly at random from some near-infinite “superpopulation.” We let p_1 be the proportion of subjects in the superpopulation who whose outcome would be 1 under treatment, and similarly let p_2 be the proportion whose outcome would be 1 under control. Again define the average causal effect of treatment as $p_1 - p_2$.

Robins notes that Assumption (1) is untenable in most contexts, since it does not account for between-subject variation. For example, disease risk may vary among subjects in a clinical trial according to pre-existing medical conditions and demographic characteristics. Further, Assumption (2) is almost always false, since typical subject recruitment strategies do not result in a uniform sample from a well-defined superpopulation. Consider, for example, a clinical trial where subjects are recruited through newspaper advertisements. It is generally implausible that all members of the target population will read and respond to such advertisements at identical rates. Further, when both Assumption (1) and Assumption (2) fail, it is unclear what unknown parameter θ the binomial confidence intervals are supposed to estimate, and why they should satisfy (1.1).

Given these deficiencies of the binomial model, we adopt instead a finite population model, which we analyze using the *potential outcomes* framework. This perspective originated in work of Neyman and has subsequently been developed by many researchers (see [IR15] for details). Suppose we have a group of n subjects, and label them arbitrarily from 1 to n . Our fundamental assumption is that the observed outcome for subject i depends only on the index i and whether that subject is assigned to treatment or control. This is known as the *stable unit treatment value assumption*, and abbreviated SUTVA. In the particular, SUTVA rules out between-subject interaction effects. We let $\mathbf{y}_i = (y_i(0), y_i(1))$ with $y_i(0), y_i(1) \in \{0, 1\}$ denote the *potential outcomes* for the i -th subject, where $y_i(1)$ is the outcome that would be observed if the subject were assigned to treatment, and $y_i(0)$ is the outcome that would be observed if the subject were assigned to control. In the experiment, only *one* of these potential outcomes is observed; the other remains unknown. We aim to estimate the *sample average treatment effect*, defined as

$$\tau(\mathbf{y}) = \frac{1}{n} \sum_{j=1}^n (y_j(1) - y_j(0)). \tag{1.2}$$

The quantity $\tau(\mathbf{y})$ is defined using only the subjects in the experiment, and says nothing about the causal effect in any larger population. Such a generalization would require additional assumptions. However, precise estimates of $\tau(\mathbf{y})$ are still broadly useful for studying causal claims. For example, if a company asserts that their latest drug greatly reduces the chance of stroke in elderly people, and a study of the average causal effect in a group of such people produces a precise estimate centered at 0, it is reasonable to reject the company’s claim.

Since we view the potential outcomes of the subjects as fixed, the only randomness in the experiment comes from the assignment of subjects to treatment and control. We retain our previous setup where m subjects are assigned to treatment, the remainder are assigned to control, and all such assignments are equally likely. Under this design, Robins gives a confidence interval for $\tau(\mathbf{y})$ that relies on a large-sample normal approximation, which is explicitly justified by randomization. While this interval is guaranteed to cover $\tau(\mathbf{y})$ at its nominal rate in the limit as n goes to infinity, there is no quantitative estimate available for its accuracy.

1.2 Algorithmic Perspective and Main Results

We now explain how the construction of confidence intervals for the causal effect in a randomized experiment with a binary outcome, where n subjects are randomized to groups of fixed sizes m and $n - m$, may be cast as a computational problem. Specifically, we seek confidence intervals that cover $\tau(\mathbf{y})$ with probability at least $1 - \alpha$ regardless of the sample size or potential outcomes (in the sense of (1.1)).

Define a random vector $\mathbf{Z} \in \{0, 1\}^n$ by letting $Z_i = 0$ if subject i is assigned to control, and $Z_i = 1$ if subject i is assigned to treatment. We let \mathbf{Y} be the vector of observed outcomes, where

$$Y_j = Z_j y_j(1) + (1 - Z_j) y_j(0) \quad (1.3)$$

represents the outcome that the experimenter observes for subject j after the experiment is completed. Using the known vectors \mathbf{Z} and \mathbf{Y} , the unknown parameter $\tau(\mathbf{y})$ may be estimated using the *Neyman estimator* $T(\mathbf{Y}, \mathbf{Z})$. Informally speaking, this estimator subtracts the average of the observed outcomes in the control group from the average of the observed outcomes in the treatment group. The precise definition is

$$T(\mathbf{Y}, \mathbf{Z}) = \sum_{j=1}^n \frac{Z_j Y_j}{m} - \sum_{j=1}^n \frac{(1 - Z_j) Y_j}{n - m}. \quad (1.4)$$

An elementary calculation shows that $\mathbb{E}[T(\mathbf{Y}, \mathbf{Z})] = \tau(\mathbf{y})$, where the expectation is taken over the variable \mathbf{Z} . Hence $T(\mathbf{Y}, \mathbf{Z})$ is an *unbiased* estimator for $\tau(\mathbf{y})$.

As a preliminary step toward constructing a confidence interval for $\tau(\mathbf{y})$, we consider the probability that some possible configuration \mathbf{w} of the potential outcomes could produce data at least as extreme as the observed data (\mathbf{Y}, \mathbf{Z}) . To make this notion precise, we consider an arbitrary length n vector \mathbf{w} of elements w_i such that $w_i(0), w_i(1) \in \{0, 1\}$, and compute the *p-value* given by

$$p(\mathbf{w}, \mathbf{Y}, \mathbf{Z}) = \mathbb{P}\left(|T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}}) - \tau(\mathbf{w})| \geq |T(\mathbf{Y}, \mathbf{Z}) - \tau(\mathbf{w})|\right). \quad (1.5)$$

Here $\tilde{\mathbf{Z}}$ is independent from \mathbf{Z} and has an identical distribution, and $\tilde{\mathbf{Y}}$ is the vector of observed outcomes in a hypothetical experiment with potential outcomes \mathbf{w} and randomization $\tilde{\mathbf{Z}}$, so that $\tilde{Y}_j = \tilde{Z}_j w_j(1) + (1 - \tilde{Z}_j) w_j(0)$. Calculating the *p-value* (1.5) is known as performing a *permutation test*, since the probability may be computed by averaging over all possible randomizations $\tilde{\mathbf{Z}}$.

Recall that $\mathbb{E}[T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}})] = \tau(\mathbf{w})$. Then, informally speaking, the *p-value* measures how atypical (far from the center of the distribution) the observed data $T(\mathbf{Y}, \mathbf{Z})$ is, under the assumption that \mathbf{w} is the true set of potential outcomes for the subjects in the experiment. The permutation test (1.5) was essentially known to Fisher in the 1930s [IR15, Chapter 5], and can be used to test the hypothesis that the true potential outcomes \mathbf{y} are equal to some specific \mathbf{w} (by rejecting this hypothesis if $p(\mathbf{w}, \mathbf{Y}, \mathbf{Z})$ is sufficiently small).

Recently, it was realized by Rigdon and Hudgens that exact confidence intervals for $\tau(\mathbf{y})$ could be constructed from a series of these permutation tests [RH15], in the following way. Given observed data (\mathbf{Y}, \mathbf{Z}) , consider *all* potential outcome vectors \mathbf{w} that could produce (\mathbf{Y}, \mathbf{Z}) , if \mathbf{w} were the true unknown set of potential outcomes.² (The binary outcome assumption ensures that this set is finite.) Fix $\alpha \in (0, 1)$, and declare such a \mathbf{w} *compatible* if $p(\mathbf{w}, \mathbf{Y}, \mathbf{Z}) \geq \alpha$, and *incompatible* otherwise. Let $L_\alpha(\mathbf{Y}, \mathbf{Z})$ equal the smallest value of $\tau(\mathbf{w})$ where \mathbf{w} ranges over all compatible

²This already rules out many vectors \mathbf{w} . For example, if $n = 2$ and we randomize into equal groups of size 1, and the observed data yields $T(\mathbf{Y}, \mathbf{Z}) = 1$, then the configuration $\mathbf{w}_1 = \mathbf{w}_2 = (0, 0)$ is not possible.

\mathbf{w} , and let $U_\alpha(\mathbf{Y}, \mathbf{Z})$ equal the largest value of $\tau(\mathbf{w})$ where \mathbf{w} ranges over all compatible \mathbf{w} . Set $\mathcal{I}_\alpha = [L_\alpha(\mathbf{X}), U_\alpha(\mathbf{X})]$. Then it is straightforward to show that \mathcal{I}_α will contain $\tau(\mathbf{y})$ with probability at least $1 - \alpha$, regardless of the value of \mathbf{y} .³

Constructing \mathcal{I}_α through permutation tests eliminates any reliance on uncontrolled large-sample approximations, but has the drawback of requiring significant computational resources. While a naive construction of \mathcal{I}_α might proceed by directly computing (1.5) for all possible \mathbf{w} to find L_α and U_α , this is extremely inefficient, and infeasible even for small n . The problem of efficiently determining \mathcal{I}_α then reduces to two subproblems. First, minimize the number of permutation tests required to find L_α and U_α . Second, minimize the computational effort required to perform a permutation test.

Regarding the first subproblem, Rigdon and Hudgens showed that \mathcal{I}_α can be constructed using at most $O(n^4)$ permutation tests [RH15]. Later, Li and Ding showed that the same intervals could be constructed in $O(n^2)$ permutation tests in the balanced case, where an equal number of subjects are assigned to treatment and control (that is, $n = 2m$), and conjectured that their algorithm also reproduces the Rigdon–Hudgens intervals in unbalanced trials [LD16]. Specifically, they showed that the interval their algorithm returns always contains the corresponding Rigdon–Hudgens interval \mathcal{I}_α , but were unable to exclude the possibility that theirs are strictly larger for unbalanced designs.

We now consider the second subproblem. An *exact permutation test* enumerates all $\binom{n}{m}$ possible assignments of subjects, which grows exponentially in n and quickly becomes computationally infeasible. Rigdon and Hudgens therefore recommend performing approximate permutation tests, where the permutation p -value is approximated by Monte Carlo simulation with a fixed number of samples K from the distribution of $\tilde{\mathbf{Z}}$ in (1.5). However, such intervals may fail to cover at the nominal rate due to Monte Carlo error. Further, previous literature has not provided any formal analysis regarding the question of how large K should be taken so that the resulting intervals \mathcal{I}_α cover the unknown parameter $\tau(\mathbf{y})$ with high probability. Such an analysis is not entirely straightforward, since if one uses approximate permutation tests in the Rigdon–Hudgens and Li–Ding algorithms described previously, the acceptance and rejection decisions made by successive tests may be highly correlated.

We make two contributions, corresponding to the two subproblems identified previously. We focus on the case of balanced experiments, where the treatment and control groups have equal size, as this setting is common in practice and it is minimax optimal in terms of efficiency in estimation.⁴ First, we show that the exact confidence intervals \mathcal{I}_α from [RH15] can be constructed in $O(n \log n)$ permutation tests, improving on the $O(n^2)$ tests required by [LD16]. Our proof is based on a new monotonicity property for the permutation p -values (1.5), which permits us to greatly reduce the effective search space when finding L_α and U_α . Second, we show how to construct confidence intervals using approximate permutation tests that are guaranteed to cover the sample average treatment effect with probability $1 - \alpha$ for any $\alpha > 0$. We further show that by increasing K , these intervals will approximate those of [RH15] arbitrarily well, and quantify how K must scale in order to guarantee any desired coverage rate and approximation accuracy as n increases. The Monte Carlo guarantees are proved by a careful application of standard concentration inequalities, taking into account the details of our new algorithm for constructing \mathcal{I}_α .

In fact, we provide finite-sample bounds for both subproblems, with small, explicit constants. We are able to construct exact confidence intervals in $8n \log_2(n)$ permutation tests, and guarantee

³See Lemma 3.1 in Section 3, below.

⁴See also [Kal18, Section 2] for a justification of balanced randomization via its blinding properties. In general, [Kal18] shows that balanced complete randomization enjoys a minimax property even when potentially prognostic covariates are available to the designer of the experiment.

that our Monte Carlo intervals cover the unknown parameter $\tau(\mathbf{y})$ with probability $1 - \alpha$ if

$$K \geq \frac{1}{\varepsilon^2} \log \left(\frac{16n \log_2(n)}{\varepsilon} \right), \quad (1.6)$$

where $\varepsilon < \alpha$ is an approximation parameter defined in Section 5, and smaller values of ε correspond to confidence intervals that are smaller on average. (For both of these bounds, we have assumed that $n \geq 10$ in order to state them simply.) Observe that a single Monte Carlo sample of $T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}})$ can be evaluated in $O(n)$ arithmetical operations, by using the Fisher–Yates shuffle algorithm to generate the randomization $\tilde{\mathbf{Z}}$ into equal groups. Then the overall complexity of algorithm producing the Monte Carlo intervals is

$$O \left(\frac{n^2 \log n}{\varepsilon^2} \log \left(\frac{n \log(n)}{\varepsilon} \right) \right). \quad (1.7)$$

We note that the previous best construction of permutation confidence intervals with guaranteed coverage in [LD16] requires $O(2^n n^{5/2})$ operations, due to the $O(n^2)$ permutation tests required, the $\binom{2m}{m} = O(2^n n^{-1/2})$ randomizations required for an exact permutation test, and the $O(n)$ operations required for a single randomization.

Our algorithmic contributions permit permutation inference to be applied to experiments far larger than those accessible by previous methods. Additionally, we also extend our analysis to experiments with missing data, a common statistical problem that has not been addressed previously in this context.

1.3 Open Problems

Many interesting questions remain. First, it is natural to ask whether an even more efficient construction of \mathcal{I}_α is possible, either by improving the number of permutation tests necessary or finding a more efficient way to compute (approximate) permutation p -values. We are unaware of any formal analysis of the computational hardness of such questions. Regarding the Monte Carlo permutation tests, we do not know how to implement any of the standard variance reduction techniques in this context, such as importance sampling, which would provide a practical (though non-asymptotic) improvement in run time. We note that these methods have been successfully applied to permutation inference in other contexts [BB19, MPS88].

Second, a number of questions arise from relaxing the various assumptions made earlier. For example, is efficient approximate inference possible in the case of a general bounded and discrete outcome, or a continuous outcome supported on a bounded interval? Or for unbalanced designs? Additionally, a more realistic setup is to suppose that each subject comes with additional data, such as demographic information and medical history, that is correlated with the potential outcomes. Rigdon and Hudgens show that their algorithm adapts to designs using stratified randomization, where subjects are independently randomized to treatment or control within subgroups formed using these additional covariates [RH15]. However, they observe that this approach quickly becomes computationally intractable, requiring roughly $O(n^{4k})$ permutation tests, where k is the number of strata. It would be interesting to find a fast method to analyze stratified designs, and to investigate more sophisticated randomization mechanisms involving clustering, rerandomization, and matching [MR12, GLSR04].

Finally, it is of practical importance to find an efficient algorithm for prospective sample size calculations for the exact intervals \mathcal{I}_α . Researchers often wonder how many experimental subjects they need to detect an effect with high probability (that is, have $0 \notin \mathcal{I}_\alpha$), assuming the effect is larger than some given threshold. If too few subjects are used, the experiment is uninformative,

and if too many subjects are used, the experiment is potentially wasteful. While it is well known how to perform sample size estimation in simple settings such as the binomial model discussed earlier, we are unaware of any rigorous results for the Rigdon–Hudgens permutation intervals.

1.4 Outline

We begin in Section 2 by introducing our notation and reviewing previous work in detail. In Section 3, we present basic properties of the confidence intervals constructed in [RH15]. In Section 4, we show how to construct these intervals in $O(n \log n)$ permutation tests. In Section 5, we construct exact intervals using Monte Carlo permutation tests. We consider missing data in Section 6. The remaining sections contain proofs of the assertions in the previous sections.

2 Preliminaries

2.1 Notation

Let $\mathbb{Z}_{>0}$ denote the positive integers. We consider a group of $n \in \mathbb{Z}_{>0}$ subjects who undergo an experiment with a binary outcome. We always suppose that the experiment is balanced, so that $n = 2m$ for some $m \in \mathbb{Z}_{>0}$, m subjects are assigned to treatment, and the remaining m subjects are assigned to control.

A *potential outcome table* is any vector

$$\mathbf{w} \in \{(0, 0), (0, 1), (1, 0), (1, 1)\}^n. \quad (2.1)$$

Letting $\mathbf{w}_i = (w_i(0), w_i(1))$ for $i \in \llbracket 1, n \rrbracket$, the coordinate $w_i(1)$ denotes the outcome for the i -th subject under treatment, and $w_i(0)$ denotes the outcome under control. We use \mathbf{y} to denote the potential outcome table for the subjects in the experiment, which is a fixed (but unknown) ground truth. We use \mathbf{w} to denote an arbitrary vector in $\{(0, 0), (0, 1), (1, 0), (1, 1)\}^n$. Given some table \mathbf{w} , we let the corresponding count vector $\mathbf{v} = \mathbf{v}(\mathbf{w}) \in \mathbb{Z}_{\geq 0}^{\{0,1\}^2}$ be defined by

$$\mathbf{v} = (v_{11}, v_{10}, v_{01}, v_{00}), \quad v_{ab} = \sum_{j=1}^n \mathbb{1}\{\mathbf{w}_j = (a, b)\}. \quad (2.2)$$

We now introduce notation for the randomization of the subjects. Consider the random vector

$$\mathbf{Z} = (Z_1, Z_1, \dots, Z_n) \in \{0, 1\}^n \quad (2.3)$$

whose distribution is uniform over the set

$$\mathcal{Z}(n) = \left\{ \mathbf{z} \in \{0, 1\}^n : \sum_{i=1}^n z_i = m \right\}. \quad (2.4)$$

Here, the indices i such that $Z_i = 1$ represent the subjects assigned to treatment. The vector of observed experimental outcomes $\mathbf{Y} = (Y_j)_{j=1}^n$ is then given by

$$Y_j = Z_j y_j(1) + (1 - Z_j) y_j(0). \quad (2.5)$$

We say that a table \mathbf{w} is *possible* given the observed data \mathbf{Y} if for every $j \in \llbracket 1, n \rrbracket$, we have $\mathbf{w}_j(1) = Y_j$ if $Z_j = 1$, and $\mathbf{w}_j(0) = Y_j$ if $Z_j = 0$. We say that a count vector \mathbf{v} is possible if it arises

from some possible table \mathbf{w} .

The sample average treatment effect for a potential outcome table is defined by

$$\tau(\mathbf{w}) = \frac{1}{n} \sum_{j=1}^n (w_j(1) - w_j(0)). \quad (2.6)$$

The Neyman estimator of $\tau(\mathbf{y})$ is given by

$$T(\mathbf{Y}, \mathbf{Z}) = \sum_{j=1}^n \frac{Z_j Y_j}{m} - \sum_{j=1}^n \frac{(1 - Z_j) Y_j}{n - m}. \quad (2.7)$$

This estimator depends on (\mathbf{Y}, \mathbf{Z}) only through the observed count vector $\mathbf{n} \in \mathbb{Z}_{\geq 0}^{\{0,1\}^2}$ defined by

$$\mathbf{n} = (n_{11}, n_{10}, n_{01}, n_{00}), \quad n_{zy} = \sum_{j=1}^n \mathbb{1}\{Z_j = z, Y_j = y\}, \quad (2.8)$$

for $z \in \{0, 1\}$ and $y \in \{0, 1\}$. For example, n_{10} is the number of subjects in the treatment group with observed outcome 0. We overload the notation slightly and write $T(\mathbf{n})$ for the value of the Neyman estimator $T(\mathbf{Y}, \mathbf{Z})$, if \mathbf{n} is the observed count vector for (\mathbf{Y}, \mathbf{Z}) . We will do this without further comment for all other functions of (\mathbf{Z}, \mathbf{Y}) that depend only on the associated vector \mathbf{n} , and similarly for functions of tables \mathbf{w} that depend only on the associated count vector \mathbf{v} .

Throughout the paper, $\llbracket a, b \rrbracket$ denotes the set $\{k \in \mathbb{Z} : a \leq k \leq b\}$. We write \log for the natural logarithm. When we use the base 2 logarithm, it is always denoted by \log_2 . Additionally, we often abbreviate probability mass function as pmf. Further notation, used only in the proofs, is introduced in Section 7.

2.2 Previous Work

2.2.1 Rigdon and Hudgens

We begin by recalling the confidence interval procedure proposed in [RH15]. As \mathbf{w} ranges over all tables, $\tau(\mathbf{w})$ takes every value in the set

$$\mathcal{S}(n) = \left\{ -\frac{n}{n}, -\frac{n-1}{n}, \dots, \frac{0}{n}, \dots, \frac{n-1}{n}, \frac{n}{n} \right\}, \quad (2.9)$$

which has $2n + 1$ elements. After the experiment is completed, Rigdon and Hudgens note that the set of $\tau(\mathbf{w})$ for \mathbf{w} that are possible given the observed data is further restricted to

$$\mathcal{C}(\mathbf{Y}, \mathbf{Z}) = \left\{ \frac{1}{n} \left(\sum_{j=1}^n Y_j (2Z_j - 1) - m \right), \dots, \frac{1}{n} \left(\sum_{j=1}^n Y_j (2Z_j - 1) - m + n \right) \right\}, \quad (2.10)$$

which has $n + 1$ elements. We observe that \mathcal{C} depends on (\mathbf{Y}, \mathbf{Z}) only through the associated observed count vector \mathbf{n} .

We now consider some $\alpha \in (0, 1)$ and a realization of (\mathbf{Y}, \mathbf{Z}) . We will construct a confidence interval for τ at level $1 - \alpha$ (that is, one that contains the true parameter with probability at least $1 - \alpha$). We first describe a permutation test for the sharp null hypothesis that $\mathbf{y} = \mathbf{w}$ for some given \mathbf{w} , then explain how to form a confidence interval from a series of these tests.

Definition 2.1. Let $\tilde{\mathbf{Z}}$ be a random vector independent from \mathbf{Z} with the same distribution, and let $\tilde{\mathbf{Y}}$ be the vector with entries

$$\tilde{Y}_j = \tilde{Z}_j w_j(1) + (1 - \tilde{Z}_j) w_2(0). \quad (2.11)$$

We say that \mathbf{w} is compatible with data (\mathbf{Y}, \mathbf{Z}) at level $1 - \alpha$ if⁵

$$p(\mathbf{w}, \mathbf{Y}, \mathbf{Z}) = \mathbb{P} \left(|T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}}) - \tau(\mathbf{w})| \geq |T(\mathbf{Y}, \mathbf{Z}) - \tau(\mathbf{w})| \right) \geq \alpha. \quad (2.12)$$

We say that \mathbf{w} is compatible with an observed count vector \mathbf{n} if (2.12) holds with $T(\mathbf{Y}, \mathbf{Z})$ replaced by $T(\mathbf{n})$, and define $p(\mathbf{w}, \mathbf{n})$ similarly.

Let \mathbf{v} denote the count vector associated to \mathbf{w} , as defined in (2.2), and let \mathbf{n} be as in (2.8). Then it is straightforward to see that $p(\mathbf{w}, \mathbf{Y}, \mathbf{Z}) = p(\mathbf{v}, \mathbf{n})$, so that the p -value depends on $(\mathbf{w}, \mathbf{Y}, \mathbf{Z})$ only through the count vectors \mathbf{n} and \mathbf{v} .

Algorithm 2.2. A *permutation test* takes as inputs $\alpha \in (0, 1)$, a potential outcome table \mathbf{w} , and a vector $\mathbf{n} \in \mathbb{Z}_{\geq 0}^{\{0,1\}^2}$. It returns a binary decision, accept or reject, as follows. The algorithm computes $p(\mathbf{w}, \mathbf{n})$ by direct enumeration of all $\binom{n}{m}$ elements of the set $\mathcal{Z}(\mathbf{n})$ (defined in (2.4)), over which $\tilde{\mathbf{Z}}$ is uniformly distributed. If $p(\mathbf{w}, \mathbf{n}) \geq \alpha$, it accepts. Otherwise, it rejects.

By using $p(\mathbf{v}, \mathbf{n})$ in place of $p(\mathbf{w}, \mathbf{n})$, we can generalize the previous definition to use count vectors \mathbf{v} as inputs.

We now present the confidence interval procedure from [RH15]. Given an observed count vector \mathbf{n} , set

$$\mathcal{T}(\mathbf{n}) = \{(i, j, k, l) : i \in \llbracket 0, n_{11} \rrbracket, j \in \llbracket 0, n_{10} \rrbracket, k \in \llbracket 0, n_{01} \rrbracket, l \in \llbracket 0, n_{00} \rrbracket\}. \quad (2.13)$$

We define vectors $\mathbf{v}(i, j, k, l)$ in the following way. Given $(i, j, k, l) \in \mathcal{T}(\mathbf{n})$, we set

$$\mathbf{v}(i, j, k, l) = (i + k, n_{11} - i + l, n_{01} - k + j, n_{10} + n_{00} - j - l). \quad (2.14)$$

To motivate this definition, we think of starting with the observed count vector \mathbf{n} , and imputing various combinations of the unknown potential outcomes. In the n_{11} subjects with observed outcome 1 in the treatment group, we impute 1 for the unknown outcomes in exactly i of them, to get i potential outcome pairs $(1, 1)$. Similarly, j, k, l represent the number of 1's imputed in the other three combinations of group and outcome. It is straightforward to see that the collection of $\mathbf{v}(i, j, k, l)$ such that $(i, j, k, l) \in \mathcal{T}(\mathbf{n})$ is exactly the set of count vectors \mathbf{v} possible given the observed count vector \mathbf{n} . (However, the map $(i, j, k, l) \mapsto \mathbf{v}(i, j, k, l)$ is not necessarily injective.)

Algorithm 2.3. The following algorithm takes as input $\alpha \in (0, 1)$ and a vector $\mathbf{n} \in \mathbb{Z}_{\geq 0}^{\{0,1\}^2}$. It returns an interval $\mathcal{I}_\alpha(\mathbf{n}) = [U_\alpha(\mathbf{n}), L_\alpha(\mathbf{n})]$.

By directly enumerating of all elements of \mathcal{T} , and applying Algorithm 2.2 to each one, compute the set

$$\mathcal{K}(\mathbf{n}) = \left\{ \tau(\mathbf{v}(i, j, k, l)) : (i, j, k, l) \in \mathcal{T}(\mathbf{n}) \text{ and Algorithm 2.2 accepts } \mathbf{v}(i, j, k, l) \right\}. \quad (2.15)$$

Return $U_\alpha(\mathbf{n}) = \max(\mathcal{K})$ and $L_\alpha = \min(\mathcal{K})$.

⁵In fact, the inequality $\geq \alpha$ can be replaced with the inequality $> \alpha$ in (2.12), which creates a stricter threshold for acceptance. This change leads to potentially shorter confidence intervals that still cover at the nominal rate, which can be seen by inspecting the proof of Proposition 3.2. Further, all of our results are still valid under the stricter definition. The definition (2.1) was used in the previous works [RH15, LD16], and we retain it for consistency.

As noted previously, any observed data (\mathbf{Y}, \mathbf{Z}) gives rise to an observed count vector \mathbf{n} . Then using Algorithm 2.3, we define

$$\mathcal{I}_\alpha(\mathbf{Y}, \mathbf{Z}) = [L_\alpha(\mathbf{Y}, \mathbf{Z}), U_\alpha(\mathbf{Y}, \mathbf{Z})] \quad (2.16)$$

by setting $\mathcal{I}_\alpha(\mathbf{Y}, \mathbf{Z}) = \mathcal{I}_\alpha(\mathbf{n})$. Rigdon and Hudgens observe that since $\mathcal{T}(\mathbf{n})$ has size $O(n^4)$, only $O(n^4)$ permutation tests are needed to construct the set $\mathcal{K}(\mathbf{n})$ in Algorithm 2.3. Further, this algorithm does not require our standing hypothesis that the treatment and control groups have an equal number of subjects, and may be applied to arbitrary unbalanced trials.

2.2.2 Li and Ding

As noted in the introduction, Li and Ding have proposed an algorithm that returns $\mathcal{I}_\alpha(\mathbf{n})$ in $O(n^2)$ permutation tests for balanced trials [LD16]. We do not give the details here, since they are not relevant to our work. However, we will make use of the following lemmas from [LD16]. The first provides a necessary and sufficient condition for checking whether a potential outcome table is possible given the observed data.

Lemma 2.4 ([LD16, Theorem 1]). *A potential outcome table \mathbf{w} with count vector \mathbf{v} is possible given observed data \mathbf{n} if and only if*

$$\begin{aligned} & \max(0, n_{11} - v_{10}, v_{11} - n_{01}, v_{11} + v_{01} - n_{10} - n_{01}) \\ & \leq \min(v_{11}, n_{11}, v_{11} + v_{01} - n_{01}, n - v_{10} - n_{01} - n_{10}). \end{aligned} \quad (2.17)$$

The second lemma shows that in the balanced case, the confidence set that we would obtain if we tested *every* possible value of τ_0 is indeed an interval.

Lemma 2.5 ([LD16, Theorem A.4]). *Fix $\alpha \in (0, 1)$ and observed data \mathbf{y} and \mathbf{Z} . For every $\tau_0 \in \mathcal{I}_\alpha(\mathbf{y}, \mathbf{Z}) \cap \mathcal{C}$, there exists a possible potential table \mathbf{w} such that $\tau(\mathbf{w}) = \tau_0$ and $p(\mathbf{w}) \geq \alpha$.*

3 Basic Properties

We now note three basic properties of the confidence intervals \mathcal{I}_α . While the first two have been mentioned informally in previous works, they have not been proved, and we feel there is value in giving precise statements and justifications. The third is new. The proofs appear in Section 7. The arguments for all results in this section extend to the unbalanced case, where the treatment group has size $m = cn$ for $c \in (0, 1)$, with only minor changes. We omit these extensions for brevity.

Our first lemma implies that the interval $\mathcal{I}_\alpha(\mathbf{n})$ always contains the estimate $T(\mathbf{n})$ of $\tau(\mathbf{y})$.

Lemma 3.1. *For any $\mathbf{n} \in \mathbb{Z}_{\geq 0}^{\{0,1\}^2}$, we have*

$$L_\alpha(\mathbf{n}) \leq T(\mathbf{n}) \leq U_\alpha(\mathbf{n}). \quad (3.1)$$

The next proposition states that the intervals \mathcal{I}_α are exact.

Proposition 3.2. *Fix potential outcomes \mathbf{y} and $\alpha \in (0, 1)$. Then*

$$\mathbb{P}(\tau(\mathbf{y}) \in \mathcal{I}_\alpha(\mathbf{Y}, \mathbf{Z})) \geq 1 - \alpha, \quad (3.2)$$

where the probability is with respect to the variable \mathbf{Z} .

Finally, the following proposition states that the intervals \mathcal{I}_α converge at a $n^{-1/2}$ rate. This rate is characteristic of many confidence interval procedures based on central limit theorem asymptotics. We therefore see that the large n scaling behavior of the length of the interval \mathcal{I}_α is the same as that of the intervals produced by asymptotic methods.⁶

Proposition 3.3. *For any $\mathbf{n} \in \mathbb{Z}_{\geq 0}^{\{0,1\}^2}$, we have*

$$|I_\alpha(\mathbf{n})| \leq \sqrt{\frac{32 \log(2/\alpha)}{n}}. \quad (3.3)$$

4 Fast Computation of Confidence Intervals

We now present our main result on the efficient computation of \mathcal{I}_α . Its proof is given in Section 8.2.

Theorem 4.1. *Suppose $n \geq 10$. For any $\alpha \in (0, 1)$ and observed count vector \mathbf{n} , the interval $\mathcal{I}_\alpha(\mathbf{n})$ can be constructed using at most $8n \log_2 n$ permutation tests.*

To prove Theorem 4.1, we exhibit an algorithm that constructs $\mathcal{I}_\alpha(\mathbf{n})$ in the required number of permutation tests.⁷ We begin with a heuristic argument. We first observe that the problem of finding the upper bound $U_\alpha(\mathbf{n})$ reduces to finding a method to determine whether a given value $\tau_0 \in \mathcal{C}$ is incompatible with the observed data \mathbf{n} , in the sense that $p(\mathbf{v}, \mathbf{n}) < \alpha$ for all \mathbf{v} possible given \mathbf{n} such that $\tau(\mathbf{v}) = \tau_0$. If this can be done in $O(n)$ permutation tests, then we can perform a binary search on the set $\left[T(\mathbf{n}), \max(C(\mathbf{n})) \right] \cap \mathcal{S}(n)$ to find U_α by testing each value τ_0 considered by the binary search for compatibility with \mathbf{n} . Since Lemma 2.5 guarantees that every element $\mathcal{I}_\alpha(\mathbf{n})$ is compatible with the observed data, and Lemma 3.1 guarantees that $\mathcal{I}_\alpha(\mathbf{n})$ contains $T(\mathbf{n})$, the binary search will return the maximum of

$$\left[T(\mathbf{n}), \max(C(\mathbf{n})) \right] \cap \mathcal{S}(n) \cap \mathcal{I}_\alpha(\mathbf{n}), \quad (4.1)$$

which is equal to $U_\alpha(\mathbf{n})$. Since binary search on a set of size $O(n)$ can be completed in $O(\log n)$ time, in total $O(n \log n)$ permutation tests are needed. Analogous reasoning applies for finding the lower bound L_α .

Consequently, the proof of Theorem 4.1 focuses on showing that a given $\tau_0 \in \mathcal{C}$ can be checked for compatibility in $O(n)$ permutation tests. Observe that any possible count vector $\mathbf{v} = (v_{11}, v_{10}, v_{01}, v_{00})$ such that $\tau(\mathbf{v}) = \tau_0$ satisfies the equations

$$v_{11} + v_{10} + v_{01} + v_{00} = n, \quad v_{10} - v_{01} = n\tau_0. \quad (4.2)$$

The following lemma allows us to efficiently search the space cut out by (4.2) by exploiting a certain monotonicity property of the permutation p -values defined in (2.12). It is proved in Section 8.3.

Lemma 4.2. *Fix observed data \mathbf{n} , and a count vector $\mathbf{v} = (v_{11}, v_{10}, v_{01}, v_{00})$. Suppose $\min(v_{10}, v_{01}) \geq 1$ and $\max(v_{10}, v_{01}) \geq 2$, and set*

$$\mathbf{v}' = (v_{11} + 1, v_{10} - 1, v_{01} - 1, v_{00} + 1). \quad (4.3)$$

⁶In Proposition 3.3, our focus is on establishing the correct scaling rate, not obtaining the optimal constant prefactor. This estimate is quite conservative compared to the interval lengths we find empirically through simulation. It should not be used in practice for prospective sample size calculations.

⁷Again, the estimate in Theorem 4.1 is conservative relative to the empirical performance of Algorithm 4.3 below. To streamline the arguments, we did not pursue the sharpest possible bounds.

Then $p(\mathbf{v}', \mathbf{n}) \geq p(\mathbf{v}, \mathbf{n})$.

The upshot of this lemma is that if $p(\mathbf{v}, \mathbf{n}) < \alpha$, meaning \mathbf{v} is incompatible with the observed data, then all potential outcome tables on the line segment given by the translations (4.3), with endpoint \mathbf{v} , are also incompatible. Heuristically, the transformation in (4.3) causes the distribution of the variable $T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}}) - \tau(\mathbf{w})$ in the definition (2.12) to become more spread out about its mean (zero), increasing the p -value. This heuristic is made precise in the proof of Lemma 4.2.

Given this context, we now present our main algorithm. The additional check in Step (2c) exists to handle the case where Lemma 4.2 does not apply.

Algorithm 4.3. This algorithm takes as input $\alpha \in (0, 1)$ and a vector $\mathbf{n} \in \mathbb{Z}_{\geq 0}^{\{0,1\}^2}$. It returns an interval $[L_\alpha^{(1)}(\mathbf{n}), U_\alpha^{(1)}(\mathbf{n})]$. We give only the steps for finding $U_\alpha^{(1)}(\mathbf{n})$, since finding $L_\alpha^{(1)}(\mathbf{n})$ is analogous.

We perform a binary search using a function $f: \llbracket nT(\mathbf{n}), n \max(\mathcal{C}(\mathbf{n})) \rrbracket \rightarrow \{0, 1\}$ defined below. The precise definition of this binary search is given as Algorithm 8.1 in Section 8.1. Given an input x , the function f returns 0 if x/n is compatible with the vector \mathbf{n} , and 1 otherwise, where the compatibility of a given $\tau_0 \in \mathcal{C}(\mathbf{n})$ is determined through the following steps.

1. Initialize $j \leftarrow 0$.
2. If $j \leq n$, perform the following steps.
 - (a) Construct a count vector \mathbf{v} as follows (recall (2.2)). Let v_{10} be a free parameter, and set $v_{11} = j - v_{10}$. Then (solving (4.2)) we set

$$v_{01} = v_{10} - n\tau_0, \quad v_{00} = n - j - v_{10} + n\tau_0. \quad (4.4)$$

- (b) Find the smallest v_{10} that leads to a possible table \mathbf{v} given the observed data \mathbf{n} . This can be done in constant time using Lemma 8.3, stated in Section 8.1 below. If no value of v_{10} leads to a possible table, set $j \leftarrow j + 1$ and go to (2).
 - (c) Permutation test the count vector \mathbf{v} coming from the choice of v_{10} in the previous step using Algorithm 2.2 with the given value of α . If \mathbf{v} is accepted, return that τ_0 is compatible. If \mathbf{v} is rejected and $v_{10} \geq 1$ or $v_{01} \geq 1$, set $j \leftarrow j + 1$ and go to (2).
If \mathbf{v} is rejected and $v_{10} = v_{01} = 0$, additionally permutation test the count vector given by $v_{10} = 1$. If it is possible and accepted, terminate the algorithm and declare τ_0 compatible. Otherwise, set $j \leftarrow j + 1$ and go to (2).

3. Return that τ_0 is incompatible.

Our simulations using Algorithm 5.2 show significant improvements over the previous algorithms of Li–Ding and Rigdon–Hudgen, even in small samples [RH15, LD16]. See Table 1 for examples.

5 Monte Carlo Intervals

In Section 4, we constructed the interval \mathcal{I}_α using exact permutation tests (Algorithm 2.2). In this section, we show how to construct confidence intervals using approximate Monte Carlo permutation tests, which are far less computationally demanding.

We begin by defining an approximate permutation test.

\mathbf{n}	Confidence Interval	Algorithm 4.3	[LD16]	[RH15]
(2,6,8,0)	$[-14, -5]$	24	113	189
(6,4,4,6)	$[-4, 10]$	16	308	1225
(8,4,5,7)	$[3, 13]$	26	421	2160

Table 1: 95% confidence intervals for $n\tau(\mathbf{n})$. The second column gives the confidence interval $\mathcal{I}(\mathbf{n})$, after scaling by n . The remaining columns indicate the number of permutation tests required to return \mathcal{I}_α for Algorithm 4.3 and the algorithms from [RH15, LD16]. Data for the last two columns is from [LD16, Table I].

Algorithm 5.1. Take as input $\alpha \in (0, 1)$, $\varepsilon \in (0, \alpha)$, $K \in \mathbb{Z}_{>0}$, a potential outcome table \mathbf{w} , and a vector $\mathbf{n} \in \mathbb{Z}_{\geq 0}^{\{0,1\}^2}$. The algorithm will return a binary decision, either accept or reject.

Let $\mathbf{Z}^{(1)}, \dots, \mathbf{Z}^{(K)}$ be a sequence of independent, identically distributed random vectors, with common distribution equal to that of \mathbf{Z} . Define random variables V_1, \dots, V_K by

$$V_i = \mathbb{1}\left(|T(\mathbf{w}, \mathbf{Z}^{(i)}) - \tau(\mathbf{w})| \geq |T(\mathbf{n}) - \tau(\mathbf{w})|\right), \quad i \in \llbracket 1, K \rrbracket, \quad (5.1)$$

where $T(\mathbf{w}, \mathbf{Z}^{(i)})$ is to be defined to be $T(\tilde{\mathbf{Y}}, \mathbf{Z}^{(i)})$ and $\tilde{\mathbf{Y}}$ is the observed data arising from \mathbf{w} and \mathbf{Z} . Set

$$S = \frac{1}{K} \sum_{i=1}^K V_i. \quad (5.2)$$

Simulate a realization of S by simulating realizations of the K variables $\mathbf{Z}^{(1)}, \dots, \mathbf{Z}^{(K)}$, and constructing the V_i and S accordingly. Return the decision accept if $S + \varepsilon \geq \alpha$, and return the decision reject otherwise.

Algorithm 5.2. This algorithm takes as input $\alpha \in (0, 1)$, $\varepsilon \in (0, \alpha)$, and a vector $\mathbf{n} \in \mathbb{Z}_{\geq 0}^{\{0,1\}^2}$. It returns an interval $\mathcal{I}_{\alpha, \varepsilon, K}(\mathbf{n}) = [L_{\alpha, \varepsilon, K}(\mathbf{n}), U_{\alpha, \varepsilon, K}(\mathbf{n})]$. The algorithm is identical to Section 4, except every use of an exact permutation test (Algorithm 2.2) is replaced with an approximate permutation test with the given parameters (Algorithm 5.1).

Given observed data (\mathbf{Y}, \mathbf{Z}) we define the interval $\mathcal{I}_{\alpha, \varepsilon, K}(\mathbf{Y}, \mathbf{Z})$ to be the interval $\mathcal{I}_{\alpha, \varepsilon, K}(\mathbf{n})$ returned by Algorithm 5.2 for the observed count vector \mathbf{n} corresponding to (\mathbf{Y}, \mathbf{Z}) .

The first part of the next proposition shows that the intervals generated by Algorithm 5.2 are guaranteed to cover $\tau(\mathbf{y})$ with probability $1 - \alpha$, if the parameters of the algorithm are chosen correctly. The second part shows that these intervals can be made to approximate the deterministic intervals \mathcal{I}_α arbitrarily well.

Proposition 5.3. Fix a potential outcome table \mathbf{y} and $\alpha \in (0, 1)$. Suppose $n \geq 0$.

1. Fix $\varepsilon \in (0, \alpha)$. If

$$K \geq \frac{1}{\varepsilon^2} \log \left(\frac{16n \log_2(n)}{\varepsilon} \right), \quad (5.3)$$

then

$$\mathbb{P}(\tau(\mathbf{y}) \in \mathcal{I}_{\alpha - \varepsilon, \varepsilon, K}(\mathbf{Y}, \mathbf{Z})) \geq 1 - \alpha, \quad (5.4)$$

where the probability is with respect to the variable \mathbf{Z} .

2. Fix $\varepsilon \in (0, \alpha/3)$. We have

$$\mathbb{P}(\mathcal{I}_{\alpha-\varepsilon,\varepsilon,K}(\mathbf{Y}, \mathbf{Z}) \subset \mathcal{I}_{\alpha-3\varepsilon}) \geq 1 - 16(n \log_2(n)) \exp(-K\varepsilon^2). \quad (5.5)$$

Proposition 5.3 clarifies two important points. First, if n is held fixed, guaranteeing coverage at level $1 - \alpha$ for the ε -approximate interval $\mathcal{I}_{\alpha-\varepsilon,\varepsilon,K}$ requires a number of Monte Carlo samples K that scales asymptotically as $\varepsilon^{-2} \log(\varepsilon^{-1})$ as ε tends to zero. Second, if ε is fixed, the number of samples K should increase asymptotically as $\log(16n \log_2(n))$ as n grows to guarantee that the Monte Carlo intervals are contained in the exact interval $\mathcal{I}_{\alpha-3\varepsilon}$.

We close with two notes on the practical implementation of Algorithm 5.2. First, it is trivially parallelizable, for example by dividing the computation of the V_i in (5.1) across distinct processors. Second, for clarity we have stated Algorithm 5.2 using a fixed number of K samples for each potential outcome table \mathbf{w} to be tested using Algorithm 5.1. A simple modification is to stop sampling variables V_i and reject \mathbf{w} when the probability that $S + \varepsilon > \alpha$ becomes sufficiently small in Algorithm 5.1. In general, this will require fewer samples than the lower bound for K given in Proposition 5.3, which is essentially a worst-case bound for tables \mathbf{w} such that $p(\mathbf{w}, \mathbf{Y}, \mathbf{Z}) \approx \alpha$. The validity of this adaptive algorithm can be justified by replacing the concentration inequalities used in the proof of Proposition 5.3 with more general ones that account for the correlations induced by the early stopping. Such inequalities are well-known in the bandit algorithm literature; see for example [LS20, Exercise 7.1]. For brevity, we do not give the details here.

6 Missing Data

We now consider a more general situation in which the realized outcomes Y_i may not be observed for some subjects, with no restriction on the distribution of the unobserved Y_i . More precisely, we suppose that instead of observing (\mathbf{Y}, \mathbf{Z}) , the experimenter observes (\mathbf{M}, \mathbf{Z}) , where

$$\mathbf{M} \in \{-1, 0, 1\}^n, \quad M_i = -J_i + (1 - J_i)Y_i, \quad (6.1)$$

and $\mathbf{J} \in \{1, 0\}^n$. Here \mathbf{J} is a random variable which may have an arbitrary distribution, and in particular may depend on \mathbf{Y} and \mathbf{Z} . The indices i such that $J_i = 0$ denote realized outcomes that are observed, while those such that $J_i = 1$ indicate missing data. Hence, we have $M_i = -1$ if the realized outcome for the i -th subject is not observed.

We can construct exact confidence intervals for $\tau(\mathbf{y})$ given (\mathbf{M}, \mathbf{Z}) by considering, roughly speaking, the two extremal imputations of the missing data (leading to the smallest and largest estimates of $\tau(\mathbf{y})$).

Definition 6.1. Given a observed data (\mathbf{M}, \mathbf{Z}) , we define the vector $\mathbf{Y}^{(+)} \in \{0, 1\}^n$ by

$$Y_i^{(+)} = 1 \text{ if } J_i = 1 \text{ and } Z_i = 1, \quad (6.2)$$

$$Y_i^{(+)} = 0 \text{ if } J_i = 1 \text{ and } Z_i = 0, \quad (6.3)$$

$$Y_i^{(+)} = M_i \text{ if } J_i = 0. \quad (6.4)$$

Similarly, we define $\mathbf{Y}^{(-)} \in \{0, 1\}^n$ by

$$Y_i^{(-)} = 0 \text{ if } J_i = 1 \text{ and } Z_i = 1, \quad (6.5)$$

$$Y_i^{(-)} = 1 \text{ if } J_i = 1 \text{ and } Z_i = 0, \quad (6.6)$$

$$Y_i^{(-)} = M_i \text{ if } J_i = 0. \quad (6.7)$$

Finally, for $\alpha \in (0, 1)$, we set

$$\mathcal{I}_\alpha^\circ(\mathbf{M}, \mathbf{Z}) = [L_\alpha(\mathbf{Y}^{(-)}, \mathbf{Z}), U_\alpha(\mathbf{Y}^{(+)}, \mathbf{Z})]. \quad (6.8)$$

The following proposition is proved in Section 10.

Proposition 6.2. *Fix potential outcomes \mathbf{y} and $\alpha \in (0, 1)$. Then*

$$\mathbb{P}(\tau(\mathbf{y}) \in \mathcal{I}_\alpha^\circ(\mathbf{M}, \mathbf{Z})) \geq 1 - \alpha, \quad (6.9)$$

where the probability is with respect to the variable \mathbf{Z} .

Remark 6.3. The intervals \mathcal{I}_α° can be used to lift our hypothesis that n is even, in the following way. Given a group of $2m - 1$ subjects, insert a fictitious subject to create a group of $2m$ subjects, randomize to equal groups, and construct \mathcal{I}_α° by treating the data from the fictitious subject as missing. An argument similar to the one that proves Proposition 6.2 shows that this covers the sample average treatment effect for the original group of $2m - 1$ subjects with probability at least $1 - \alpha$. Such intervals will slightly sacrifice precision relative to the intervals \mathcal{I}_α applied to an unbalanced design with groups of m and $m - 1$. (This loss of precision is straightforwardly quantified and decays with rate n^{-1} .) However, the only provably correct construction of \mathcal{I}_α in the unbalanced setting requires $O(n^4)$ permutation tests (as discussed in Section 1), while \mathcal{I}_α° can be constructed in $O(n \log n)$ permutation tests by Theorem 4.1. Hence, there is a trade-off between precision and computational practicality, which may favor using the \mathcal{I}_α° construction when n is large.

7 Proofs for Section 3

The following definition will be used in the proofs of our results. For $k \in \mathbb{Z}_{>0}$, we denote

$$k^{-1}\mathbb{Z} = \left\{ \frac{j}{k} : j \in \mathbb{Z} \right\}, \quad k^{-1}(2\mathbb{Z}) = \left\{ \frac{2j}{k} : j \in \mathbb{Z} \right\}, \quad k^{-1}(2\mathbb{Z} + 1) = \left\{ \frac{2j + 1}{k} : j \in \mathbb{Z} \right\}. \quad (7.1)$$

Definition 7.1. Let $f: \mathbb{Z} \rightarrow \mathbb{R}_{\geq 0}$ be a probability mass function. We say that f is *symmetric about k_0* for some $k_0 \in 2^{-1}\mathbb{Z}$ if

$$f(k_0 - x) = f(k_0 + x) \text{ for all } x \in 2^{-1}\mathbb{Z} \text{ such that } k_0 + x \in \mathbb{Z}. \quad (7.2)$$

If f is symmetric about k_0 , we say that f is *decreasing away from k_0* if

$$f(k_0 + y) \leq f(k_0 + x) \text{ for all } x, y \in 2^{-1}\mathbb{Z} \text{ such that } y > x \text{ and } k_0 + x, k_0 + y \in \mathbb{Z}. \quad (7.3)$$

We now suppose that $\sum_{k=-\infty}^{\infty} kf(k)$ converges absolutely, so that the mean of the distribution corresponding to f exists. In this case, we say that f is *symmetric-decreasing* (abbreviated SD) if f is symmetric about $k_0 = \sum_{k=-\infty}^{\infty} kf(k)$, and decreasing away from k_0 .

Given $k \in \mathbb{Z}_{>0}$ and $k_0 \in 2^{-1}\mathbb{Z}$, we say that a probability mass function $f: k^{-1}\mathbb{Z} \rightarrow \mathbb{R}_{\geq 0}$ is symmetric about $k^{-1}k_0 \in (2k)^{-1}\mathbb{Z}$ if the function $g: \mathbb{Z} \rightarrow \mathbb{R}_{\geq 0}$ defined by $g(x) = f(kx)$ is symmetric about k_0 . Similarly, we say that f is SD if $g(x) = f(kx)$ is.

7.1 Proof of Lemma 3.1

We begin with two preliminary lemmas.

Lemma 7.2. *For any potential outcome table \mathbf{w} , the pmf of $T(\mathbf{w}, \mathbf{Z})$ is symmetric about $\mathbb{E}[T(\mathbf{w}, \mathbf{Z})]$.*

Proof. Let \mathbf{Z}^\dagger be the random vector defined by $Z_i^\dagger = 1 - Z_i$. Then for any randomization given by \mathbf{Z} , the vector \mathbf{Z}^\dagger represents a randomization that exchanges the control and treatment groups. We have

$$\frac{1}{2} \left(T(\mathbf{w}, \mathbf{Z}) + T(\mathbf{w}, \mathbf{Z}^\dagger) \right) = \frac{1}{2m} \sum_{i=1}^n (w_j(1) - w_j(0)) = \mathbb{E}[T(\mathbf{w}, \mathbf{Z})]. \quad (7.4)$$

Let f be the pmf of $T(\mathbf{w}, \mathbf{Z})$, which is supported on $m^{-1}\mathbb{Z}$. Since \mathbf{Z} and \mathbf{Z}^\dagger have the same distribution, we deduce from (7.4) that

$$f\left(\mathbb{E}[T(\mathbf{w}, \mathbf{Z})] - j\right) = f\left(\mathbb{E}[T(\mathbf{w}, \mathbf{Z})] + j\right) \quad (7.5)$$

for all $j \in n^{-1}\mathbb{Z}$ such that $\mathbb{E}[T(\mathbf{w}, \mathbf{Z})] + j \in m^{-1}\mathbb{Z}$, which completes the proof. \square

Lemma 7.3. *Let $n = 2m$ for some $m \in \mathbb{Z}_{>0}$, and fix observed data \mathbf{Y}, \mathbf{Z} . Then there exists a potential outcome table \mathbf{w} such that $\tau(\mathbf{w}) = T(\mathbf{Y}, \mathbf{Z})$ and \mathbf{w} is possible given \mathbf{Y}, \mathbf{Z} .*

Proof. We have $T(\mathbf{Y}, \mathbf{Z}) = m^{-1}(n_{11} - n_{01})$. Let \mathbf{w} be a potential outcome table with count vector $\mathbf{v} = (0, n_{11} + n_{00}, n_{10} + n_{01}, 0)$. Then using

$$n_{11} + n_{10} = n_{01} + n_{00} = m, \quad (7.6)$$

we compute

$$\tau(\mathbf{w}) = \frac{1}{n}(n_{11} + n_{00} - n_{10} - n_{01}) \quad (7.7)$$

$$= \frac{1}{n}(n_{11} + (m - n_{01}) - (m - n_{11}) - n_{01}) = \frac{2n_{11} - 2n_{01}}{n} = T(\mathbf{Y}, \mathbf{Z}). \quad (7.8)$$

This completes the proof. \square

Proof of Lemma 3.1. Let \mathbf{w} be the potential outcome table provided by Lemma 7.3, so that $\tau(\mathbf{w}) = T(\mathbf{Y}, \mathbf{Z})$ and $\mathcal{I}_\alpha(\mathbf{n}) = \mathcal{I}_\alpha(\mathbf{Y}, \mathbf{Z})$. Then the distribution of $T(\mathbf{w}, \mathbf{Z})$ is symmetric about $T(\mathbf{Y}, \mathbf{Z})$, so $T(\mathbf{Y}, \mathbf{Z})$ is the median of the distribution. Then by (2.12), we have $p(\mathbf{w}, \mathbf{Y}, \mathbf{Z}) = 1$, which shows that $\tau(\mathbf{w})$ will always be accepted in Algorithm 2.3. This completes the proof. \square

7.2 Proof of Proposition 3.2

Proof of Proposition 3.2. The event $\{\tau(\mathbf{y}) \notin \mathcal{I}_\alpha(\mathbf{y}, \mathbf{Z})\}$ is contained in the event where Algorithm 2.3 rejects $\tau(\mathbf{y})$, which is in turn contained in the event where

$$\mathbb{P}\left(|T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}}) - \tau(\mathbf{y})| \geq |T(\mathbf{Y}, \mathbf{Z}) - \tau(\mathbf{y})|\right) < \alpha. \quad (7.9)$$

Here probability is with respect to $\tilde{\mathbf{Z}}$, and we consider the left side of (7.9) as a function of $T(\mathbf{Y}, \mathbf{Z})$. Let x denote the smallest real number such that

$$\mathbb{P}\left(|T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}}) - \tau(\mathbf{y})| \geq x\right) < \alpha, \quad (7.10)$$

which exists because $T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}})$ is a discrete random variable. Then the event where (7.9) occurs is the same as the event where $|T(\mathbf{Y}, \mathbf{Z}) - \tau(\mathbf{y})| \geq x$. By the definition of x , and the fact that (\mathbf{Y}, \mathbf{Z}) and $(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}})$ are identically distributed, this event has probability at most α . This completes the proof. \square

7.3 Proof of Proposition 3.3

Proof of Proposition 3.3. Let \mathcal{P} be the set of all sequences of pairs $((q_i, r_i))_{i=1}^m$ such that $q_i, r_i \in \llbracket 1, n \rrbracket$ for all $i \in \llbracket 1, m \rrbracket$ and

$$\{q_1, q_2, \dots, q_m, r_1, r_2, \dots, r_m\} = \{1, 2, \dots, n\}. \quad (7.11)$$

Observe that every element of $\llbracket 1, n \rrbracket$ appears exactly once in such a sequence as some q_i or r_i .

We now perform complete randomization on the group of n subjects in the following way. First, let $P = ((Q_i, R_i))_{i=1}^m$ be a random variable uniformly distributed on the set \mathcal{P} . Hence, for all $i \in \llbracket 1, m \rrbracket$, Q_i and R_i are random variables taking values in $\llbracket 1, n \rrbracket$. We define a random vector $\tilde{\mathbf{Z}} \in \{0, 1\}^n$ as follows. For each $i \in \llbracket 1, N \rrbracket$, set

$$(\tilde{Z}_{Q_i}, \tilde{Z}_{R_i}) = (1, 0) \text{ if } B_i = 1, \quad (\tilde{Z}_{Q_i}, \tilde{Z}_{R_i}) = (0, 1) \text{ if } B_i = 0. \quad (7.12)$$

One straightforwardly checks that the resulting $\tilde{\mathbf{Z}}$ has a uniform distribution on randomizations of the n subjects to two equal groups. That is, $\tilde{\mathbf{Z}}$ has the same distribution as \mathbf{Z} .

Next, let \mathbf{w} by any potential outcome table. We will condition on P and derive a concentration inequality for the resulting conditional distribution of $T(\mathbf{w}, \tilde{\mathbf{Z}})$. First note that the conditional mean of $T(\mathbf{w}, \tilde{\mathbf{Z}})$ is

$$\mathbb{E}[T(\mathbf{w}, \tilde{\mathbf{Z}}) | P] = \tau(\mathbf{w}).$$

Next, observe that the variables

$$A_i = Z_{Q_i} y_{Q_i}(1) - (1 - Z_{Q_i}) y_{Q_i}(0) + Z_{R_i} y_{R_i}(1) - (1 - Z_{R_i}) y_{R_i}(0) \quad (7.13)$$

defined for $i \in \llbracket 1, m \rrbracket$ are conditionally independent, after conditioning on P , and the conditional distribution of $T(\mathbf{w}, \tilde{\mathbf{Z}})$ is given by

$$\frac{1}{m} \sum_{i=1}^m A_i. \quad (7.14)$$

Observe that we have $A_i \in \{-1, 0, 1\}$.

Now recall that Hoeffding's inequality states that given independent random variables $\{X_i\}_{i=1}^n$ such that $a \leq X \leq b$, we have

$$\mathbb{P}\left(|S_n - \mathbb{E}[S_n]| \geq t\right) \leq 2 \exp\left(-\frac{2t^2}{n(b-a)^2}\right) \quad (7.15)$$

for all $t > 0$, where

$$S_n = X_1 + \dots + X_n. \quad (7.16)$$

We apply Hoeffding's inequality to the sequence $\{m^{-1}A_i\}_{i=1}^m$ with $a = -m^{-1}$ and $b = m^{-1}$, after conditioning on P , yielding

$$\mathbb{P}\left(|T(\mathbf{w}, \tilde{\mathbf{Z}}) - \tau(\mathbf{w})| \geq t \mid P\right) \leq 2 \exp\left(-\frac{t^2 n}{8}\right). \quad (7.17)$$

Because this inequality is true conditional on every realization of \mathcal{P} , we have

$$\mathbb{P}\left(|T(\mathbf{w}, \tilde{\mathbf{Z}}) - \tau(\mathbf{w})| \geq t\right) \leq 2 \exp\left(-\frac{t^2 n}{8}\right). \quad (7.18)$$

Inserting $t = |T(\mathbf{Y}, \mathbf{Z}) - \tau(\mathbf{w})|$, we get

$$\mathbb{P}\left(|T(\mathbf{w}, \tilde{\mathbf{Z}}) - \tau(\mathbf{w})| \geq |T(\mathbf{Y}, \mathbf{Z}) - \tau(\mathbf{w})|\right) \leq 2 \exp\left(-\frac{|T(\mathbf{Y}, \mathbf{Z}) - \tau(\mathbf{w})|^2 n}{8}\right). \quad (7.19)$$

To reject \mathbf{w} , we must have

$$\mathbb{P}\left(|T(\mathbf{w}, \tilde{\mathbf{Z}}) - \tau(\mathbf{w})| \geq |T(\mathbf{Y}, \mathbf{Z}) - \tau(\mathbf{w})|\right) < \alpha. \quad (7.20)$$

We introduce the shorthand

$$\delta = |T(\mathbf{Y}, \mathbf{Z}) - \tau(\mathbf{w})|. \quad (7.21)$$

Hence, we solve for

$$2 \exp\left(-\frac{\delta^2 n}{8}\right) < \alpha, \quad (7.22)$$

giving

$$\delta \geq \sqrt{\frac{8 \log(2/\alpha)}{n}}. \quad (7.23)$$

Therefore the permutation test will reject all values of τ_0 with distance at least

$$\sqrt{\frac{8 \log(2/\alpha)}{n}} \quad (7.24)$$

from $T(\mathbf{Y}, \mathbf{Z})$. This implies the conclusion. \square

8 Proofs for Section 4

8.1 Preliminary Results

For completeness, and to facilitate the proofs of our finite-sample bounds, we begin by specifying the version of binary search that we use.

Algorithm 8.1. This *binary search algorithm* takes as input $k_1, k_2 \in \mathbb{Z}_{>0}$ such that $k_2 > k_1$, and a function $f: \llbracket k_1, k_2 \rrbracket \rightarrow \{0, 1\}$ such that $f(x) = 0$ if $x \leq r$ and $f(x) = 1$ if $x > r$, where $r \in \llbracket k_1 - 1, k_2 \rrbracket$ is unknown. The algorithm returns r through the following steps.

1. Initialize $a \leftarrow k_1$ and $b \leftarrow k_2$.
2. Set $c \leftarrow \lfloor (a + b)/2 \rfloor$ and evaluate $f(c)$. If $f(c) = 0$, set $a \leftarrow c$. If $f(c) = 1$, set $b \leftarrow c$.

3. Repeat the previous step until $b = a + 1$. In this case, return a if $a > k_1$ and $b < k_2$. If $a = k_1$, evaluate $f(k_1)$ and return $k_1 - 1$ if $f(k_1) = 1$, and k_1 if $f(k_1) = 0$. If $b = k_2$, evaluate $f(k_2)$, and return k_2 if $f(k_2) = 0$ and $k_2 - 1$ if $f(k_2) = 1$.

The following fact is well known.

Lemma 8.2. *Algorithm 8.1 terminates in at most $\lfloor \log_2(k_2 - k_1 + 1) \rfloor + 2$ evaluations of f .*

Proof. Without loss of generality, we may suppose that $k_1 = 1$. Further, by setting $f(x) = 1$ for $x > k_2$, we may consider f as a function on $\llbracket 1, 2^d + 1 \rrbracket$, where $d = \lfloor \log_2(k_2) \rfloor + 1$ satisfies $2^d \geq k_2$. Then the number of integers in the interval $[a, b]$ after each evaluation of f follows the sequence $2^{d-1} + 1, \dots, 2^1 + 1, 2$ as the algorithm progresses, with a possible additional evaluation of f if $a = 1$ or $b = 2^d + 1$ at the final step. Then at at worst $d + 1$ evaluations are needed in total, which proves the theorem. \square

The following lemma determines the possible potential outcome tables, given some observed data, in a certain one-parameter family of tables. We remark that it shows Step (2b) of Algorithm 4.3 can be completed in constant time.

Lemma 8.3. *Fix $j, v_{10} \in \mathbb{Z}_{\geq 0}$ and $\tau_0 \in \mathcal{C}$, and consider the potential outcome vector*

$$\mathbf{v} = (j - v_{10}, v_{10}, v_{10} - n\tau_0, n - j - v_{10} + n\tau_0). \quad (8.1)$$

Given observed data \mathbf{n} , a necessary condition for \mathbf{v} to be possible is that all of the following inequalities hold:

$$j \geq n\tau_0 + n_{01}, \quad j \geq n_{11}, \quad n \geq j + n_{10}, \quad n_{11} + n\tau_0 + n_{10} + n_{01} \geq j. \quad (8.2)$$

In this case, the possible values of v_{10} are given by the interval

$$\max(0, n\tau_0, j - n_{11} - n_{01}, n_{11} + n_{01} + n\tau_0 - j) \leq v_{10} \leq \min(j, n_{11} + n_{00}, n_{10} + n_{01} + n\tau_0, n + n\tau_0 - j), \quad (8.3)$$

which may be empty.

Proof of Lemma 8.3. We apply the criterion of Lemma 2.4. The maximum in that statement becomes

$$\max(0, n_{11} - v_{10}, j - v_{10} - n_{01}, j - n\tau_0 - n_{10} - n_{01}). \quad (8.4)$$

Similarly, the minimum becomes

$$\min(j - v_{10}, n_{11}, j - n\tau_0 - n_{01}, n - v_{10} - n_{10} - n_{01}). \quad (8.5)$$

Each argument of the maximum function must be less than each argument of the minimum function. We test the arguments of the maximum from left to right. Starting with 0, we must have

$$j \geq v_{10}, \quad n_{11} \geq 0, \quad j \geq n\tau_0 + n_{01}, \quad n - n_{01} - n_{10} \geq v_{10}, \quad (8.6)$$

where we note that the condition $n_{11} \geq 0$ is always true. Considering $n_{11} - v_{10}$, we get

$$j \geq n_{11}, \quad v_{10} \geq 0, \quad v_{10} \geq n_{11} + n_{01} + n\tau_0 - j, \quad n \geq n_{11} + n_{10} + n_{01}. \quad (8.7)$$

Note that the last condition is always satisfied. Considering $j - v_{10} - n_{01}$, we get

$$n_{10} \geq 0, \quad v_{10} \geq j - n_{11} - n_{01}, \quad v_{10} \geq n\tau_0, \quad n \geq j + n_{10}, \quad (8.8)$$

and the first condition is always true. Finally, considering $j - n\tau_0 - n_{10} - n_{01}$,

$$n\tau_0 + n_{10} + n_{01} \geq v_{10}, \quad n_{11} + n\tau_0 + n_{10} + n_{01} \geq j, \quad n_{10}, \quad n + n\tau_0 - j \geq v_{10}. \quad (8.9)$$

Collecting the previous four displays completes the proof. \square

8.2 Proof of Theorem 4.1

We now prove Theorem 4.1 assuming Lemma 4.2, which is proved below.

Proof of Theorem 4.1. We will show that Algorithm 4.3 returns \mathcal{I}_α in the required number of permutation tests. It suffices to show that Algorithm 4.3 returns $U_\alpha^{(1)}(\mathbf{n}) = U_\alpha(\mathbf{n})$ for the upper bound in $2(n+1)\lceil \log_2(n+1) + 2 \rceil$ permutation tests, since the same estimate holds for the lower bound L_α by analogous reasoning. Then the conclusion follows after using

$$8n \log_2(n) \geq 4(n+1)\lceil \log_2(n+1) + 2 \rceil \quad (8.10)$$

for $n \geq 10$.

Let $C_+ = \max(\mathcal{C}(\mathbf{n}))$, where we recall that $\mathcal{C}(\mathbf{n})$ was defined in (2.10). By Lemma 3.1, we have $U_\alpha(\mathbf{n}) \in [T(\mathbf{n}), C_+]$. Recall that by Lemma 2.5, the set of elements

$$\mathcal{J}_\alpha(\mathbf{n}) = \{\tau_0 \in \mathcal{C}(\mathbf{n}) : \text{there exists a table } \mathbf{u} \text{ such that } \tau(\mathbf{u}) = \tau_0 \text{ and } p(\mathbf{n}, \mathbf{u}) \geq \alpha\} \quad (8.11)$$

is an interval equal to $\mathcal{I}_\alpha(\mathbf{n})$, so that (by Lemma 3.1)

$$[T(\mathbf{n}), C_+] \cap \mathcal{J}_\alpha(\mathbf{n}) = [T(\mathbf{n}), U_\alpha(\mathbf{n})]. \quad (8.12)$$

Then it suffices to show that the f given in Algorithm 4.3 satisfies $f(x) = 0$ if $x/n \in \mathcal{J}_\alpha$ and $f(x) = 1$ otherwise, and that $f(x)$ can be computed in at most $2(n+1)$ permutation tests for every x . In this case, the binary search indicated in Algorithm 4.3 will return $U_\alpha(\mathbf{n})$ in a total $2(n+1)\lceil \log_2(n+1) + 2 \rceil$ permutation tests.

We begin with the claim that $f(x) = 0$ if $x/n \in \mathcal{J}_\alpha$ and $f(x) = 1$ otherwise. Set $\tau_0 = x/n$. If a vector \mathbf{v} permutation tested in Step (2c) satisfies $p(\mathbf{n}, \mathbf{v}) \geq \alpha$, then $\tau_0 = \tau(\mathbf{v})$ lies in $\mathcal{J}_\alpha(\mathbf{n})$, and f correctly declares τ_0 compatible.

Otherwise, if $j = n+1$ and the last step declares τ_0 incompatible, we must show that every possible potential outcome count vector \mathbf{u} with $\tau(\mathbf{u}) = \tau_0$ satisfies $p(\mathbf{n}, \mathbf{u}) < \alpha$. To this end, fix an arbitrary \mathbf{u} such that $\tau(\mathbf{u}) = \tau_0$. If \mathbf{u} was tested for compatibility in Step (2c), then we must have $p(\mathbf{n}, \mathbf{u}) < \alpha$, since the permutation test of \mathbf{u} must have rejected for f to have declared τ_0 incompatible.

Therefore, we may suppose \mathbf{u} was not tested, let $j = u_{11} + u_{01}$, and let v_{01} be the value determined in Step (2b). Let \mathbf{v} be the vector determined by v_{01} through the equalities in (2a). If τ_0 was tested declared incompatible by f , then \mathbf{v} was permutation tested and rejected. Hence $\mathbf{u} \neq \mathbf{v}$ and $p(\mathbf{n}, \mathbf{v}) < \alpha$. Suppose $v_{01} + v_{10} \geq 1$. Using this inequality and Lemma 4.2, we deduce that

$$p(\mathbf{n}, \mathbf{u}) \leq p(\mathbf{n}, \mathbf{v}) < \alpha, \quad (8.13)$$

since $\mathbf{v} = \mathbf{u} + k(1, -1, -1, 1)$ for some $k \in \mathbb{Z}_{>0}$, as desired. Now suppose that $v_{10} = v_{01} = 0$. Let \mathbf{v}' be the vector constructed from $v'_{10} = 1$ using (4.4) and $j = u_{11} + u_{01}$. We claim that \mathbf{v}' is possible, and hence tested and rejected. Indeed, if \mathbf{v}' is not possible given the observed data, then

by Lemma 8.3 we find that $\mathbf{u} \neq \mathbf{v}$ is not possible, since since the set of possible coordinates in (8.3) is an interval. Then $p(\mathbf{v}', \mathbf{n}) < \alpha$, and as before we have $p(\mathbf{n}, \mathbf{u}) \leq p(\mathbf{n}, \mathbf{v}') < \alpha$.

Hence, \mathbf{u} is rejected if \mathbf{v} is, and we have shown that evaluating f provides a valid method for checking the compatibility of τ_0 . Finally, the claim that f can be evaluated using $2(n+1)$ permutation tests follows from the fact that examining a given $j \in \llbracket 0, n \rrbracket$ requires at most two permutation tests, as described in Step (2c). This completes the proof. \square

8.3 Proof of Lemma 4.2

We first provide some intuition for the proof of Lemma 4.2. We will see that the sampling distribution for the Neyman estimator $T(\mathbf{w})$ of a generic table \mathbf{w} is SD on $m^{-1}\mathbb{Z}$. If all tables \mathbf{w} had this property, the proof of Lemma 4.2 would be fairly straightforward (by applying Lemma 8.4 below). However, there are exceptional tables that do not, due to parity issues. Whenever we have a table of the form $\mathbf{v} = (a, 0, 0, b)$ with $a + b = n = 2m$, the estimator $T(\mathbf{w}, \mathbf{Z})$ is supported on $m^{-1}(2\mathbb{Z})$ or $m^{-1}(2\mathbb{Z} + 1)$ instead of $m^{-1}\mathbb{Z}$. Most of our effort goes into showing that this is the only such bad case.

Lemma 8.4. *Fix observed data \mathbf{n} , and a count vector*

$$\mathbf{v} = (v_{11}, v_{10}, v_{01}, v_{00}). \quad (8.14)$$

Suppose $\min(v_{10}, v_{01}) \geq 1$, and define

$$\mathbf{v}' = (v_{11} + 1, v_{10} - 1, v_{01} - 1, v_{00} + 1). \quad (8.15)$$

Set

$$\mathbf{v}^\circ = (v_{11}, v_{10} - 1, v_{01} - 1, v_{00}), \quad (8.16)$$

and suppose that the pmf of $T(\mathbf{v}^\circ, \mathbf{Z})$ is SD on the lattice $(m-1)^{-1}\mathbb{Z}$, where here \mathbf{Z} is uniformly distributed over $\mathcal{Z}(n-2)$. Then

$$p(\mathbf{n}, \mathbf{v}') \geq p(\mathbf{n}, \mathbf{v}). \quad (8.17)$$

Proof. We couple the distributions of $T(\mathbf{v}, \mathbf{Z})$ and $T(\mathbf{v}', \mathbf{Z})$ in the following way. Let \mathbf{w} be a potential outcome table with count vector \mathbf{v} , and label the subjects so that $\mathbf{w}_1 = (1, 0)$ and $\mathbf{w}_2 = (0, 1)$. Let \mathbf{w}' be a potential outcome table such that $\mathbf{w}'_1 = (1, 1)$ and $\mathbf{w}'_2 = (0, 0)$, and $\mathbf{w}'_i = \mathbf{w}_i$ for $i \geq 3$. Then \mathbf{w}' has the count vector \mathbf{v}' . Let $\tilde{\mathbf{Y}}$ be the observed outcome vector for \mathbf{w} , let $\tilde{\mathbf{Y}}'$ be the observed outcome vector for \mathbf{w}' , and recall from (2.12) that

$$p(\mathbf{w}, \mathbf{Y}, \mathbf{Z}) = \mathbb{P} \left(|T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}}) - \tau_0| \geq |T(\mathbf{Y}, \mathbf{Z}) - \tau_0| \right). \quad (8.18)$$

by definition.

We may suppose $T(\mathbf{Y}, \mathbf{Z}) - \tau_0 \neq 0$, otherwise $p(\mathbf{n}, \mathbf{v}') = p(\mathbf{n}, \mathbf{v}) = 1$ and the claim (8.17) is trivially true. Then (8.18) implies

$$p(\mathbf{w}, \mathbf{Y}, \mathbf{Z}) = \mathbb{P} \left(T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}}) - \tau_0 \geq |T(\mathbf{Y}, \mathbf{Z}) - \tau_0| \right) \quad (8.19)$$

$$+ \mathbb{P} \left(T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}}) - \tau_0 \leq -|T(\mathbf{Y}, \mathbf{Z}) - \tau_0| \right). \quad (8.20)$$

To complete the proof, it suffices to show that each of the probabilities (8.19) and (8.20) increases

when \mathbf{w} is replaced \mathbf{w}' . Observe that

$$\mathbb{E}[T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}})] = \tau(\mathbf{w}) = \tau_0, \quad (8.21)$$

and that $T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}})$ is symmetric about its mean, by Lemma 7.2. Then by symmetry, it suffices to show that (8.19) increases when \mathbf{w} is replaced with \mathbf{w}' .

As a preliminary observation, note that $\tau_0 \in \mathcal{C}$ takes on values in the lattice $n^{-1}\mathbb{Z}$, while T takes on values in $m^{-1}\mathbb{Z} = 2n^{-1}\mathbb{Z}$. However, we are considering the probability

$$\mathbb{P}\left(T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}}) \geq |T(\mathbf{Y}, \mathbf{Z}) - \tau_0| + \tau_0\right), \quad (8.22)$$

from (8.19), and one checks that $|T(\mathbf{Y}, \mathbf{Z}) - \tau_0| + \tau_0$ is an element of $m^{-1}\mathbb{Z}$ regardless of the value of τ_0 , since $T \in m^{-1}\mathbb{Z}$.

Now, given $T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}})$ for some fixed realization of $\tilde{\mathbf{Z}}$, we consider what happens to its value as \mathbf{w} changes to \mathbf{w}' , so that $\tilde{\mathbf{Y}}$ is replaced by $\tilde{\mathbf{Y}}'$. There are three cases.

1. $\tilde{Z}_1 = \tilde{Z}_2$. Then the first two subjects are in the same group. The change $(1, 0) \mapsto (1, 1)$ and $(0, 1) \mapsto (0, 0)$ leaves T invariant in this case.
2. $\tilde{Z}_1 = 1$ and $\tilde{Z}_2 = 0$. The change of $(1, 0) \mapsto (1, 1)$ for the first subject does not affect T . The change $(0, 1) \mapsto (0, 0)$ for the second subject increases T by m^{-1} .
3. $\tilde{Z}_1 = 0$ and $\tilde{Z}_2 = 1$. The change of $(1, 0) \mapsto (1, 1)$ for the first subject decreases T by m^{-1} . The change $(0, 1) \mapsto (0, 0)$ for the second subject does not affect T .

In the conclusion, the value of T changes by 0 or $\pm m^{-1}$. Then for $a \in \mathbb{Z}$, we have

$$\mathbb{P}\left(T(\tilde{\mathbf{Y}}', \tilde{\mathbf{Z}}) = \frac{a}{m}\right) = \mathbb{P}\left(T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}}) = \frac{a-1}{m} \text{ and } (\tilde{Z}_1, \tilde{Z}_2) = (1, 0)\right) \quad (8.23)$$

$$+ \mathbb{P}\left(T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}}) = \frac{a}{m} \text{ and } \tilde{Z}_1 = \tilde{Z}_2\right) \quad (8.24)$$

$$+ \mathbb{P}\left(T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}}) = \frac{a+1}{m} \text{ and } (\tilde{Z}_1, \tilde{Z}_2) = (0, 1)\right). \quad (8.25)$$

As noted below (8.20), it suffices to show that

$$\mathbb{P}\left(T(\tilde{\mathbf{Y}}', \tilde{\mathbf{Z}}) \geq |T(\mathbf{Y}, \mathbf{Z}) - \tau_0| + \tau_0\right) \geq \mathbb{P}\left(T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}}) \geq |T(\mathbf{Y}, \mathbf{Z}) - \tau_0| + \tau_0\right), \quad (8.26)$$

since this implies that (8.19) increases if \mathbf{v} is replaced by \mathbf{v}' . Summing (8.23), we get

$$\mathbb{P}\left(T(\tilde{\mathbf{Y}}', \tilde{\mathbf{Z}}) \geq |T(\mathbf{Y}, \mathbf{Z}) - \tau_0| + \tau_0\right) \quad (8.27)$$

$$= \mathbb{P}\left(T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}}) \geq |T(\mathbf{Y}, \mathbf{Z}) - \tau_0| + \tau_0 + m^{-1}\right) \quad (8.28)$$

$$+ \mathbb{P}\left(T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}}) = |T(\mathbf{Y}, \mathbf{Z}) - \tau_0| + \tau_0 \text{ and } \tilde{Z}_1 = \tilde{Z}_2\right) \quad (8.29)$$

$$+ \mathbb{P}\left(T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}}) = |T(\mathbf{Y}, \mathbf{Z}) - \tau_0| + \tau_0 \text{ and } (\tilde{Z}_1, \tilde{Z}_2) = (1, 0)\right) \quad (8.30)$$

$$+ \mathbb{P}\left(T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}}) = |T(\mathbf{Y}, \mathbf{Z}) - \tau_0| + \tau_0 - m^{-1} \text{ and } (\tilde{Z}_1, \tilde{Z}_2) = (1, 0)\right). \quad (8.31)$$

Using the previous equality, we see that (8.26) is equivalent to the statement that

$$\mathbb{P}\left(T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}}) = |T(\mathbf{Y}, \mathbf{Z}) - \tau_0| + \tau_0 \text{ and } \tilde{Z}_1 = \tilde{Z}_2\right) \quad (8.32)$$

$$+ \mathbb{P}\left(T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}}) = |T(\mathbf{Y}, \mathbf{Z}) - \tau_0| + \tau_0 \text{ and } (\tilde{Z}_1, \tilde{Z}_2) = (1, 0)\right) \quad (8.33)$$

$$+ \mathbb{P}\left(T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}}) = |T(\mathbf{Y}, \mathbf{Z}) - \tau_0| + \tau_0 - m^{-1} \text{ and } (\tilde{Z}_1, \tilde{Z}_2) = (1, 0)\right) \quad (8.34)$$

$$\geq \mathbb{P}\left(T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}}) = |T(\mathbf{Y}, \mathbf{Z}) - \tau_0| + \tau_0\right), \quad (8.35)$$

which rearranges to

$$\mathbb{P}\left(T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}}) = |T(\mathbf{Y}, \mathbf{Z}) - \tau_0| - m^{-1} + \tau_0 \text{ and } (\tilde{Z}_1, \tilde{Z}_2) = (1, 0)\right) \quad (8.36)$$

$$\geq \mathbb{P}\left(T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}}) = |T(\mathbf{Y}, \mathbf{Z}) - \tau_0| + \tau_0 \text{ and } (\tilde{Z}_1, \tilde{Z}_2) = (0, 1)\right). \quad (8.37)$$

Let $\tilde{\mathbf{Z}}^\circ = (\tilde{Z}_3^\circ, \dots, \tilde{Z}_n^\circ)$ be uniformly distributed on $\mathcal{Z}(n-2)$ (recall (2.4)), let $\mathbf{w}^\circ = (\mathbf{w}_3, \dots, \mathbf{w}_n)$, and let $\tilde{\mathbf{Y}}^\circ$ be the vector of observed outcomes corresponding to $\tilde{\mathbf{Z}}^\circ$ and \mathbf{w}° . Observe that \mathbf{v}° is the count vector for \mathbf{w}° . Then (8.36) may be rewritten as

$$\mathbb{P}\left(T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}}) = |T(\mathbf{Y}, \mathbf{Z}) - \tau_0| - m^{-1} + \tau_0 \mid (\tilde{Z}_1, \tilde{Z}_2) = (1, 0)\right) \mathbb{P}((\tilde{Z}_1, \tilde{Z}_2) = (1, 0)) \quad (8.38)$$

$$\geq \mathbb{P}\left(T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}}) = |T(\mathbf{Y}, \mathbf{Z}) - \tau_0| + \tau_0 \mid (\tilde{Z}_1, \tilde{Z}_2) = (0, 1)\right) \mathbb{P}((\tilde{Z}_1, \tilde{Z}_2) = (0, 1)), \quad (8.39)$$

and further as

$$\mathbb{P}\left(T(\tilde{\mathbf{Y}}^\circ, \tilde{\mathbf{Z}}^\circ) = A - (m-1)^{-1}\right) \geq \mathbb{P}\left(T(\tilde{\mathbf{Y}}^\circ, \tilde{\mathbf{Z}}^\circ) = A\right), \quad (8.40)$$

where

$$A = \frac{m}{m-1} \cdot \left(|T(\mathbf{Y}, \mathbf{Z}) - \tau_0| + \tau_0\right) \in (m-1)^{-1}\mathbb{Z}. \quad (8.41)$$

In the previous computations, we used the fact that the net contribution to $T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}})$ from the first two subjects is zero, since we assumed that $\mathbf{w}_1 = (1, 0)$ and $\mathbf{w}_2 = (0, 1)$. We also used that the conditional distribution of $(\tilde{Z}_3, \dots, \tilde{Z}_n)$, conditional on either $(\tilde{Z}_1, \tilde{Z}_2) = (1, 0)$ or $(\tilde{Z}_1, \tilde{Z}_2) = (0, 1)$, is uniform on the set $\mathcal{Z}(n-2)$, and that the probabilities of these two events are equal. Observe that

$$A > \frac{m}{m-1} \cdot \tau_0 = \mathbb{E}\left[T(\tilde{\mathbf{Y}}^\circ, \tilde{\mathbf{Z}}^\circ)\right], \quad (8.42)$$

since we assumed earlier that $T(\mathbf{Y}, \mathbf{Z}) - \tau_0 \neq 0$. By hypothesis, the pmf of $T(\tilde{\mathbf{Y}}^\circ, \tilde{\mathbf{Z}}^\circ)$ is SD, since it corresponds to the count vector \mathbf{v}° . This fact, and the fact that A is strictly greater than the mean by (8.42), together imply that (8.40) holds. Since (8.40) is equivalent to (8.26), we see that (8.26) holds, which completes the proof. \square

We next prove some lemmas to set up an induction.

Lemma 8.5. *Let X_1 be SD on $m^{-1}\mathbb{Z}$, and let X_2 be uniformly distributed on $\{0, m^{-1}\}$. Then $X_1 + X_2$ is SD on $m^{-1}\mathbb{Z}$.*

Proof. It suffices to show that $m(X_1 + X_2)$ is SD on \mathbb{Z} . Denote the pdf of mX_1 by f_1 . Then the pdf g of $m(X_1 + X_2)$ is given by $g(x) = \frac{1}{2}(f_1(x) + f_1(x-1))$. Suppose that mX_1 is symmetric

about $a \in \mathbb{Z}$. Then it is clear that $g(x)$ is symmetric about $a + \frac{1}{2}$. To show that $g(x)$ is decreasing, it suffices to show that for any $b \geq a + 1/2$, we have $g(b) \geq g(b + 1)$. We find

$$g(b) = \frac{1}{2}(f_1(b) + f_1(b - 1)) \geq \frac{1}{2}(f_1(b + 1) + f_1(b)) = g(b + 1), \quad (8.43)$$

as desired. In the previous equation, we used $f(b) \geq f(b + 1)$, which is true since f is SD, and $f_1(b - 1) \geq f(b)$. The latter is true by the SD property of f_1 if $b \geq a + 1$. If $b = a + 1/2$, then this follows from the symmetric of f about a . \square

Lemma 8.6. *Fix a potential outcome count vector*

$$\mathbf{v} = (v_{11}, v_{10}, v_{01}, v_{00}), \quad (8.44)$$

and suppose at least one of the following conditions holds.

1. We have $v_{10} \geq 2$, and the pmf for $T(\mathbf{v} - \mathbf{a}, \mathbf{Z}')$ is SD on $(m - 1)^{-1}\mathbb{Z}$ for every choice of

$$\mathbf{a} \in \{(1, 1, 0, 0), (0, 2, 0, 0), (0, 1, 1, 0), (0, 1, 0, 1)\} \quad (8.45)$$

such that $\mathbf{v} - \mathbf{a}$ has nonnegative entries. Here, \mathbf{Z}' denotes a random variable uniformly distributed on $\mathcal{Z}(n - 2)$.

2. We have $v_{01} \geq 2$, and the pmf for $T(\mathbf{v} - \mathbf{a}, \mathbf{Z}')$ is SD on $(m - 1)^{-1}\mathbb{Z}$ for every choice of

$$\mathbf{a} \in \{(1, 0, 1, 0), (0, 1, 1, 0), (0, 0, 2, 0), (0, 0, 1, 1)\} \quad (8.46)$$

such that $\mathbf{v} - \mathbf{a}$ has nonnegative entries.

Then the pmf of $T(\mathbf{v}, \mathbf{Z})$ is SD on $m^{-1}\mathbb{Z}$.

Proof. We prove the claim only for the case $v_{10} \geq 2$, since the proof in the case that $v_{01} \geq 2$ is similar. By relabeling, we may suppose that the first subject has the potential outcome vector $\mathbf{w}_1 = (1, 0)$. We sample \mathbf{Z} in the following way. First, choose subject j from $\{2, \dots, n\}$ uniformly at random and consider the pair $(\mathbf{w}_1, \mathbf{w}_j)$. Then sample (Z_1, Z_j) uniformly at random from the assignments $(Z_1 = 1, Z_j = 0)$ and $(Z_1 = 0, Z_j = 1)$. Finally, assign the remaining $n - 2$ subjects to groups by an independent randomization to equal groups (chosen uniformly at random). We will show that conditional distribution of $T(\mathbf{v}, \mathbf{Z})$, conditional on the choice of j , is SD with a mean of $\mathbb{E}[T(\mathbf{v}, \mathbf{Z})]$. This suffices to prove the theorem, since the unconditional pmf for $T(\mathbf{v}, \mathbf{Z})$ is the weighted sum of such conditional pmfs. There are four cases.

1. $\mathbf{w}_j = (1, 1)$. Then $w_1(1) - w_j(0) = 0$, and $w_j(1) - w_1(0) = 1$, so the distribution of

$$Z_1 w_1(1) + (1 - Z_1) w_1(0) + Z_j w_j(1) + (1 - Z_j) w_j(0) \quad (8.47)$$

is uniform on the set $\{0, 1\}$. Further, the contribution to T from the other $n - 2$ subjects, after the independent randomization is also symmetric decreasing on $m^{-1}\mathbb{Z}$, by assumption. Then the sum of these (conditionally) independent variables is SD with mean $\mathbb{E}[T(\mathbf{v}, \mathbf{Z})]$, by Lemma 8.5.

2. $\mathbf{w}_j = (1, 0)$. One checks that the distribution (8.47) is constant and equal to 1, and that the same reasoning as in the previous case applies, since the sum of a random variable that is SD and the constant 1 is still SD.

3. $\mathbf{w}_j = (0, 1)$. The distribution of (8.47) is constant and equal to 0, and the same reasoning as in the previous point applies.
4. $\mathbf{w}_j = (0, 0)$. The distribution of (8.47) is uniform on $\{0, 1\}$, and we can apply Lemma 8.5.

This completes the proof. \square

Lemma 8.7. *Fix any potential outcome count vector of the form*

$$\mathbf{v} = (v_{11}, 0, 0, v_{00}). \quad (8.48)$$

Then the pmf of $T(\mathbf{v}, \mathbf{Z})$ is SD on $m^{-1}(2\mathbb{Z})$ if v_{11} is even. It is SD on $m^{-1}(2\mathbb{Z} + 1)$ if v_{11} is odd.

Remark 8.8. Note that this statement concerns the lattices $m^{-1}(2\mathbb{Z})$ and $m^{-1}(2\mathbb{Z} + 1)$, instead of the lattice $m^{-1}\mathbb{Z}$, due to the parity issue mentioned in the introduction to this subsection.

Proof. We prove the claim for all even $n \in \mathbb{Z}_{\geq 0}$ by induction. The base case $n = 0$ is trivial. For the induction step, we suppose that the claim holds for $n - 2$ and will show it is true for n . Fix an arbitrary \mathbf{v} of the form (8.48) with $v_{11} + v_{00} = n$. We suppose that $v_{11} \geq 2$. When $v_{11} = 0$, the conclusion is trivial, and when $v_{11} = 1$, we can apply the reasoning below with v_{00} instead of v_{11} .

By relabeling, we may suppose $\mathbf{w}_1 = (1, 1)$. We sample \mathbf{Z} as in the previous proof by choosing a partner \mathbf{w}_j for \mathbf{w}_1 uniformly at random, randomizing the pair \mathbf{w}_1 and \mathbf{w}_j to treatment and control, and independently randomizing the other $n - 2$ subjects to equal groups.

We consider first the case that $\mathbf{w}_2 = (1, 1)$. Then the distribution of

$$Z_1 w_1(1) + (1 - Z_1) w_1(0) + Z_j w_j(1) + (1 - Z_j) w_j(0) \quad (8.49)$$

conditional on the choice of partner and that fact that subjects 1 and 2 are assigned to different groups is constant and equal to 0. This implies we are done by induction, since $v_{11} - 2$, the number of $(1, 1)$ potential outcome vectors in the remaining $n - 2$ subjects, has the same parity as v_{11} , and the sampling distributions of

$$\sum_{i=1}^n Z_i w_i(1) + (1 - Z_i) w_i(0) \quad (8.50)$$

and

$$\sum_{i=3}^n Z_i w_i(1) + (1 - Z_i) w_i(0) \quad (8.51)$$

are then both supported and SD on either $m^{-1}(2\mathbb{Z})$ or $m^{-1}(2\mathbb{Z} + 1)$, as desired.

We next consider the case where $\mathbf{w}_2 = (0, 0)$. Then the conditional distribution of (8.49) is uniform on $\{-1, 1\}$. We are again done by induction, since there are $v_{11} - 1$ potential outcome vectors of the form $(1, 1)$ in the remaining $n - 2$ subjects, and this number has the opposite parity to v_{11} . Reasoning as in the proof of Lemma 8.5, we see that the sums (8.50) and (8.51) are such that one is supported and SD on $m^{-1}(2\mathbb{Z})$ and one is supported and SD on $m^{-1}(2\mathbb{Z} + 1)$, as desired.

Combining the conclusions of these two cases completes the proof. \square

Lemma 8.9. *Let $n = 2m$ for some $m \in \mathbb{Z}_{>0}$. Fix a potential outcome count vector*

$$\mathbf{v} = (v_{11}, 1, 0, v_{00}) \text{ or } (v_{11}, 0, 1, v_{00}). \quad (8.52)$$

Then the pmf of $T(\mathbf{v}, \mathbf{Z})$ is SD on $m^{-1}\mathbb{Z}$.

Proof. We consider only the first case, where $v_{10} = 1$ and $v_{01} = 0$, since the argument for the other case is similar. By relabeling, we may suppose $\mathbf{w}_1 = (1, 0)$, and we resample \mathbf{Z} as in the proof of Lemma 8.6 by first choosing uniformly at random a partner \mathbf{w}_j for \mathbf{w}_1 from the set $\{\mathbf{w}_2, \dots, \mathbf{w}_n\}$, sampling (Z_1, Z_j) uniformly at random from the assignments $(Z_1 = 1, Z_j = 0)$ and $(Z_1 = 0, Z_j = 1)$, and assigning the remaining $n - 2$ subjects to groups by an independent, uniform randomization to equal groups. Let \mathbf{v}° denote the (random) count vector for the potential outcome table given by $\{\mathbf{w}_2, \dots, \mathbf{w}_n\} \setminus \{\mathbf{w}_j\}$, and let \mathbf{Z}' denote a vector uniformly distributed on $\mathcal{Z}(n - 2)$.

We now condition on the choice of \mathbf{w}_j . By Lemma 8.7, the conditional distribution of $T(\mathbf{v}^\circ, \mathbf{Z}')$ for the remaining $n - 2$ people is SD on $(m - 1)^{-1}(2\mathbb{Z})$ or $(m - 1)^{-1}(2\mathbb{Z} + 1)$. Similarly, the conditional distribution of (8.47) for the partnership is always uniform on $\{0, 1\}$, since $\mathbf{w}_j = (0, 0)$ or $\mathbf{w}_j = (1, 1)$. Then the conditional pmf for $T(\mathbf{v}, \mathbf{Z})$ is SD on $m^{-1}\mathbb{Z}$ by Lemma 8.5, with mean independent of the choice of \mathbf{w}_j and equal to the unconditional mean $\mathbb{E}[T(\mathbf{v}, \mathbf{Z})]$. Since this holds for every choice of \mathbf{w}_j , it holds for the unconditional distribution of $T(\mathbf{v}, \mathbf{Z})$, as desired. \square

The proof of the following lemma is somewhat computational, so we defer it to the next subsection.

Lemma 8.10. *Let $n = 2m$ for some $m \in \mathbb{Z}_{>0}$. Fix a potential outcome count vector*

$$\mathbf{v} = (v_{11}, 2, 0, v_{00}). \quad (8.53)$$

Then the pmf of $T(\mathbf{v}, \mathbf{Z})$ is SD on $m^{-1}\mathbb{Z}$.

Corollary 8.11. *Let $n = 2m$ for some $m \in \mathbb{Z}_{>0}$. Fix a potential outcome count vector*

$$\mathbf{v} = (v_{11}, 1, 1, v_{00}) \text{ or } (v_{11}, 0, 2, v_{00}). \quad (8.54)$$

Then the pmf of $T(\mathbf{v}, \mathbf{Z})$ is SD on $m^{-1}\mathbb{Z}$.

Proof. The pmfs of $T(\mathbf{v}, \mathbf{Z})$ for the given vectors are translates of the pmf corresponding to $T(\mathbf{v}, \mathbf{Z})$ for

$$\mathbf{v} = (v_{11}, 2, 0, v_{00}), \quad (8.55)$$

so we are done by the previous lemma. \square

Lemma 8.12. *Let $n = 2m$ for some $m \in \mathbb{Z}_{>0}$. Fix any potential outcome count vector*

$$\mathbf{v} = (v_{11}, v_{10}, v_{01}, v_{00}), \quad (8.56)$$

and suppose $v_{10} + v_{01} \geq 1$. Then the pmf of $T(\mathbf{v}, \mathbf{Z})$ is SD on $m^{-1}\mathbb{Z}$.

Proof. We induct on n . The base case $n = 0$ is trivial. For the induction step, suppose that the claim is true for $n - 2$, and fix any \mathbf{v} corresponding to a potential outcome table with n subjects such that $\min(v_{10}, v_{01}) \geq 1$. If $(v_{10}, v_{01}) = (1, 0)$ or $(v_{10}, v_{01}) = (0, 1)$, we are done by Lemma 8.9. If $v_{01} + v_{01} = 2$, we are done by Lemma 8.10 and Corollary 8.11.

In the case $v_{01} + v_{01} > 2$, we will use Lemma 8.6 to show the claim is true for the given \mathbf{v} . In this case, suppose first that $v_{10} \geq 2$ and $v_{01} \geq 1$. Then by the induction hypothesis, Lemma 8.9, Lemma 8.10, and Corollary 8.11, the first condition given in the statement of Lemma 8.6 holds. The case where $v_{10} \geq 1$ and $v_{01} \geq 2$ is analogous. This completes the proof. \square

Proof of Lemma 4.2. This is an immediate consequence of combining Lemma 8.12 and Lemma 8.4. \square

8.4 Proof of Lemma 8.10

We first recall a formula of Copas [Cop73]. Given a potential outcome table \mathbf{w} and a randomization \mathbf{Z} , we define a *treatment count vector* \mathbf{x} by

$$\mathbf{x} = (x_{11}, x_{10}, x_{01}, x_{00}), \quad x_{ab} = \sum_{j=1}^n Z_j \mathbb{1}\{\mathbf{w}_j = (a, b)\}. \quad (8.57)$$

Note that

$$x_{11} + x_{10} + x_{01} + x_{00} = m, \quad (8.58)$$

if $n = 2m$ and we randomize into equal groups.

Lemma 8.13 ([Cop73]). *Fix a potential outcome count vector \mathbf{v} . For any $s_0, s_1 \in \mathbb{Z}_{\geq 0}$, the probability of observing a treatment count vector \mathbf{x} such that*

$$s_1 = x_{11} + x_{10}, \quad s_0 = (v_{11} - x_{11}) + (v_{01} - x_{01}) \quad (8.59)$$

is

$$p(s_1, s_0) = C_{\mathbf{v}} \sum_{x=-\infty}^{\infty} \binom{v_{11}}{x} \binom{v_{10}}{s_1 - x} \binom{v_{01}}{v_{11} + v_{01} - s_0 - x} \binom{v_{00}}{m - v_{11} - s_1 - v_{01} + s_0 + x}, \quad (8.60)$$

where $C_{\mathbf{v}} > 0$ is a constant depending only on \mathbf{v} .

Proof of Lemma 8.10. We explicitly compute the distribution of $T(\mathbf{v}, \mathbf{Z})$. We first note that direct computation shows that $m\mathbb{E}[T(\mathbf{v}, \mathbf{Z})] = 1$, and we know that the distribution of $T(\mathbf{v}, \mathbf{Z})$ is symmetric about its mean by Lemma 7.2.

Adopt the notation of Lemma 8.13. Then

$$s_1 = x_{11} + x_{10}, \quad s_0 = v_{11} - x_{11}. \quad (8.61)$$

With p defined as in (8.60), and using the assumed form of \mathbf{v} , we have

$$p(s_1, s_0) = C_{\mathbf{v}} \sum_x \binom{v_{11}}{x} \binom{2}{s_1 - x} \binom{0}{v_{11} - s_0 - x} \binom{v_{00}}{m - v_{11} - s_1 - v_{01} + s_0 + x} \quad (8.62)$$

$$= C_{\mathbf{v}} \sum_x \binom{v_{11}}{x} \binom{2}{x_{11} + x_{10} - x} \binom{0}{x_{11} - x} \binom{v_{00}}{m - v_{11} - s_1 + s_0 + x}. \quad (8.63)$$

For this to be nonzero, we must have $x = x_{11}$. Using this and (8.61), we get

$$p(s_1, s_0) = p(x_{11}, x_{10}) = C_{\mathbf{v}} \binom{v_{11}}{x_{11}} \binom{2}{x_{10}} \binom{v_{00}}{m - x_{11} - x_{10}}. \quad (8.64)$$

We abbreviate $j = mT(\mathbf{v}, \mathbf{Z}) - 1$ (centering by subtracting the expectation). Then (using $v_{01} = 0$)

$$j = s_1 - s_0 - 1 = x_{11} + x_{10} - v_{11} + x_{11} - 1 = 2x_{11} + x_{10} - v_{11} - 1. \quad (8.65)$$

This yields

$$x_{11} = \frac{j + v_{11} - x_{10} + 1}{2}. \quad (8.66)$$

We also have

$$\frac{v_{00} + v_{11}}{2} + 1 = m. \quad (8.67)$$

Using (8.66) and (8.67) in (8.64), we get

$$p(x_{11}, x_{10}) = C_{\mathbf{v}} \binom{2}{x_{10}} \binom{v_{11}}{\frac{v_{11}}{2} + \frac{1-x_{10}+j}{2}} \binom{v_{00}}{\frac{v_{00}}{2} + \frac{1-x_{10}-j}{2}}. \quad (8.68)$$

We now consider two cases, depending on the parity of v_{11} .

Case I: v_{11} is even. Then v_{00} is also even, since n is even. There are two subcases, depending on the parity of j . Suppose that j is even. Then parity considerations from (8.66) force $x_{10} = 1$. Then the probability mass function becomes

$$p(j) = 2C_{\mathbf{v}} \binom{v_{11}}{\frac{v_{11}}{2} + \frac{j}{2}} \binom{v_{00}}{\frac{v_{00}}{2} - \frac{j}{2}}. \quad (8.69)$$

When j is odd, both $x_{10} = 0$ and $x_{10} = 2$ are possible, and the pmf is

$$p(j) = C_{\mathbf{v}} \binom{v_{11}}{\frac{v_{11}}{2} + \frac{1+j}{2}} \binom{v_{00}}{\frac{v_{00}}{2} + \frac{1-j}{2}} + C_{\mathbf{v}} \binom{v_{11}}{\frac{v_{11}}{2} + \frac{j-1}{2}} \binom{v_{00}}{\frac{v_{00}}{2} + \frac{-1-j}{2}}. \quad (8.70)$$

Now it is a straightforward computation to show that $p(j)$ is decreasing for $j \geq 0$. We will show that $p(j) \geq p(j+1)$ for all $j \geq 0$. First, suppose that j is even. Then, dividing through by $C_{\mathbf{v}}$, we must show that

$$2 \binom{v_{11}}{\frac{v_{11}}{2} + \frac{j}{2}} \binom{v_{00}}{\frac{v_{00}}{2} - \frac{j}{2}} \geq \binom{v_{11}}{\frac{v_{11}}{2} + \frac{j+2}{2}} \binom{v_{00}}{\frac{v_{00}}{2} - \frac{j}{2}} + \binom{v_{11}}{\frac{v_{11}}{2} + \frac{j}{2}} \binom{v_{00}}{\frac{v_{00}}{2} + \frac{-2-j}{2}}. \quad (8.71)$$

This is clear, since $\binom{n}{k}$ is symmetric and decreasing away from $k = n/2$.

Next, we suppose that j is odd. We want to show

$$\binom{v_{11}}{\frac{v_{11}}{2} + \frac{1+j}{2}} \binom{v_{00}}{\frac{v_{00}}{2} + \frac{1-j}{2}} + \binom{v_{11}}{\frac{v_{11}}{2} + \frac{j-1}{2}} \binom{v_{00}}{\frac{v_{00}}{2} + \frac{-1-j}{2}} \geq 2 \binom{v_{11}}{\frac{v_{11}}{2} + \frac{j+1}{2}} \binom{v_{00}}{\frac{v_{00}}{2} - \frac{j+1}{2}}. \quad (8.72)$$

This is again clear, for the same reason.

Case II: v_{11} is odd. Then v_{11} is odd, since n is even. There are two subcases, depending on the parity of j . Suppose that j is odd. Then parity considerations from (8.66) force $x_{10} = 1$. Then the probability mass function becomes

$$p(j) = 2C_{\mathbf{v}} \binom{v_{11}}{\frac{v_{11}}{2} + \frac{j}{2}} \binom{v_{00}}{\frac{v_{00}}{2} - \frac{j}{2}}. \quad (8.73)$$

When j is even, both $x_{10} = 0$ and $x_{10} = 2$ are possible, and the pmf is

$$p(j) = C_{\mathbf{v}} \binom{v_{11}}{\frac{v_{11}}{2} + \frac{1+j}{2}} \binom{v_{00}}{\frac{v_{00}}{2} + \frac{1-j}{2}} + C_{\mathbf{v}} \binom{v_{11}}{\frac{v_{11}}{2} + \frac{j-1}{2}} \binom{v_{00}}{\frac{v_{00}}{2} + \frac{-1-j}{2}}. \quad (8.74)$$

The same verification as in the previous case shows that $p(j)$ is decreasing for $j \geq 0$. \square

9 Proof of Proposition 5.3

Lemma 9.1. Fix $\varepsilon > 0$ and $K \in \mathbb{Z}_{>0}$. Then

$$\mathbb{P}\left(|p(\mathbf{w}, \mathbf{Y}, \mathbf{Z}) - S| > \varepsilon\right) \leq 2 \exp(-K\varepsilon^2). \quad (9.1)$$

Proof. We have $\mathbb{E}[V_i] = p(\mathbf{w}, \mathbf{Y}, \mathbf{Z})$ and $0 \leq V_i \leq 1$ by definition, so Hoeffding's inequality yields

$$\mathbb{P}\left(\left|K \cdot p(\mathbf{w}, \mathbf{Y}, \mathbf{Z}) - \sum_{i=1}^K V_i\right| > K\varepsilon\right) \leq 2 \exp\left(-\frac{(K\varepsilon)^2}{K}\right) = 2 \exp(-K\varepsilon^2), \quad (9.2)$$

as desired. \square

Lemma 9.2. Fix a potential outcome table \mathbf{y} and $\alpha \in (0, 1)$.

1. With probability at least $1 - 8(n+1)[\log_2(n+1) + 2] \exp(-K\varepsilon^2)$, we have

$$\mathcal{I}_\alpha(\mathbf{Y}, \mathbf{Z}) \subset \mathcal{I}_{\alpha, \varepsilon, K}(\mathbf{Y}, \mathbf{Z}). \quad (9.3)$$

2. With probability at least $1 - 8(n+1)[\log_2(n+1) + 2] \exp(-K\varepsilon^2)$, we have

$$\mathcal{I}_{\alpha, \varepsilon, K}(\mathbf{Y}, \mathbf{Z}) \subset \mathcal{I}_{\alpha - 2\varepsilon}(\mathbf{Y}, \mathbf{Z}). \quad (9.4)$$

Proof. We begin with the first claim. The event where $\mathcal{I}_\alpha(\mathbf{Y}, \mathbf{Z}) \subset \mathcal{I}_{\alpha, \varepsilon, K}(\mathbf{Y}, \mathbf{Z})$ is contained in the event

$$\{U_\alpha(\mathbf{Y}, \mathbf{Z}) \leq U_{\alpha, \varepsilon, K}(\mathbf{Y}, \mathbf{Z})\} \cap \{L_\alpha(\mathbf{Y}, \mathbf{Z}) \geq L_{\alpha, \varepsilon, K}(\mathbf{Y}, \mathbf{Z})\}. \quad (9.5)$$

We begin by bounding the probability that $U_\alpha(\mathbf{Y}, \mathbf{Z}) \leq U_{\alpha, \varepsilon, K}(\mathbf{Y}, \mathbf{Z})$ does not hold. Let N be the total number of permutation tests performed by Algorithm 4.3 in the course of finding $U_\alpha(\mathbf{Y}, \mathbf{Z})$. Letting 0 represent acceptance and 1 represent rejection, we let $\mathbf{a} = (a_i)_{i=1}^N \in \{0, 1\}^N$ denote the sequences of acceptances and rejections made by the permutation tests in Algorithm 4.3. We condition on the observed data (\mathbf{Y}, \mathbf{Z}) , so that \mathbf{a} becomes deterministic. Given this observed data, let M denote the (random) number of approximate permutations performed by Algorithm 5.2 before returning $U_{\alpha, \varepsilon, K}(\mathbf{Y}, \mathbf{Z})$. Let $\mathbf{b} = (b_i)_{i=1}^M \in \{0, 1\}^M$ denote the corresponding (random) sequence of acceptances and rejections. Let J be a random variable denoting the smallest positive integer such that $a_J \neq b_J$. If $a_i = b_i$ for all $i \in \llbracket 1, N \rrbracket$, in which case $U_\alpha = U_{\alpha, \varepsilon, K}$, we set $J = N + 1$ and define $b_J = 0$. Note that J is well defined, since it is not possible for Algorithm 5.2 to terminate unless \mathbf{b} differs from \mathbf{a} at some location, or $M = N$ and $\mathbf{b} = \mathbf{a}$. Further, we have the deterministic bound $J \leq N + 1$.

The event $\{U_\alpha(\mathbf{Y}, \mathbf{Z}) \leq U_{\alpha, \varepsilon, K}(\mathbf{Y}, \mathbf{Z})\}$ is equal to the event that $b_J = 0$. That is, in order that $U_\alpha \leq U_{\alpha, \varepsilon, K}$, at the first time the result of an approximate permutation test in Algorithm 5.2 differs from the result of an exact permutation test, the approximate test must accept. We therefore

compute

$$\mathbb{P}\left(\{U_\alpha(\mathbf{Y}, \mathbf{Z}) \leq U_{\alpha, \varepsilon, K}(\mathbf{Y}, \mathbf{Z})\}^c\right) = \mathbb{P}(b_J = 1) \quad (9.6)$$

$$= \sum_{i=1}^N \mathbb{P}(b_i = 1 \text{ and } J = i) \quad (9.7)$$

$$\leq \sum_{i=1}^N \mathbb{P}(b_i = 1 \text{ and } a_i = 0) \quad (9.8)$$

$$\leq N \cdot 2 \exp(-K\varepsilon^2) \quad (9.9)$$

$$\leq 4(n+1)[\log_2(n+1) + 2] \exp(-K\varepsilon^2). \quad (9.10)$$

To justify (9.9), we note that $a_i = 0$ implies $\mathbf{p}(\mathbf{w}, \mathbf{Y}, \mathbf{Z}) \geq \alpha$ for the table \mathbf{w} considered at the i -th step of Algorithm 4.3. For Algorithm 5.2 to reject \mathbf{w} , it must be rejected by Algorithm 5.1. In the notation of Algorithm 5.2, this happens when $S + \varepsilon < \alpha$. Then by Lemma 9.1, Algorithm 5.1 rejects \mathbf{w} with probability at most $2 \exp(-K\varepsilon^2)$. In (9.10), we used the fact that at most

$$2(n+1)[\log_2(n+1) + 2] \quad (9.11)$$

permutation tests are required in Algorithm 5.2 to find U_α , which was shown in the proof of Theorem 4.1. The bound on the probability that $L_\alpha \geq L_{\alpha, \varepsilon, K}$ fails to hold follows by a nearly identical argument. From a union bound over the complements of the two events in (9.5), we deduce (9.3).

For the proof of the second claim, we keep the notation \mathbf{a} and \mathbf{a} , but now define \mathbf{a} using the interval $\mathcal{I}_{\alpha-2\varepsilon}(\mathbf{Y}, \mathbf{Z})$. The event $\mathcal{I}_{\alpha, \varepsilon, K}(\mathbf{Y}, \mathbf{Z}) \subset \mathcal{I}_{\alpha-2\varepsilon}(\mathbf{Y}, \mathbf{Z})$ is contained in the event

$$\{U_{\alpha-2\varepsilon}(\mathbf{Y}, \mathbf{Z}) \geq U_{\alpha, \varepsilon, K}(\mathbf{Y}, \mathbf{Z})\} \cap \{L_{\alpha-2\varepsilon}(\mathbf{Y}, \mathbf{Z}) \leq L_{\alpha, \varepsilon, K}(\mathbf{Y}, \mathbf{Z})\}. \quad (9.12)$$

The event $\{U_{\alpha-2\varepsilon}(\mathbf{Y}, \mathbf{Z}) \geq U_{\alpha, \varepsilon, K}(\mathbf{Y}, \mathbf{Z})\}$ is equal to the event

$$\{b_J = 1\} \cup \{J = N + 1\}. \quad (9.13)$$

We compute

$$\mathbb{P}\left(\{U_{\alpha-2\varepsilon}(\mathbf{Y}, \mathbf{Z}) \geq U_{\alpha, \varepsilon, K}(\mathbf{Y}, \mathbf{Z})\}^c\right) = \mathbb{P}(b_J = 0 \text{ and } J \leq N) \quad (9.14)$$

$$= \sum_{i=1}^N \mathbb{P}(b_i = 0 \text{ and } J = i) \quad (9.15)$$

$$\leq \sum_{i=1}^N \mathbb{P}(b_i = 0 \text{ and } a_i = 1) \quad (9.16)$$

$$\leq N \cdot 2 \exp(-K\varepsilon^2) \quad (9.17)$$

$$\leq 4(n+1)[\log_2(n+1) + 2] \exp(-K\varepsilon^2). \quad (9.18)$$

To justify (9.17), we note that $a_i = 1$ implies $\mathbf{p}(\mathbf{w}, \mathbf{Y}, \mathbf{Z}) < \alpha - 2\varepsilon$ for the table \mathbf{w} considered at the i -th step of Algorithm 4.3. For Algorithm 5.2 to accept \mathbf{w} , it must be accepted by Algorithm 5.1.

We bound the acceptance probability as

$$\mathbb{P}(S + \varepsilon \geq \alpha) \leq \mathbb{P}(|p(\mathbf{w}, \mathbf{Y}, \mathbf{Z}) - S| > \varepsilon) \leq 2 \exp(-K\varepsilon^2) \quad (9.19)$$

using Lemma 9.1. The estimate for the event $\{L_{\alpha-2\varepsilon}(\mathbf{Y}, \mathbf{Z}) \leq L_{\alpha,\varepsilon,K}(\mathbf{Y}, \mathbf{Z})\}$ is identical, and is shown by essentially the same argument. This completes the proof of (9.4) after using a union bound. \square

Proof of Proposition 5.3. For the first part, it suffices to choose K large enough so that

$$\mathbb{P}(\mathcal{I}_{\alpha-\varepsilon} \subset \mathcal{I}_{\alpha-\varepsilon,\varepsilon,K}) \geq 1 - \varepsilon. \quad (9.20)$$

Then the claim follows from Lemma 9.2 and the fact that

$$16n \log_2(n) \geq 8(n+1)[\log_2(n+1) + 2] \quad (9.21)$$

for $n \geq 10$. The second claim is immediate from Lemma 9.2 and the previous inequality. \square

10 Proof of Proposition 6.2

Proof of Proposition 6.2. First, note that by Proposition 3.2, we have

$$\mathbb{P}(\tau(\mathbf{y}) \in \mathcal{J}_\alpha(\mathbf{M}, \mathbf{Z})) \geq 1 - \alpha, \quad (10.1)$$

where

$$J_\alpha(\mathbf{M}, \mathbf{Z}) = \bigcup_{\mathbf{Y}' \in \mathcal{Y}(\mathbf{M})} \mathcal{I}_\alpha(\mathbf{Y}', \mathbf{Z}), \quad (10.2)$$

and $\mathcal{Y}(\mathbf{M})$ consists of all vectors of realized outcomes \mathbf{Y}' that could arise from the partially observed vector \mathbf{M} . To complete the proof of the theorem it then suffices to show that

$$J_\alpha(\mathbf{M}, \mathbf{Z}) \subset I_\alpha^\circ(\mathbf{M}, \mathbf{Z}), \quad (10.3)$$

or equivalently,

$$I_\alpha^\circ(\mathbf{M}, \mathbf{Z})^c \subset J_\alpha(\mathbf{M}, \mathbf{Z})^c = \bigcap_{\mathbf{Y}' \in \mathcal{Y}(\mathbf{M})} \mathcal{I}_\alpha(\mathbf{Y}', \mathbf{Z})^c. \quad (10.4)$$

Fix some $\tau_0 \in I_\alpha^\circ(\mathbf{M}, \mathbf{Z})^c$. We may suppose that $\tau_0 > U_\alpha(\mathbf{Y}^{(+)}, \mathbf{Z})$, since the analogous argument in the complementary case is similar. Then by definition, for all \mathbf{w} such that $\tau(\mathbf{w}) = \tau_0$, \mathbf{w} is incompatible with the data $(\mathbf{Y}^{(+)}, \mathbf{Z})$. We must show that all such \mathbf{w} are incompatible with all choices of $(\mathbf{Y}', \mathbf{Z})$ with $\mathbf{Y}' \in \mathcal{Y}(\mathbf{M})$. Fix such a \mathbf{w} and a choice of \mathbf{Y}' .

By the definition of $\mathbf{Y}^{(+)}$, and since $T(\mathbf{Y}^{(+)}, \mathbf{Z})$ and $T(\mathbf{Y}', \mathbf{Z})$ are coupled through \mathbf{Z} , we have

$$\tau_0 \geq T(\mathbf{Y}^{(+)}, \mathbf{Z}) \geq T(\mathbf{Y}', \mathbf{Z}), \quad (10.5)$$

where the first inequality follows from Lemma 3.1 and our assumption that $\tau_0 > U_\alpha(\mathbf{Y}^{(+)}, \mathbf{Z})$. Then

$$\mathbb{P}\left(|T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}}) - \tau_0| \geq |T(\mathbf{Y}^{(+)}, \mathbf{Z}) - \tau_0|\right) \leq \mathbb{P}\left(|T(\tilde{\mathbf{Y}}, \tilde{\mathbf{Z}}) - \tau_0| \geq |T(\mathbf{Y}', \mathbf{Z}) - \tau_0|\right) < \alpha, \quad (10.6)$$

so \mathbf{w} is incompatible with $(\mathbf{Y}', \mathbf{Z})$, as desired. \square

References

- [BB19] Zach Branson and Marie-Abèle Bind. Randomization-based inference for Bernoulli trial experiments and implications for observational studies. *Statistical Methods in Medical Research*, 28(5):1378–1398, 2019.
- [Cop73] J. B. Copas. Randomization models for the matched and unmatched 2×2 tables. *Biometrika*, 60(3):467–476, 1973.
- [GLSR04] Robert Greevy, Bo Lu, Jeffrey Silber, and Paul Rosenbaum. Optimal multivariate matching before randomization. *Biostatistics*, 5(2):263–275, 2004.
- [IR15] Guido Imbens and Donald Rubin. *Causal Inference in Statistics, Social, and Biomedical Sciences*. Cambridge University Press, 2015.
- [Kal18] Nathan Kallus. Optimal a priori balance in the design of controlled experiments. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 80(1):85–112, 2018.
- [LD16] Xinran Li and Peng Ding. Exact confidence intervals for the average causal effect on a binary outcome. *Statistics in Medicine*, 35(6):957–960, 2016.
- [LS20] Tor Lattimore and Csaba Szepesvári. *Bandit Algorithms*. Cambridge University Press, 2020.
- [MPS88] Cyrus R. Mehta, Nitin R. Patel, and Pralay Senchaudhuri. Importance sampling for estimating exact probabilities in permutational inference. *Journal of the American Statistical Association*, 83(404):999–1005, 1988.
- [MR12] Kari Lock Morgan and Donald Rubin. Rerandomization to improve covariate balance in experiments. *The Annals of Statistics*, 40(2):1263–1282, 2012.
- [RH15] Joseph Rigdon and Michael Hudgens. Randomization inference for treatment effects on a binary outcome. *Statistics in Medicine*, 34(6):924–935, 2015.
- [Rob88] James Robins. Confidence intervals for causal parameters. *Statistics in Medicine*, 7(7):773–785, 1988.
- [SS80] Thomas Santner and Mark Snell. Small-sample confidence intervals for $p_1 - p_2$ and p_1/p_2 in 2×2 contingency tables. *Journal of the American Statistical Association*, 75(370):386–394, 1980.
- [Was04] Larry Wasserman. *All of Statistics: A Concise Course in Statistical Inference*. Springer, 2004.