

Testing conditional independence under isotonicity

Rohan Hore¹, Jake A. Soloff¹, Rina Foygel Barber¹ and Richard J. Samworth²

¹*Department of Statistics, University of Chicago*

²*Statistical Laboratory, University of Cambridge*

November 10, 2024

Abstract

In this paper, we investigate the problem of testing conditional independence (CI) between two random variables X and Y given Z , under the assumption that X is stochastically increasing in Z . The hardness of testing the CI hypothesis is well documented in the literature. While existing approaches often rely on parametric models, smoothness assumptions, or approximations to the conditional distribution of X given Z and/or Y given Z , our test requires no such knowledge beyond a shape constraint. Inspired by the standard permutation method for unconditional independence testing, our procedure determines the significance of a statistic by randomly swapping the X values within ordered pairs of Z samples. The matched pairs and the test statistic may depend on both Y and Z , providing the analyst with significant flexibility in constructing a powerful test. Our test not only achieves finite-sample Type I error control, but also has non-trivial asymptotic power against alternatives that are not too close to the null models. We validate our theoretical findings through a series of simulations and real data experiments.

1 Introduction

Consider the problem of testing the conditional independence (CI) hypothesis

$$H_0^{\text{CI}} : X \perp\!\!\!\perp Y \mid Z,$$

where X and Y are variables of interest (such as a treatment X and an outcome Y), while Z represents a (potentially high-dimensional) confounder. Our available data consist of a sample $(X_1, Y_1, Z_1), \dots, (X_n, Y_n, Z_n) \stackrel{\text{iid}}{\sim} P$, where P is an unknown distribution on $(X, Y, Z) \in \mathcal{X} \times \mathcal{Y} \times \mathcal{Z}$. Throughout, we write $P_{X|Z}$ and $P_{Y|Z}$ to denote the conditional distributions of X given Z and Y given Z respectively.

In the case where the distribution of Z is continuous, [Shah and Peters \(2020\)](#) established that, without further assumptions, there is no universally valid test of H_0^{CI} that achieves non-trivial power for *any* alternative distribution $P \notin H_0^{\text{CI}}$; see also [Neykov, Balakrishnan and Wasserman \(2021\)](#) and [Kim et al. \(2022\)](#). Existing approaches to testing conditional independence have therefore sought to guarantee validity (Type I error control) over restricted classes of null distributions that impose one of the following additional structures:

- (a) a parametric model, such as joint Gaussianity of (X, Y, Z) , or a Gaussian linear model for $Y \mid (X, Z)$ (Kalisch and Bühlmann, 2007);
- (b) a known or well-estimated conditional distribution $P_{X|Z}$ (Candès et al., 2018; Barber, Candès and Samworth, 2020; Berrett et al., 2020; Niu et al., 2024); or
- (c) smoothness of the conditional distribution $P_{X|Z}$ (Shah and Peters, 2020; Lundborg, Shah and Peters, 2022; Kim et al., 2022; Lundborg et al., 2024+).

1.1 Our contributions

In this work, we introduce a nonparametric structure under which we can test conditional independence: we assume a shape constraint—specifically, a form of stochastic monotonicity—for the conditional distribution of $X \mid Z$. Such a constraint is motivated by several applications, particularly in biomedicine, where for instance incidence of diabetes becomes more prevalent with age (Yan et al., 2023), and left ventricular wall thickness is a known risk factor for hypertrophic cardiomyopathy (O’Mahony et al., 2014). We observe a similar trend in economics, where greater professional experience is often linked to higher salaries. In agriculture, crop yields typically increase with optimal rainfall levels. Drawing insights from these real-world examples, we consider the following constraint:

Assumption 1 (Monotonicity of the conditional distribution $P_{X|Z}$). *Let $\mathcal{X} \subseteq \mathbb{R}$ and let \preceq be a partial order on \mathcal{Z} . We assume X is stochastically increasing in Z , meaning that*

$$\text{if } z \preceq z' \text{ then } \mathbb{P}_P \{X \geq x \mid Z = z\} \leq \mathbb{P}_P \{X \geq x \mid Z = z'\} \text{ for all } x.$$

This assumption does not fall into any of the categories (a)–(c) above. We will often consider the case where the control variable Z is univariate, $\mathcal{Z} \subseteq \mathbb{R}$, under the usual total order \leq . Our framework also allows for multivariate $Z \in \mathbb{R}^d$, in which case the most common partial order is the coordinatewise order.

Our main contribution is to introduce a broad strategy for testing the isotonic conditional independence (ICI) null hypothesis

$$H_0^{\text{ICI}} : X \perp\!\!\!\perp Y \mid Z, \text{ and } P_{X|Z} \text{ satisfies Assumption 1.} \quad (1)$$

Naturally, this test should only be applied in settings where the monotonicity condition of Assumption 1 is well-motivated, so that a rejection of H_0^{ICI} can reasonably be interpreted as evidence that $X \not\perp\!\!\!\perp Y \mid Z$. However, we emphasize again that some additional assumption beyond H_0^{ICI} is essential for any valid test to have non-trivial power at any alternative.

1.2 Background: testing independence

To set the stage for some of the notation and ideas underlying our methodology, we briefly review the simpler framework of permutation testing of the null hypothesis of marginal independence, $X \perp\!\!\!\perp Y$.

Given a joint distribution P on $\mathcal{X} \times \mathcal{Y}$, let $(X_i, Y_i)_{i \in [n]} \stackrel{\text{iid}}{\sim} P$, and write $\mathbf{X} = (X_i)_{i=1}^n$ and $\mathbf{Y} = (Y_i)_{i=1}^n$. We can reframe the problem of testing marginal independence as testing

whether the entries of \mathbf{X} are i.i.d. given \mathbf{Y} . Specifically, permutation tests look for violations of exchangeability of \mathbf{X} given \mathbf{Y} . The general approach proceeds as follows: based on \mathbf{Y} , the analyst chooses any statistic $T : \mathcal{X}^n \rightarrow \mathbb{R}$, with larger values of $T(\mathbf{X})$ indicating evidence against the independence null.¹ Write \mathcal{S}_n for the set of permutations of $[n]$, and for $\sigma \in \mathcal{S}_n$, let $T_\sigma = T(\mathbf{X}^\sigma)$ denote the value of the statistic when the entries of \mathbf{X} are permuted according to σ —that is, $\mathbf{X}^\sigma = (X_{\sigma(1)}, \dots, X_{\sigma(n)})$. Finally, define a p -value

$$p = \frac{1}{n!} \sum_{\sigma \in \mathcal{S}_n} \mathbb{1}\{T_\sigma \geq T\}.$$

This construction produces a valid p -value for *any* choice of test statistic T , and the statistic T can be tailored to have power against certain specific alternatives. Indeed, this strategy has been successfully employed to construct independence tests via nearest neighbour distances and mutual information (Berrett and Samworth, 2019), moment methods (Kim, Balakrishnan and Wasserman, 2022), kernels (Pfister et al., 2018), the Hilbert–Schmidt independence criterion (Albert et al., 2022) and U -statistics (Berrett and Samworth, 2021; Berrett, Kontoyiannis and Samworth, 2021),

In our work, since we are interested in conditional (rather than marginal) independence, we will follow a similar strategy, but will use a restricted class of functions T and a (data-dependent) subgroup of permutations $\sigma \in \mathcal{S}_n$, both of which respect the stochastic monotonicity Assumption 1. Our framework still affords the analyst a great deal of flexibility in designing their test, while controlling Type I error across the more challenging null class H_0^{ICI} .

2 Methodology

In this section we give a general procedure, called the **PairSwap-ICI** test, for testing the isotonic conditional independence null H_0^{ICI} . Intuitively, it is plausible that we should be able to construct powerful tests against some alternatives. For example, if $Z_i \preceq Z_j$, then the shape constraint ensures that $X_i \leq X_j$ should hold *at least* half of the time; if we instead observe $X_i \gg X_j$, then this may be due to the influence of Y . Our test builds on and formalizes this intuition: after observing \mathbf{Y} and \mathbf{Z} , the analyst specifies pairs (i, j) such that $Z_i \preceq Z_j$ and then may use large differences $X_i - X_j$ as evidence against the null.

More formally, based on \mathbf{Y} and \mathbf{Z} , and without access to \mathbf{X} , the analyst chooses:

- (i) A sequence of ordered pairs

$$(i_1, j_1), \dots, (i_L, j_L)$$

of indices in $[n] = \{1, \dots, n\}$, where all $2L$ entries are distinct. We require the pairs to be ordered in the sense that

$$Z_{i_\ell} \preceq Z_{j_\ell} \tag{2}$$

for each $\ell \in [L] = \{1, \dots, L\}$. We refer to such a choice of ordered pairs (with any $L \leq \lfloor n/2 \rfloor$) as a *matching*, and denote the set of all possible matchings in $[n]$ satisfying (2) as $\mathcal{M}_n(\mathbf{Z})$.

¹We emphasize that $T(\mathbf{X})$ is allowed to depend on both \mathbf{X} and \mathbf{Y} —for instance we may define the function as $T(\mathbf{x}) = |\sum_{i=1}^n x_i Y_i|$, though we suppress the dependence on \mathbf{Y} in our notation.

- (ii) A sequence of functions ψ_1, \dots, ψ_L , where each $\psi_\ell : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ satisfies the *anti-monotonicity* property

$$\psi_\ell(x + \Delta, x' - \Delta') - \psi_\ell(x' - \Delta', x + \Delta) \geq \psi_\ell(x, x') - \psi_\ell(x', x), \quad (3)$$

for all $\Delta, \Delta' \geq 0$. An example of a class of functions ψ that satisfy anti-monotonicity is $\psi(x, x') = f(x - x')$ for any anti-symmetric and monotone nondecreasing function f , such as $f(x) = x$ or $f(x) = \text{sign}(x)$.

With these choices in place, our test statistic $T : \mathcal{X}^n \rightarrow \mathbb{R}$ is defined as²

$$T(\mathbf{x}) = \sum_{\ell=1}^L \psi_\ell(x_{i_\ell}, x_{j_\ell}). \quad (4)$$

In order to calibrate the test, the analyst compares the observed test statistic $T = T(\mathbf{X})$ with versions of T where indices within pairs (i_ℓ, j_ℓ) are randomly swapped. Specifically, for $\mathbf{s} \in \{\pm 1\}^L$, define $T_{\mathbf{s}} = T(\mathbf{X}^{\mathbf{s}})$, where $\mathbf{X}^{\mathbf{s}}$ is a swapped version of the data vector \mathbf{X} , with entries

$$\begin{cases} (X_{i_\ell}^{\mathbf{s}}, X_{j_\ell}^{\mathbf{s}}) = (X_{i_\ell}, X_{j_\ell}) & s_\ell = +1, \\ (X_{i_\ell}^{\mathbf{s}}, X_{j_\ell}^{\mathbf{s}}) = (X_{j_\ell}, X_{i_\ell}) & s_\ell = -1. \end{cases}$$

That is, $s_\ell = -1$ indicates that the random variables X_{i_ℓ} and X_{j_ℓ} are swapped, while $s_\ell = +1$ indicates no swap. Informally, the two constraints (2) and (3) ensure that, under the null, each $\psi_\ell(X_{i_\ell}, X_{j_\ell})$ is likely to be no larger than its swapped version, $\psi_\ell(X_{j_\ell}, X_{i_\ell})$ —and thus, the statistic $T = T(\mathbf{X})$ is likely to be no larger than its swapped copies, $T_{\mathbf{s}} = T(\mathbf{X}^{\mathbf{s}})$. If instead $T > T_{\mathbf{s}}$ for many swaps \mathbf{s} , this indicates evidence against the null. To formalize this intuition, we define p -value for PairSwap-ICI test as

$$p := \frac{1}{2^L} \sum_{\mathbf{s} \in \{\pm 1\}^L} \mathbb{1}\{T_{\mathbf{s}} \geq T\}. \quad (5)$$

In words, we are comparing the observed value T of the statistic, against all possible permuted statistic values $T_{\mathbf{s}}$ that we would obtain by swapping indices within matched pairs of our observed data \mathbf{X} .

Our method is quite flexible, in that the analyst may select any pairs (i_ℓ, j_ℓ) subject to monotonicity (2), and any functions ψ_ℓ satisfying (3), to construct a valid test. In particular, they can decide on these aspects of the test *after* exploring the data \mathbf{Y}, \mathbf{Z} , choosing to include a pair (i_ℓ, j_ℓ) if the data observed in \mathbf{Y} indicates that $X_{i_\ell} > X_{j_\ell}$ would be likely under the alternative. However, the quality of the matches (i_ℓ, j_ℓ) and functions ψ_ℓ affects the power of our test; we discuss effective strategies for designing a statistic T in Section 3.

Example: a linear test statistic. Before proceeding, we give a simple example of a test statistic that we might choose to use: consider a linear test statistic,

$$T(\mathbf{x}) = \sum_{i=1}^n \beta_i x_i$$

²Again, as for marginal permutation tests in Section 1.2, here we suppress dependence on \mathbf{Y} and \mathbf{Z} in the notation $T(\mathbf{x})$, even though this statistic does depend on \mathbf{Y} and \mathbf{Z} through the choices of the matched pairs (i_ℓ, j_ℓ) and functions ψ_ℓ .

for some coefficients $\beta_i \in \mathbb{R}$. This function can be used as the test statistic for the PairSwap-ICI test, as long as the coefficients satisfy

$$\beta_{i_\ell} \geq \beta_{j_\ell} \text{ for each } \ell \in [L]. \quad (6)$$

To see why, first note that without loss of generality we can take $\beta_i = 0$ for all $i \in [n] \setminus \{i_1, j_1, \dots, i_L, j_L\}$, i.e., all data points not belonging to any of the L pairs. This is because the indicator $\mathbb{1}\{T_{\mathbf{s}} \geq T\}$, appearing in the computation of the p -value, is invariant to these terms. Next, define

$$\psi_\ell(x, x') = \beta_{i_\ell}x + \beta_{j_\ell}x',$$

which satisfies (3) because for any $\Delta, \Delta' \geq 0$ with $x, x', x + \Delta, x' - \Delta' \in \mathcal{X}$, we have

$$\psi_\ell(x + \Delta, x' - \Delta') - \psi_\ell(x' - \Delta', x + \Delta) - \psi_\ell(x, x') + \psi_\ell(x', x) = (\beta_{i_\ell} - \beta_{j_\ell})(\Delta + \Delta') \geq 0$$

by our assumption that $\beta_{i_\ell} \geq \beta_{j_\ell}$. We then have $T(\mathbf{x})$ equal to the test statistic defined in (4).

We remark that choosing such a test statistic is by no means implying an assumption that the dependence between X and Y follows a linear model—it may be the case that the statistic $T(\mathbf{x}) = \beta^\top \mathbf{x}$ has good power for distinguishing the null from the alternative even if a linear model is only a coarse approximation to the true model.

2.1 Validity

Our first main result is that our method yields a valid test of H_0^{ICI} .

Theorem 1. *Under H_0^{ICI} , the conditional Type I error of the PairSwap-ICI test satisfies $\mathbb{P}\{p \leq \alpha \mid \mathbf{Y}, \mathbf{Z}\} \leq \alpha$ for all $\alpha \in [0, 1]$. In particular, the test enjoys marginal error control: $\mathbb{P}\{p \leq \alpha\} \leq \alpha$ for all α .*

Our proof of Theorem 1 formalizes our intuition at the start of this section, making use of the fact that, under the null, $\psi_\ell(X_{i_\ell}, X_{j_\ell})$ tends to be smaller than its swapped version $\psi_\ell(X_{j_\ell}, X_{i_\ell})$.

Proof of Theorem 1. Our proof is split into three steps. First we derive some deterministic properties of the p -value p . Next, we compare to the *sharp null*, where the pair X_{i_ℓ}, X_{j_ℓ} are identically distributed (rather than stochastically ordered). Finally, we examine the validity of the test under the sharp null.

Step 1: some deterministic properties of the p -value. First, fix $\alpha \in [0, 1]$ and define a function $p : \mathbb{R}^n \rightarrow [0, 1]$ as

$$p(\mathbf{x}) = \frac{1}{2^L} \sum_{\mathbf{s} \in \{\pm 1\}^L} \mathbb{1}\{T(\mathbf{x}^{\mathbf{s}}) \geq T(\mathbf{x})\}, \quad (7)$$

so that the p -value p in (5) can be written as $p = p(\mathbf{X})$. For each $\mathbf{s} \in \{\pm 1\}^L$, we can observe that the value $p(\mathbf{x}^{\mathbf{s}})$ is simply computing the quantile of $T(\mathbf{x}^{\mathbf{s}})$ among all possible swapped

statistics, $(T(\mathbf{x}^{\mathbf{s}'}) : \mathbf{s}' \in \{\pm 1\}^L)$. Consequently, it holds deterministically that

$$\frac{1}{2^L} \sum_{\mathbf{s} \in \{\pm 1\}^L} \mathbb{1} \{p(\mathbf{x}^{\mathbf{s}}) \leq \alpha\} \leq \alpha. \quad (8)$$

(Lemma D.24 in the Appendix verifies this bound, for completeness.)

In addition, we claim that $p(\cdot)$ is monotone in its coordinates, namely, $p(\mathbf{x})$ is nonincreasing in each x_{i_ℓ} , and nondecreasing in each x_{j_ℓ} . To see why, for each $\mathbf{s} \in \{\pm 1\}^L$, we can calculate

$$\begin{aligned} \mathbb{1} \{T(\mathbf{x}^{\mathbf{s}}) \geq T(\mathbf{x})\} &= \mathbb{1} \left\{ \sum_{\ell: s_\ell = +1} \psi_\ell(x_{i_\ell}, x_{j_\ell}) + \sum_{\ell: s_\ell = -1} \psi_\ell(x_{j_\ell}, x_{i_\ell}) \geq \sum_{\ell=1}^L \psi_\ell(x_{i_\ell}, x_{j_\ell}) \right\} \\ &= \mathbb{1} \left\{ \sum_{\ell: s_\ell = -1} (\psi_\ell(x_{i_\ell}, x_{j_\ell}) - \psi_\ell(x_{j_\ell}, x_{i_\ell})) \leq 0 \right\}. \end{aligned}$$

By the anti-monotonicity condition (3) on ψ_ℓ , this function is nonincreasing in each x_{i_ℓ} , and nondecreasing in each x_{j_ℓ} , and therefore the same is true for $p(\mathbf{x})$ as well.

Step 2: compare to the sharp null. For each $i \in [n]$, let $P_i = P_{X_i|Z}(\cdot | Z_i)$ denote the null distribution of X_i (after conditioning on \mathbf{Y}, \mathbf{Z}). By Assumption 1, we know that $P_{i_\ell} \preceq_{\text{st}} P_{j_\ell}$ for each pair $\ell \in [L]$, where \preceq_{st} denotes the stochastic ordering on distributions. Next, we also define distributions \bar{P}_ℓ for each $\ell \in [L]$, given by the mixture

$$\bar{P}_\ell = \frac{1}{2}P_{i_\ell} + \frac{1}{2}P_{j_\ell}.$$

In particular, then,

$$P_{i_\ell} \preceq_{\text{st}} \bar{P}_\ell \preceq_{\text{st}} P_{j_\ell}, \quad \ell \in [L]. \quad (9)$$

We will now compare the observed data values, whose distribution (conditional on \mathbf{Y}, \mathbf{Z}) is given by

$$\mathbf{X} = (X_1, \dots, X_n) \sim P_1 \times \dots \times P_n,$$

against a different distribution,

$$\mathbf{X}_\# = ((X_\#)_1, \dots, (X_\#)_n) \sim (P_\#)_1 \times \dots \times (P_\#)_n,$$

where the distributions $(P_\#)_i$ are defined by setting

$$(P_\#)_{i_\ell} = (P_\#)_{j_\ell} = \bar{P}_\ell$$

for each $\ell \in [L]$ (and, for any index $i \in [n] \setminus \{i_1, j_1, \dots, i_L, j_L\}$ that does not belong to any of the L matched pairs, we simply take $(P_\#)_i = P_i$). We can think of this alternative vector of observations as being drawn from a *sharp null*, because for each pair ℓ , the random variables $(X_\#)_{i_\ell}, (X_\#)_{j_\ell}$ are identically distributed (rather than stochastically ordered, as for X_{i_ℓ}, X_{j_ℓ}). In particular, this implies that for any $\mathbf{s} \in \{\pm 1\}^L$,

$$(\mathbf{X}_\#)^{\mathbf{s}} \stackrel{\text{d}}{=} \mathbf{X}_\# \quad (10)$$

after conditioning on \mathbf{Y}, \mathbf{Z} .

In Step 1, we verified that the function $p(\mathbf{x})$ is nonincreasing in each x_{i_ℓ} , and nondecreasing in each x_{j_ℓ} . In particular, combined with the stochastic ordering (9), this means that $p(\mathbf{X}_\#) \preceq_{\text{st}} p(\mathbf{X})$ (conditional on \mathbf{Y}, \mathbf{Z}). We therefore have

$$\mathbb{P}\{p \leq \alpha \mid \mathbf{Y}, \mathbf{Z}\} = \mathbb{P}\{p(\mathbf{X}) \leq \alpha \mid \mathbf{Y}, \mathbf{Z}\} \leq \mathbb{P}\{p(\mathbf{X}_\#) \leq \alpha \mid \mathbf{Y}, \mathbf{Z}\}.$$

From this point on, then, we only need to verify the validity of the p -value $p(\mathbf{X}_\#)$ computed under the sharp null.

Step 3: validity under the sharp null. For data $\mathbf{X}_\#$ drawn under a sharp null, we have

$$\begin{aligned} \mathbb{P}\{p(\mathbf{X}_\#) \leq \alpha \mid \mathbf{Y}, \mathbf{Z}\} &= \frac{1}{2^L} \sum_{\mathbf{s} \in \{\pm 1\}^L} \mathbb{P}\{p((\mathbf{X}_\#)^{\mathbf{s}}) \leq \alpha \mid \mathbf{Y}, \mathbf{Z}\} \\ &= \mathbb{E} \left[\frac{1}{2^L} \sum_{\mathbf{s} \in \{\pm 1\}^L} \mathbb{1}\{p((\mathbf{X}_\#)^{\mathbf{s}}) \leq \alpha\} \mid \mathbf{Y}, \mathbf{Z} \right] \leq \alpha, \end{aligned}$$

where the first step holds by (10), while the last step holds by the deterministic calculation (8) from Step 1. \square

The p -value constructed in (5) requires computing $T_{\mathbf{s}} = T(\mathbf{X}^{\mathbf{s}})$ for all 2^L values of $\mathbf{s} \in \{-1, 1\}^L$, which may be computationally prohibitive for moderate or large L . In practice, it is common to use a Monte Carlo approximation to the p -value: we sample $\mathbf{s}^{(1)}, \dots, \mathbf{s}^{(M)} \stackrel{\text{iid}}{\sim} \text{Unif}(\{\pm 1\}^L)$, and then compute

$$\hat{p}_M = \frac{1 + \sum_{m=1}^M \mathbb{1}\{T_{\mathbf{s}^{(m)}} \geq T\}}{1 + M}$$

The extra ‘1+’ term appearing in the numerator and denominator is necessary to ensure error control for this Monte Carlo version of our test (Davison and Hinkley, 1997; Phipson and Smyth, 2010); in particular, this correction ensures we cannot have $\hat{p}_M = 0$. The following theorem verifies that this version of the test also controls the Type I error.

Theorem 2. *Fix any $M \in \mathbb{N}$. Under H_0^{ICI} , it holds that $\mathbb{P}\{\hat{p}_M \leq \alpha \mid \mathbf{Y}, \mathbf{Z}\} \leq \alpha$ for all $\alpha \in [0, 1]$, and consequently, $\mathbb{P}\{\hat{p}_M \leq \alpha\} \leq \alpha$.*

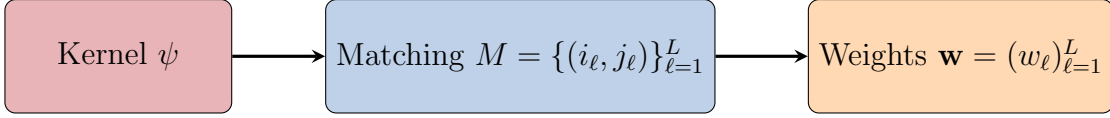
3 Designing a powerful PairSwap-ICI test

In this section, we construct a principled, powerful implementation of our test, by designing concrete choices for the pairs (i_ℓ, j_ℓ) and the functions ψ_ℓ introduced in Section 2. Throughout, we will restrict our attention to test statistics $T(\mathbf{x})$ of the form

$$T(\mathbf{x}) = \sum_{\ell=1}^L w_\ell \psi(x_{i_\ell}, x_{j_\ell}). \tag{11}$$

That is, in the original definition of the test statistic (4), we take $\psi_\ell(\cdot) = w_\ell\psi(\cdot)$ for some sequence of non-negative weights $\mathbf{w} = (w_\ell)_{\ell=1}^L$ and some *fixed* kernel ψ , which is required to satisfy the anti-monotonicity condition (3).

With this simplification, designing a test statistic now requires specifying the kernel ψ , deciding which pairs (i_ℓ, j_ℓ) are matched, and finally, how much weight w_ℓ to assign to each pair, as depicted in this flowchart:



As guaranteed by Theorem 1, our test controls the Type I error for any choice of ψ , M and \mathbf{w} , subject to the conditions (2) and (3) outlined at the start of Section 2. However, for the test to be effective, we need to tailor these choices to the specific application of interest. Of course, all of these choices interact with each other: what constitutes a good matching depends on how we choose the weights, and vice versa.

3.1 Specifying the kernel ψ

We begin by considering several simple options for the kernel ψ . As a first example, consider $\psi(x, x') = x - x'$. This choice of ψ means that $\psi(X_{i_\ell}, X_{j_\ell})$ is likely to be ≤ 0 under the null (since $Z_{i_\ell} \preceq Z_{j_\ell}$), but under the alternative, may be likely to be large (if the pair (i_ℓ, j_ℓ) is chosen wisely). Of course, we also allow for nonlinear test statistics to handle a broader range of settings. If X has heavy tails, then the distribution of a linear statistic T can be very sensitive to extreme values. We can ameliorate this sensitivity by using $\psi(x, x') = \text{sign}(x - x')$, or $\psi(x, x') = (-K) \vee (x - x') \wedge K$ (i.e., the truncation of $x - x'$ to some bounded interval $[-K, K]$) for some constant $K > 0$.

Example: a linear test statistic, revisited. To give more motivation for these simple choices, we will now see that a PairSwap-ICI test run with any *linear* test statistic $T_{\text{lin}}(\mathbf{x}) = \sum_{i=1}^n \beta_i x_i$, can always be expressed in the form (11) with the linear kernel $\psi(x, x') = x - x'$. To see why, define

$$w_\ell = \frac{\beta_{i_\ell} - \beta_{j_\ell}}{2},$$

(and note that we must have $w_\ell \geq 0$ due to (6)). Then we can write

$$T_{\text{lin}}(\mathbf{x}) = T(\mathbf{x}) + T_{\text{sym}}(\mathbf{x}),$$

where T is defined as in (11), and where the term

$$T_{\text{sym}}(\mathbf{x}) = \sum_{\ell=1}^L \frac{\beta_{i_\ell} + \beta_{j_\ell}}{2} (x_{i_\ell} + x_{j_\ell}) + \sum_{i \in [n] \setminus \{i_1, j_1, \dots, i_L, j_L\}} \beta_i x_i$$

is symmetric in the pair (x_{i_ℓ}, x_{j_ℓ}) for each ℓ . Thus $T_{\text{sym}}(\mathbf{x}^s) = T_{\text{sym}}(\mathbf{x})$ for any \mathbf{x} and any $\mathbf{s} \in \{\pm 1\}^L$. It follows that

$$\mathbb{1} \{(T_{\text{lin}})_s \geq T_{\text{lin}}\} = \mathbb{1} \{T_s \geq T\}$$

for every \mathbf{s} —that is, the p -value p defined in (5) is *identical* if we use the test statistic T of the form (11) in place of the original linear test statistic T_{lin} .

3.2 Oracle strategies for choosing the matching and weights

We now build intuition for how to choose the matching M and weights \mathbf{w} effectively by sketching the asymptotics of our test, assuming some oracle knowledge (or estimates) of the data distribution. Let us consider any statistic T of the form (11). Throughout this section, the kernel ψ is a fixed function, and we wish to choose the weights $\mathbf{w} = (w_\ell)_{\ell \in [L]}$ and matching $M = \{(i_\ell, j_\ell)\}_{\ell \in [L]}$ to maximize the power of our test. Given the data $(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$, the reference statistic $T_{\mathbf{s}}$ is a sum of L independent random variables. Under some regularity conditions on the weights \mathbf{w} and the function ψ , a central limit theorem (CLT) approximation gives, for large L , that

$$p \approx \bar{\Phi}(\hat{T}) \quad \text{where} \quad \hat{T} = \hat{T}(\mathbf{w}, M) := \frac{\sum_{\ell=1}^L w_\ell \psi(X_{i_\ell}, X_{j_\ell})}{\sqrt{\sum_{\ell=1}^L w_\ell^2 \psi(X_{i_\ell}, X_{j_\ell})^2}},$$

and where $\bar{\Phi}$ denotes the standard Gaussian survival function, i.e., $\bar{\Phi}(t) = 1 - \Phi(t)$, where Φ is the standard Gaussian distribution function. The above approximation holds for fixed $\mathbf{X}, \mathbf{Y}, \mathbf{Z}$ and relies only on the CLT approximation for a weighted sum of L independent signs $s_1, \dots, s_L \in \{\pm 1\}$.

The above calculation tells us that we should aim to choose weights that *maximize* the approximate probability of rejection, $\mathbb{P}\left\{\bar{\Phi}(\hat{T}) \leq \alpha \mid \mathbf{Y}, \mathbf{Z}\right\}$, in order to achieve the best possible power. Under some conditions this conditional power can be further approximated as

$$\Phi\left(\frac{\sum_{\ell=1}^L w_\ell \mathbb{E}[\psi(X_{i_\ell}, X_{j_\ell}) \mid \mathbf{Y}, \mathbf{Z}]}{\sqrt{\sum_{\ell=1}^L w_\ell^2 \text{Var}(\psi(X_{i_\ell}, X_{j_\ell}) \mid \mathbf{Y}, \mathbf{Z})}} - \bar{\Phi}^{-1}(\alpha)\right) \quad (12)$$

(see Theorem C.12 in the Appendix for a closer look at this approximation). The following lemma shows how to maximize this approximation over the weights, treating the matching M as fixed.

Lemma 3. *Assume that $\text{Var}(\psi_\ell(X_{i_\ell}, X_{j_\ell}) \mid \mathbf{Y}, \mathbf{Z}) > 0$ for $\ell \in [L]$. Considered as a function of $\mathbf{w} = (w_1, \dots, w_L) \in [0, \infty)^L$, the function in (12) is maximised by the choice*

$$w_\ell^* = \frac{\max\{\mathbb{E}[\psi(X_{i_\ell}, X_{j_\ell}) \mid \mathbf{Y}, \mathbf{Z}], 0\}}{\text{Var}(\psi(X_{i_\ell}, X_{j_\ell}) \mid \mathbf{Y}, \mathbf{Z})}, \quad (13)$$

for $\ell \in [L]$.

Using a plug-in estimate for the moments. In Lemma 3, the oracle weight vector \mathbf{w}^* depends on the conditional expected value and variance, $\mathbb{E}[\psi(X_{i_\ell}, X_{j_\ell}) \mid \mathbf{Y}, \mathbf{Z}]$ and $\text{Var}(\psi(X_{i_\ell}, X_{j_\ell}) \mid \mathbf{Y}, \mathbf{Z})$. In practice, we will assume that we have access to estimates \hat{E}_{ij} of $\mathbb{E}[\psi(X_{i_\ell}, X_{j_\ell}) \mid \mathbf{Y}, \mathbf{Z}]$ and $\hat{V}_{ij} > 0$ of $\text{Var}(\psi(X_{i_\ell}, X_{j_\ell}) \mid \mathbf{Y}, \mathbf{Z})$, for each $i, j \in [n]$, constructed from data that are independent of \mathbf{X} . For example, these might be obtained from a fitted

model for the conditional distribution of X given Y, Z constructed using a separate data set (e.g. via sample splitting). For the linear kernel $\psi(x, x') = x - x'$, this would amount to estimating the first two moments of $X_i | Y_i, Z_i$ for $i \in [n]$, whereas if $\psi(x, x') = \text{sign}(x - x')$, then we would need to estimate $\mathbb{P}\{X_i > X_j | Y_i, Z_i, Y_j, Z_j\}$ for distinct indices i, j . With these estimates in place, we can seek to maximize the power of our test by choosing weights

$$\hat{w}_\ell = \frac{\max\{\hat{E}_{i_\ell j_\ell}, 0\}}{\hat{V}_{i_\ell j_\ell}}$$

for $\ell \in [L]$.

Choosing the matching. In the oracle setting, once we fix the choice of weights as in (13), the estimator (12) of the test’s conditional power is maximized by solving a maximum-weight matching problem, namely that of finding

$$M^* \in \operatorname{argmax}_{M \in \mathcal{M}_n(\mathbf{Z})} \sum_{(i,j) \in M} (W_{ij}^*)^2 \quad \text{where} \quad W_{ij}^* = \frac{\max\{\mathbb{E}[\psi(X_i, X_j) | \mathbf{Y}, \mathbf{Z}], 0\}}{\operatorname{Var}(\psi(X_i, X_j) | \mathbf{Y}, \mathbf{Z})}. \quad (14)$$

We then run `PairSwap-ICI` with this oracle matching M^* , and with weights $w_\ell = W_{i_\ell j_\ell}^*$ for each pair $(i_\ell, j_\ell) \in M^*$.

Similarly, using the plug-in estimates, the conditional power is approximately maximized by finding

$$\hat{M} \in \operatorname{argmax}_{M \in \mathcal{M}_n(\mathbf{Z})} \sum_{(i,j) \in M} \hat{W}_{ij}^2 \quad \text{where} \quad \hat{W}_{ij} = \frac{\max\{\hat{E}_{ij}, 0\}}{\hat{V}_{ij}}. \quad (15)$$

To run `PairSwap-ICI`, we then take weights $w_\ell = \hat{W}_{i_\ell j_\ell}$ for each pair $(i_\ell, j_\ell) \in \hat{M}$.

The plug-in matching \hat{M} can be computed in polynomial time (Edmonds, 1965; Duan and Pettie, 2014). Specifically, if $m = |\{(i, j) : Z_i \leq Z_j\}|$, then \hat{M} can be computed in time $O(mn + n^2 \log n)$ using an algorithm of Gabow (1985).³

3.3 Heuristic strategies for choosing the matching and weights

Above, in Section 3.2, the test statistic $T(\mathbf{x})$ was designed with the aim of maximizing the power of our test, but with the assumption that we have access to substantial knowledge about the distribution of the data—for instance, we have been able to fit a model for the distribution of (X, Y, Z) using a separate data set. Of course, an accurate estimate of the true model can enable a very powerful test (since we have a good approximation of the alternative that we are testing against), but in some settings the approach of Section 3.2 may not be practical, either because we are not able to reliably approximate the distribution of the data, or because the required computations are too costly (in particular, the oracle matching strategy).

³For our experiments in Sections 5 and 6, we use the Python package `networkx` (Hagberg, Swart and Schult, 2008), which uses the Blossom algorithm (Edmonds, 1965) and runs in time $O(n^3)$.

In this section, therefore, we take a completely different approach: we will propose an extremely simple scheme for designing weights \mathbf{w} , and two easy strategies for choosing a matching M , that do not require extensive prior knowledge or costly calculations. Of course, this will come at some cost in terms of the resulting power of the test, since we are no longer mimicking an oracle test—but, as we will see in both our theoretical guarantees and our empirical results below, these simple strategies can often attain high power nonetheless.

Throughout this section, we will restrict our attention to the one-dimensional setting, $\mathcal{Z} \subseteq \mathbb{R}$, and will also choose the kernel $\psi(x, x') = x - x'$ when defining our test statistic as in (11). We will work in the setting where we hypothesize that, under the alternative, there is a *positive* association between X and Y even after controlling for Z (of course, if our hypothesis is a negative association, we can follow an analogous strategy). The idea is simple: we will take pairs (i_ℓ, j_ℓ) such that

- $Z_{i_\ell} \leq Z_{j_\ell}$ (as required for validity), but $Z_{i_\ell} \approx Z_{j_\ell}$; and
- $Y_{i_\ell} > Y_{j_\ell}$ (so that, under the alternative, we expect $X_{i_\ell} > X_{j_\ell}$).

3.3.1 A simple weighting scheme

We begin by defining a mechanism for choosing the weights: we will take

$$w_\ell = (Y_{i_\ell} - Y_{j_\ell})_+,$$

the difference in Y values (if this difference is positive for the pair (i_ℓ, j_ℓ)).

Why is this simple strategy a reasonable choice across broad settings? Our test statistic $T(\mathbf{X})$ is given by

$$T(\mathbf{X}) = \sum_{\ell: Y_{i_\ell} > Y_{j_\ell}} (Y_{i_\ell} - Y_{j_\ell}) \cdot (X_{i_\ell} - X_{j_\ell}).$$

Since under the alternative, we expect $X_{i_\ell} > X_{j_\ell}$, this means that the expected value of $T(\mathbf{X})$ is large and positive under the alternative (but, of course, is non-positive under the null).

To consider another motivation, let us examine a specific model. Suppose that Y has a linear effect on the mean of X , so that

$$\mathbb{E}[X \mid Y = y, Z = z] = \beta^* y + \mu_Z^*(z),$$

and the conditional variance is constant,

$$\text{Var}(X \mid Y = y, Z = z) = \sigma^{*2}.$$

Then, following the oracle strategy of Section 3.2, as in (13) we have oracle weights

$$\begin{aligned} w_\ell^* &= \frac{\max\{\mathbb{E}[X_{i_\ell} - X_{j_\ell} \mid \mathbf{Y}, \mathbf{Z}], 0\}}{\text{Var}(X_{i_\ell} - X_{j_\ell} \mid \mathbf{Y}, \mathbf{Z})} = \frac{\max\{\beta^*(Y_{i_\ell} - Y_{j_\ell}) + (\mu_Z^*(Z_{i_\ell}) - \mu_Z^*(Z_{j_\ell})), 0\}}{2\sigma^{*2}} \\ &\approx \frac{\beta^*}{2\sigma^{*2}} \max\{(Y_{i_\ell} - Y_{j_\ell}), 0\}, \end{aligned}$$

where the last step holds since we have assumed $Z_{i_\ell} \approx Z_{j_\ell}$ in our choice of the pair (and so $\mu_Z^*(Z_{i_\ell}) \approx \mu_Z^*(Z_{j_\ell})$ likely holds). But crucially, our test is invariant to rescaling the weights—that is, choosing weights $w_\ell = \max\{Y_{i_\ell} - Y_{j_\ell}, 0\}$ is equivalent to choosing weights $\frac{\beta^*}{2\sigma^{*2}} \max\{(Y_{i_\ell} - Y_{j_\ell}), 0\}$, and thus is nearly equivalent to the oracle weights.

3.3.2 Two simple matching schemes

Next, we propose two matching strategies that, again, do not require any knowledge or estimate of the model.

Neighbour matching. Our first simple matching strategy is to choose pairs that are nearest neighbours in the list of sorted Z values.

Algorithm 1 neighbour matching

Preliminaries: sort Z values, i.e., find a permutation π of $\{1, \dots, n\}$ such that

$$Z_{\pi(1)} \leq Z_{\pi(2)} \leq \dots \leq Z_{\pi(n-1)} \leq Z_{\pi(n)}.$$

for $m = 1, \dots, \lfloor n/2 \rfloor$ **do**

 If $Y_{\pi(2m-1)} > Y_{\pi(2m)}$, then add a new pair to the matching,

$$(i_\ell, j_\ell) = (\pi(2m - 1), \pi(2m)).$$

end for

This strategy ensures that, for all pairs $(i_\ell, j_\ell) = (\pi(2m - 1), \pi(2m))$ that are included in the matching, we have $Z_{i_\ell} \leq Z_{j_\ell}$ and $Y_{i_\ell} > Y_{j_\ell}$ by construction, and moreover, it likely holds that $Z_{i_\ell} \approx Z_{j_\ell}$ (since we have chosen two consecutive Z values in the sorted list).

However, an obvious limitation of this naïve matching strategy is that many consecutive pairs $(\pi(2m - 1), \pi(2m))$ in the sorted list can fail to have $Y_{\pi(2m-1)} > Y_{\pi(2m)}$ just by chance. In particular, if the Z values in this pair are approximately equal (as we might expect), then the values $Y_{\pi(2m-1)}, Y_{\pi(2m)}$ are expected to be approximately i.i.d.—which means that the event $Y_{\pi(2m-1)} > Y_{\pi(2m)}$ will fail half the time. In other words, we are discarding approximately half of the data—we will expect to have $L \approx n/4$ pairs in this matching (meaning that $2L \approx n/2$ many data points have been assigned to a matched pair).

Cross-bin matching. Our next strategy is cross-bin matching, which aims to avoid the inefficiency of neighbour matching—we aim to use (nearly) all of the data, rather than discarding half the data as is likely the case for neighbour matching. To implement this strategy, we will partition the list of sorted Z values into K bins, and will allow a pair of data points to be matched as long as the Z values are in adjacent bins (rather than requiring consecutive Z values, as for neighbour matching).

Algorithm 2 Cross-bin matching

Preliminaries: sort Z values, i.e., find a permutation π of $\{1, \dots, n\}$ such that

$$Z_{\pi(1)} \leq Z_{\pi(2)} \leq \dots \leq Z_{\pi(n-1)} \leq Z_{\pi(n)},$$

and define K bins of indices,

$$\begin{aligned} A_1 &= \{\pi(1), \dots, \pi(m)\}, \\ A_2 &= \{\pi(m+1), \dots, \pi(2m)\}, \\ &\dots \\ A_K &= \{\pi((K-1)m+1), \dots, \pi(Km)\}, \end{aligned}$$

where $m = \lfloor n/K \rfloor$.

for $k = 1, \dots, K$ **do**

Define $r_{k,1}, \dots, r_{k,m}$ as a permutation of A_k such that

$$Y_{r_{k,1}} \geq \dots \geq Y_{r_{k,m}}.$$

end for

for $k = 1, \dots, K-1$ **do**

for $s = 1, \dots, \lfloor m/2 \rfloor$ **do**

If $Y_{r_{k,s}} > Y_{r_{k+1,m+1-s}}$, then add a new pair to the matching,

$$(i_\ell, j_\ell) = (r_{k,s}, r_{k+1,m+1-s}).$$

end for

end for

To explain this procedure in words:

- First we group the Z values into K many bins, with each bin A_k containing m many consecutive Z values.
- Then, we attempt to match the largest values of Y in bin k with the smallest values of Y in bin $k+1$ —that is, in the inner “for loop”, at step $s = 1$ we are attempting to match the largest Y value in bin k (i.e., $Y_{r_{k,1}}$) with the smallest Y value in bin $k+1$ (i.e., $Y_{r_{k+1,m}}$), and then at step $s = 2$ we proceed to matching the second-largest and second-smallest, and so on.

This scheme is illustrated in Figure 1.

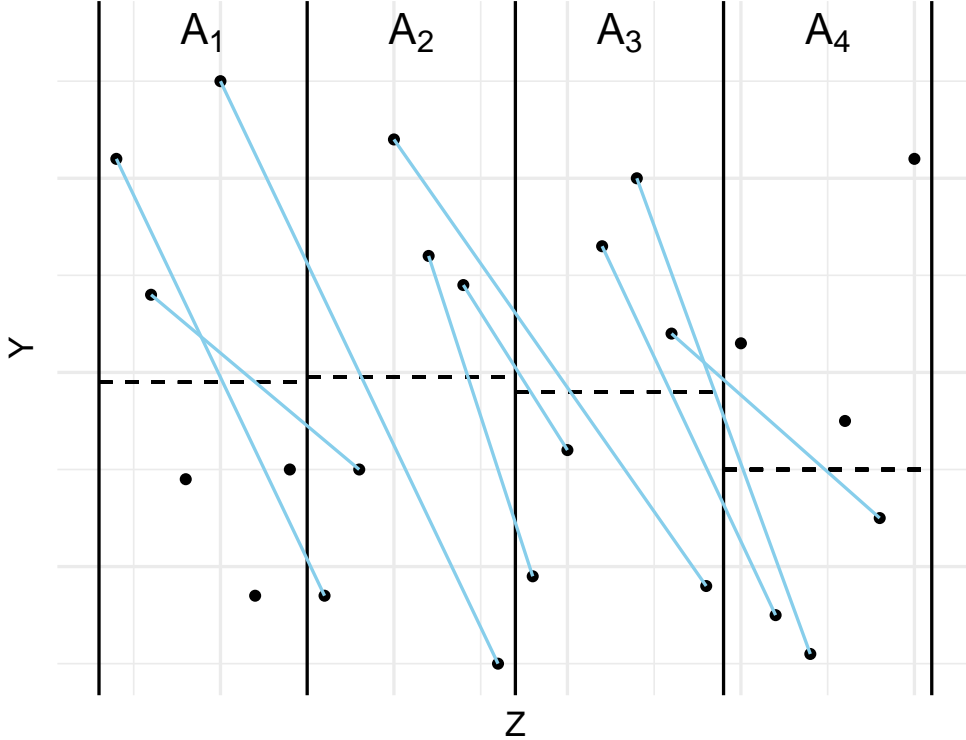


Figure 1: Demonstration of the cross-bin matching scheme described in Algorithm 2.

Why do we expect that this strategy will be more powerful than the simpler neighbour matching strategy? At a high level, while neighbour matching is expected to discard around half of the data, the cross-bin matching strategy can potentially assign nearly all data points to a matched pair. However, there is a potential tradeoff: while the pairs produced by both strategies will likely satisfy $Z_{i_\ell} \approx Z_{j_\ell}$, this approximation will be closer to equality for neighbour matching (where pairs consist of consecutive Z values) than cross-bin matching (where pairs consist of Z values in neighbouring bins, i.e., they may be up to $2m$ positions apart in the sorted list). If m is not too large (i.e., the number of bins K is not too small), though, this difference is hopefully negligible. We will examine both methods theoretically in the following section and will see these tradeoffs in more detail.

4 Power analysis

In this section, we study the power of PairSwap-ICI test under the following general model. We assume that the data $(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) = (X_i, Y_i, Z_i)_{i \in [n]} \stackrel{\text{iid}}{\sim} P$ is drawn according to

$$\mathbf{X} = \mu(\mathbf{Y}, \mathbf{Z}) + \boldsymbol{\zeta}, \quad (16)$$

where $\mu : \mathcal{Y} \times \mathcal{Z} \rightarrow \mathbb{R}$ (applied componentwise) is a measurable function, and with $(Y_i, Z_i)_{i \in [n]} \stackrel{\text{iid}}{\sim} P_{Y,Z}$ drawn independently from $(\zeta_i)_{i \in [n]} \stackrel{\text{iid}}{\sim} P_\zeta$. We suppose that P_ζ has mean 0 and unknown variance $\sigma^2 > 0$. Throughout this section, we also assume that the statistic T admits the form in (11) with $\psi(x, x') = x - x'$, and we restrict our attention to the setting $\mathcal{Y} = \mathcal{Z} = \mathbb{R}$; the partial ordering \preceq for Z will simply be the usual ordering \leq on \mathbb{R} .

Defining isotonic signal strength (ISS). Before stating the results on power under this signal plus noise model, we first introduce some notation. We denote by \mathcal{C}_{ISO} the set of non-decreasing functions on \mathbb{R} and then, we define

$$\text{ISS}_n = \inf_{g \in \mathcal{C}_{\text{ISO}}} \mathbb{E}_{P_{Y,Z}^n} [\|\mu(\mathbf{Y}, \mathbf{Z}) - g(\mathbf{Z})\|_2], \quad \widehat{\text{ISS}}_n = \inf_{g \in \mathcal{C}_{\text{ISO}}} \|\mu(\mathbf{Y}, \mathbf{Z}) - g(\mathbf{Z})\|_2, \quad (17)$$

referring to them as the oracle and empirical *isotonic signal strength* respectively (here $g(\mathbf{Z})$ is applied componentwise). In particular, under the null, we would have $\text{ISS}_n = \widehat{\text{ISS}}_n = 0$ (since we can simply take g to be the true mean function, which does not depend on Y), while under the alternative these quantities might be large.

Preview of results. To analyze the performance of PairSwap-ICI testing procedure, we study the power conditional on \mathbf{Y} and \mathbf{Z} , i.e. $\mathbb{P}\{p \leq \alpha \mid \mathbf{Y}, \mathbf{Z}\}$, and we demonstrate how this power is characterized by the isotonic signal strength quantities ISS_n and $\widehat{\text{ISS}}_n$. The organization of this section is as follows.

- In Section 4.1 we study the asymptotic upper and lower bound on the power of oracle matching (as defined in Section 3.2), conditional on \mathbf{Y} and \mathbf{Z} , and establish that the empirical quantity $\widehat{\text{ISS}}_n$ governs their behaviour. The same guarantees also hold for the plug-in matching strategy defined in Section 3.2, as long as the estimate $\hat{\mu}$ is consistent. Qualitatively, this implies that we can have non-trivial power guarantees against the alternatives with large isotonic signal strength.
- Next, in Section 4.2, we show that the converse also holds, i.e., one can not distinguish the null class from the alternatives with small isotonic signal strength. More precisely, if the oracle quantity ISS_n is too small, then no valid testing procedure for H_0^{ICI} can have non-trivial power.
- Finally, in Section 4.3, we specialize our power guarantees to the special case of partially linear Gaussian models. We show that even without the knowledge of μ or an approximation for the same, we can achieve near-optimal power guarantees with some of the more practical matching algorithms from Section 3.3.

4.1 ISS dictates the power of oracle matching

Below, we study the best-case performance of PairSwap-ICI procedure, specifically the asymptotic behaviour of the conditional power for oracle matching. In particular Theorem 4 provides asymptotic upper and lower bounds on the conditional power $\mathbb{P}\{p \leq \alpha \mid \mathbf{Y}, \mathbf{Z}\}$ of oracle matching under the following framework: we assume that the distributions $P_{Y,Z}$ and P_ζ in (16) are independent of the sample size n , while the regression function $\mu(\cdot, \cdot)$ does depend on n ; however, we suppress the dependence of n in our notation for simplicity.

Theorem 4. *Suppose that $\widehat{\text{ISS}}_n > 0$ and that $\|\mu(\mathbf{Y}, \mathbf{Z})\|_\infty \vee \|\mu(\mathbf{Y}, \mathbf{Z})\|_\infty^4 = o_P(\widehat{\text{ISS}}_n)$. Then, under the model (16), the conditional power of oracle matching with $\alpha \in (0, 1/2)$ satisfies*

$$\Phi\left(\frac{\widehat{\text{ISS}}_n}{\sqrt{2}\sigma} - \bar{\Phi}^{-1}(\alpha)\right) - o_P(1) \leq \mathbb{P}\{p \leq \alpha \mid \mathbf{Y}, \mathbf{Z}\} \leq \Phi\left(\frac{\widehat{\text{ISS}}_n}{\sigma} - \bar{\Phi}^{-1}(\alpha)\right) + o_P(1). \quad (18)$$

We observe that $\widehat{\text{ISS}}_n/\sigma$, which can be interpreted here as the *signal-to-noise ratio*, dictates both the upper and lower bound. We also notice that the upper and lower bounds on conditional power match up to a factor of 2—in fact, this factor is unavoidable without imposing further model assumptions (e.g., symmetry of $\mu(Y, Z)$ conditional on Z , as we will discuss in more detail in Appendix C.4.3).

Remark 1. *While the asymptotic lower and upper bounds require that $(\|\mu(\mathbf{Y}, \mathbf{Z})\|_\infty \vee \|\mu(\mathbf{Y}, \mathbf{Z})\|_\infty^4)/\widehat{\text{ISS}}_n = o_P(1)$, we remark that this is a natural assumption on the quality of matching. For example, if X is bounded, then $\|\mu(\mathbf{Y}, \mathbf{Z})\|_\infty \vee \|\mu(\mathbf{Y}, \mathbf{Z})\|_\infty^4 = O(1)$ —and so for this assumption to hold, it is sufficient to require $\widehat{\text{ISS}}_n \rightarrow \infty$, i.e., the amount of signal in the data increases with n .*

Power guarantees for plug-in matching with an estimate of μ . Now, we shift our attention to the more practical setting, where we do not have the oracle knowledge of μ . A natural solution is data splitting—i.e., we learn an estimate $\hat{\mu}$ on one random split of the data, and then implement PairSwap-ICI test with the plug-in matching \hat{M} . While the aforementioned data splitting approach is more accurate and practical, for the sake of simplicity, in order to state the following result, we assume that we have an independent data where we can learn $\hat{\mu}$. Finally, under suitable consistency assumptions on $\hat{\mu}$, we can recover the power guarantees in Theorem 4.

Theorem 5. *Consider the setting and assumptions of Theorem 4. Suppose we use plug-in matching \hat{M} with an estimate $\hat{\mu}$ constructed based on independent data, where $\hat{\mu}$ satisfies*

$$\|\hat{\mu}(\mathbf{Y}, \mathbf{Z}) - \mu(\mathbf{Y}, \mathbf{Z})\|_2 = o_P(\widehat{\text{ISS}}_n), \quad \|\hat{\mu}(\mathbf{Y}, \mathbf{Z}) - \mu(\mathbf{Y}, \mathbf{Z})\|_\infty^4 = o_P(\widehat{\text{ISS}}_n).$$

Then the conditional power of the PairSwap-ICI satisfies (18) for $\alpha \in (0, 1/2)$.

4.2 ISS characterizes hardness of testing the null H_0^{ICI}

While $\widehat{\text{ISS}}_n$ governs the asymptotic upper and lower bounds on power in 4 and 5, the relationship between ISS_n and power extends beyond the PairSwap-ICI test. Here, we establish this connection formally. In particular, we will see that the quantity ISS_n determines which alternative models are distinguishable from the class of null models, via *any* valid testing procedure. Towards this goal, we start with a simple total-variation calculation to give an upper bound on the power function of any valid test.

Proposition 6. *Fix $\alpha \in (0, 1)$. Fix any test ϕ that controls false positives at level α , i.e.,⁴*

$$\phi : (\mathcal{X} \times \mathbb{R} \times \mathbb{R})^n \rightarrow [0, 1] \text{ such that } \sup_{P \in H_0^{\text{ICI}}} \mathbb{E}_P[\phi(\mathbf{X}, \mathbf{Y}, \mathbf{Z})] \leq \alpha.$$

Then for any distribution $P_{X,Y,Z}$,

$$\mathbb{E}_{P_{X,Y,Z}}[\phi(\mathbf{X}, \mathbf{Y}, \mathbf{Z})] \leq \alpha + \inf_{Q_{X,Y,Z} \in H_0^{\text{ICI}}} \text{d}_{\text{TV}}(P_{X,Y,Z}^n, Q_{X,Y,Z}^n).$$

⁴Traditionally, we think of a hypothesis test ϕ as a map from data to a decision, i.e., $\phi(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) \in \{0, 1\}$. Why, then, do we define tests ϕ as mapping to the space $[0, 1]$? This is because, in some settings, we may want to consider randomized tests—for instance, in the notation above where ϕ maps to $[0, 1]$, an outcome $\phi(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) = 0.75$ represents that, given the data, our randomized test rejects the null with probability 0.75. Of course, a nonrandomized test is simply a special case, obtained by restricting the output of ϕ to lie in $\{0, 1\}$.

This last total variation term is the distance of $P_{X,Y,Z}$ from the null class H_0^{ICI} , meaning the closer P is to null models, the harder it will be to get non-trivial power against P . While in general it is hard to derive exact expressions for the total variation term, under some additional model assumptions on P , we can come up with interpretable upper bounds for the same. Two such examples are listed below.

- **Gaussian setting:** In the first example, we consider a special class of Gaussian alternatives, which is given by (16) with $P_\zeta = \mathcal{N}(0, \sigma^2)$ for some $\sigma > 0$. In this special case, ISS_n gives a meaningful upper bound on the offset total-variation term, up to a constant.

Corollary 7. *Suppose, $P_\zeta = \mathcal{N}(0, \sigma^2)$ and $(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$ satisfy (16). Then*

$$\inf_{Q_{X,Y,Z} \in H_0^{\text{ICI}}} \text{d}_{\text{TV}}(P_{X,Y,Z}^n, Q_{X,Y,Z}^n) \leq \frac{\text{ISS}_n}{2\sigma},$$

and consequently, for any test ϕ that controls Type I error at level α ,

$$\mathbb{E}_P[\phi(\mathbf{X}, \mathbf{Y}, \mathbf{Z})] \leq \alpha + \frac{\text{ISS}_n}{2\sigma}.$$

- **Binary setting:** Now, suppose $\mathbf{X}, \mathbf{Y}, \mathbf{Z}$ are generated from the model $P_{X,Y,Z}$ given by

$$X_i \sim \text{Ber}(\mu(Y_i, Z_i)), \quad (Y_1, Z_1), \dots, (Y_n, Z_n) \stackrel{\text{iid}}{\sim} P_{Y,Z}. \quad (19)$$

Here as well, ISS_n leads to a very interpretable upper bound on the power of our test, as long as $\mu(Y, Z)$ is almost surely away from the extremities, i.e., 0 or 1.

Corollary 8. *Suppose $(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$ satisfy (19) where $\mu(Y, Z) \in (\epsilon, 1 - \epsilon)$ almost surely for some $\epsilon > 0$. Then*

$$\inf_{Q_{X,Y,Z} \in H_0^{\text{ICI}}} \text{d}_{\text{TV}}(P_{X,Y,Z}^n, Q_{X,Y,Z}^n) \leq \left(\frac{1}{\epsilon(1 - \epsilon)} \right)^{1/2} \text{ISS}_n,$$

and consequently, for any test ϕ that controls Type I error at level α ,

$$\mathbb{E}_P[\phi(\mathbf{X}, \mathbf{Y}, \mathbf{Z})] \leq \alpha + \left(\frac{1}{\epsilon(1 - \epsilon)} \right)^{1/2} \text{ISS}_n.$$

To summarize, an appropriate *signal-to-noise ratio* determines the extent to which a test can outperform the trivial test that rejects the null H_0^{ICI} with probability α without using data. Consequently, in both examples, no valid test can achieve non-trivial power when ISS_n is negligible relative to the noise in $P_{X|Y,Z}$.

Note that this characterization relies on the oracle quantity ISS_n , while the conditional power of the PairSwap-ICI test is determined by the empirical version, $\widehat{\text{ISS}}_n$. We would expect that, under mild conditions, $\text{ISS}_n \approx \widehat{\text{ISS}}_n$ (i.e., a concentration property should hold, as long as n is large)—we can therefore interpret ISS_n as characterizing the power, both in terms of upper and lower bounds, at least for certain cases.

4.3 Near-optimal power guarantees without knowledge of μ

The asymptotic upper and lower bounds on conditional power from Theorems 4 and 5 require either oracle knowledge of μ or an independent estimate of μ whose Euclidean norm error is small by comparison with $\widehat{\text{ISS}}_n$. If this side information is not available, then a natural question arises: can near-optimal power guarantees be achieved with some of the simpler and more practical matching schemes from Section 3.3?

The aim of this subsection is to show that this is indeed possible, under some additional assumptions on the model class (16). Specifically, consider the class of partially linear Gaussian models given by

$$(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) \text{ satisfy (16) with } P_\zeta = \mathcal{N}(0, \sigma^2), \mu(\mathbf{Y}, \mathbf{Z}) = \mu_0(\mathbf{Z}) + \beta_n \mathbf{Y} \quad (20)$$

for some $\mu_0 \in \mathcal{C}_{\text{ISO}}$ (applied componentwise) and some $\beta_n \in [0, \infty)$. We also assume that Y is a bounded and mean-zero random variable, with $|Y| \leq 1$. Note that μ_0 is a fixed function, and so the dependence on n of the distribution of (X, Y, Z) is solely through β_n . The following lemma relates ISS_n to a more interpretable quantity, namely the expected conditional variance of Y given Z .

Lemma 9. *Under the model class (20),*

$$\sqrt{n}\beta_n \cdot (\mathbb{E}[\text{Var}(Y | Z)])^{1/2} (1 + o_P(1)) \leq \text{ISS}_n \leq \sqrt{n}\beta_n.$$

We can draw several conclusions from this result. By Lemma 9 and Corollary 7, no valid test can achieve asymptotically non-trivial power against an alternative in the model class (20) when $\beta_n = o(1/\sqrt{n})$. From this point on, then, we will consider the regime $\beta_n \gtrsim 1/\sqrt{n}$. Another consequence of Lemma 9 is that the asymptotic lower bound on the conditional power of oracle matching in (18) further simplifies to give

$$\mathbb{P}\{p \leq \alpha | \mathbf{Y}, \mathbf{Z}\} \geq \Phi\left(\sqrt{n}\beta_n \left\{\frac{\mathbb{E}[\text{Var}(Y | Z)]}{2\sigma^2}\right\}^{1/2} - \bar{\Phi}^{-1}(\alpha)\right) - o_P(1). \quad (21)$$

Next, we aim to show that practical matching schemes, such as neighbour matching and cross-bin matching, can achieve conditional power at least as good as the lower bound above.

Theorem 10. *Let $\beta_n \gtrsim 1/\sqrt{n}$, and assume that $\mu_0(Z)$ is a sub-Gaussian random variable. Then the conditional power of neighbour matching (Algorithm 1), implemented with kernel $\psi(x, x') = x - x'$ and weights $w_\ell = \max\{Y_{i_\ell} - Y_{j_\ell}, 0\}$, satisfies*

$$\left|\mathbb{P}\{p \leq \alpha | \mathbf{Y}, \mathbf{Z}\} - \Phi\left(\sqrt{n}\beta_n \left\{\frac{\mathbb{E}[\text{Var}(Y | Z)]}{4\sigma^2}\right\}^{1/2} - \bar{\Phi}^{-1}(\alpha)\right)\right| = o_P(1).$$

Notably, the conditional power of neighbour matching matches the lower bound from (21) up to a factor of 2. This factor arises due to the inherent inefficiency of neighbour matching: as discussed in Section 3.3, we have seen that neighbour matching discards roughly half the data, which explains the extra factor of 2.

To avoid this loss of a factor of 2, we now consider the cross-bin matching strategy, which will typically assign nearly all data points to a matched pair. Consequently, the power of cross-bin matching can then meet the lower bound in (21).

Theorem 11. *Consider the setting as in Theorem 10. Then, the conditional power of cross-bin matching (Algorithm 2) with $K \propto \sqrt{n}$ many bins, implemented with kernel $\psi(x, x') = x - x'$, weights $w_\ell = \max\{Y_{i_\ell} - Y_{j_\ell}, 0\}$, satisfies*

$$\mathbb{P}\{p \leq \alpha \mid \mathbf{Y}, \mathbf{Z}\} \geq \Phi \left(\sqrt{n}\beta_n \left\{ \frac{\mathbb{E}[\text{Var}(Y \mid Z)]}{2\sigma^2} \right\}^{1/2} - \bar{\Phi}^{-1}(\alpha) \right) - o_P(1)$$

under suitable smoothness assumptions (stated formally in Theorem C.14).

In particular, this result matches the lower bound in (21), without an additional factor of 2 as for neighbour matching. (While the above result provides only an asymptotic lower bound on the conditional power of cross-bin matching, under additional assumptions we can characterize the power more exactly—see Appendix C.4.)

5 Simulations

In this section, we evaluate the performance of our method on simulated data, and compare the matching strategies from Section 3.3. For simplicity, we focus on the univariate case $Z = \mathbb{R}$. We will test two versions of the PairSwap-ICI method:

- Neighbour matching (Algorithm 1), with the linear kernel $\psi(x, x') = x - x'$ and weights $w_\ell = \max\{Y_{i_\ell} - Y_{j_\ell}, 0\}$ as discussed in Section 3.3;
- Cross-bin matching (Algorithm 2) with $K = 2 \lfloor n^{1/2} \rfloor$ bins, again with the linear kernel $\psi(x, x') = x - x'$ and weights $w_\ell = \max\{Y_{i_\ell} - Y_{j_\ell}, 0\}$. This particular choice of K is motivated by Theorem 11, which suggests choosing $K \propto \sqrt{n}$ bins in order to achieve competitive power with neighbour matching.

5.1 Conservativeness under the null H_0^{ICI}

Theorem 1 establishes valid, finite-sample Type I error control for our method. The purpose of this section is to evaluate how conservative the Type I error is under various null distributions. Because our inference relies on the fact that matched pairs (X_{i_ℓ}, X_{j_ℓ}) are stochastically ordered under the null, intuitively the conservativeness of our test depends on the strength of monotonicity in the conditional distribution.

To see how the dependence between X and Z affects the rejection probability, we sample X from an additive noise model

$$X \mid Y, Z \sim \mathcal{N}(\mu(\gamma Z), 1),$$

where Y, Z are independent standard normal random variables. As long as μ is nondecreasing and $\gamma \geq 0$, this joint distribution belongs to the null H_0^{ICI} . The scalar γ controls the strength of the monotonicity of $X \mid Z$. In particular, as $\gamma \downarrow 0$ we expect the Type I error $\mathbb{P}\{p \leq \alpha\}$ to approach α .

In our simulations, we consider two functions μ , the identity $\mu(z) = z$ and the Gaussian CDF $\mu(z) = \Phi(z)$. Figure 2 shows the Type I error as a function of γ for two levels of α . We observe similar results for each α , where the test typically becomes more conservative as γ increases, as expected. Under the null, our test is more conservative for cross-bin matching than for neighbour matching, since the Z values are further apart in cross-bin matching.

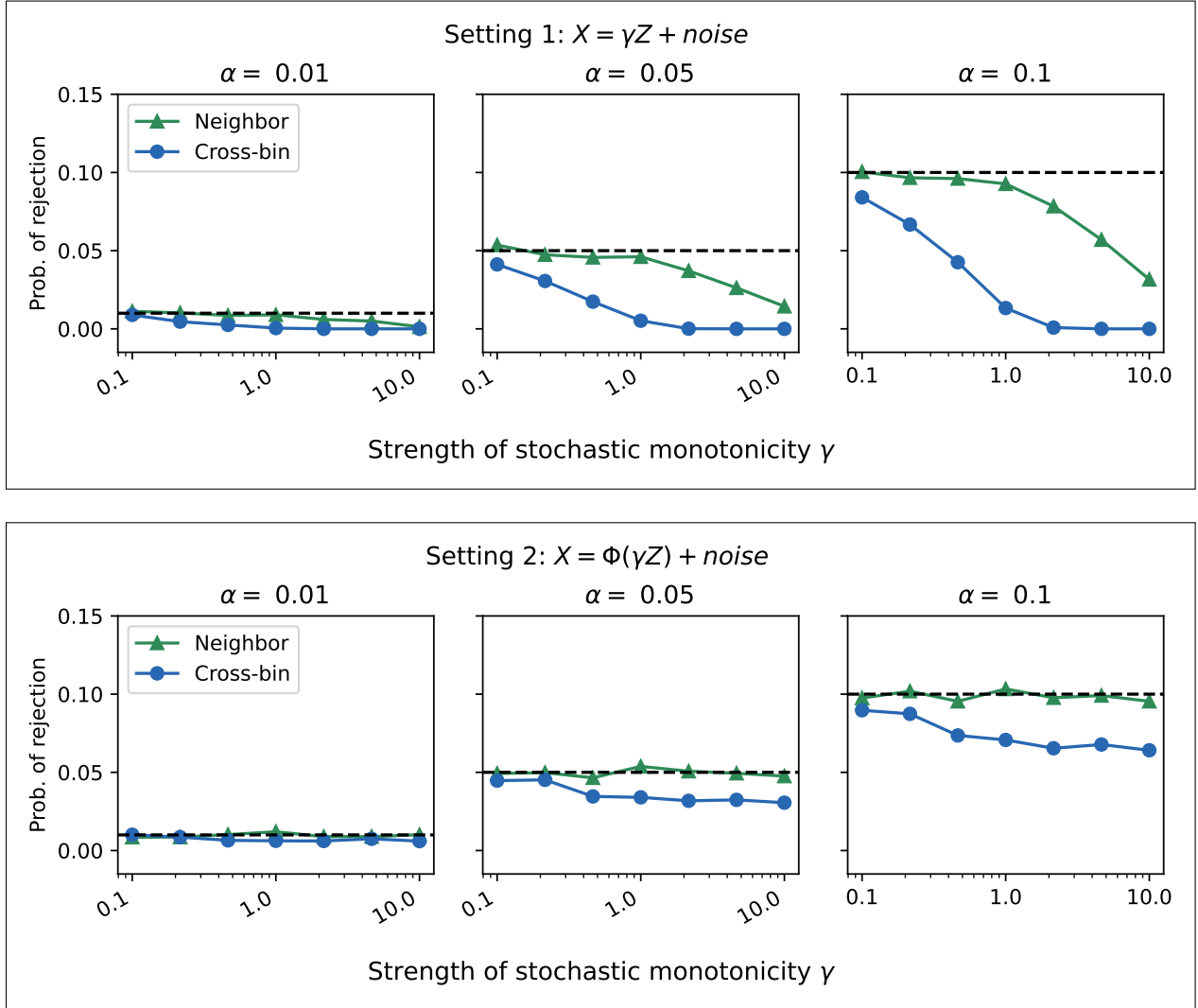


Figure 2: Simulation results illustrating Type I error control under the null H_0^{ICI} for two forms of the conditional mean $\mathbb{E}[X | Z]$. Each subplot shows the rejection probability of PairSwap-ICI test on data set of size 1000, averaged over 10^4 simulation trials, as a function of the strength of stochastic monotonicity γ .

5.2 Power under alternatives

In Section 4.3 we showed theoretically that our heuristic methods—neighbour matching and cross-bin matching—achieve high power (in fact, power that tends to 1) in the partial linear model (20) provided the signal β_n exceeds the detection threshold $n^{-1/2}$. In this section, we

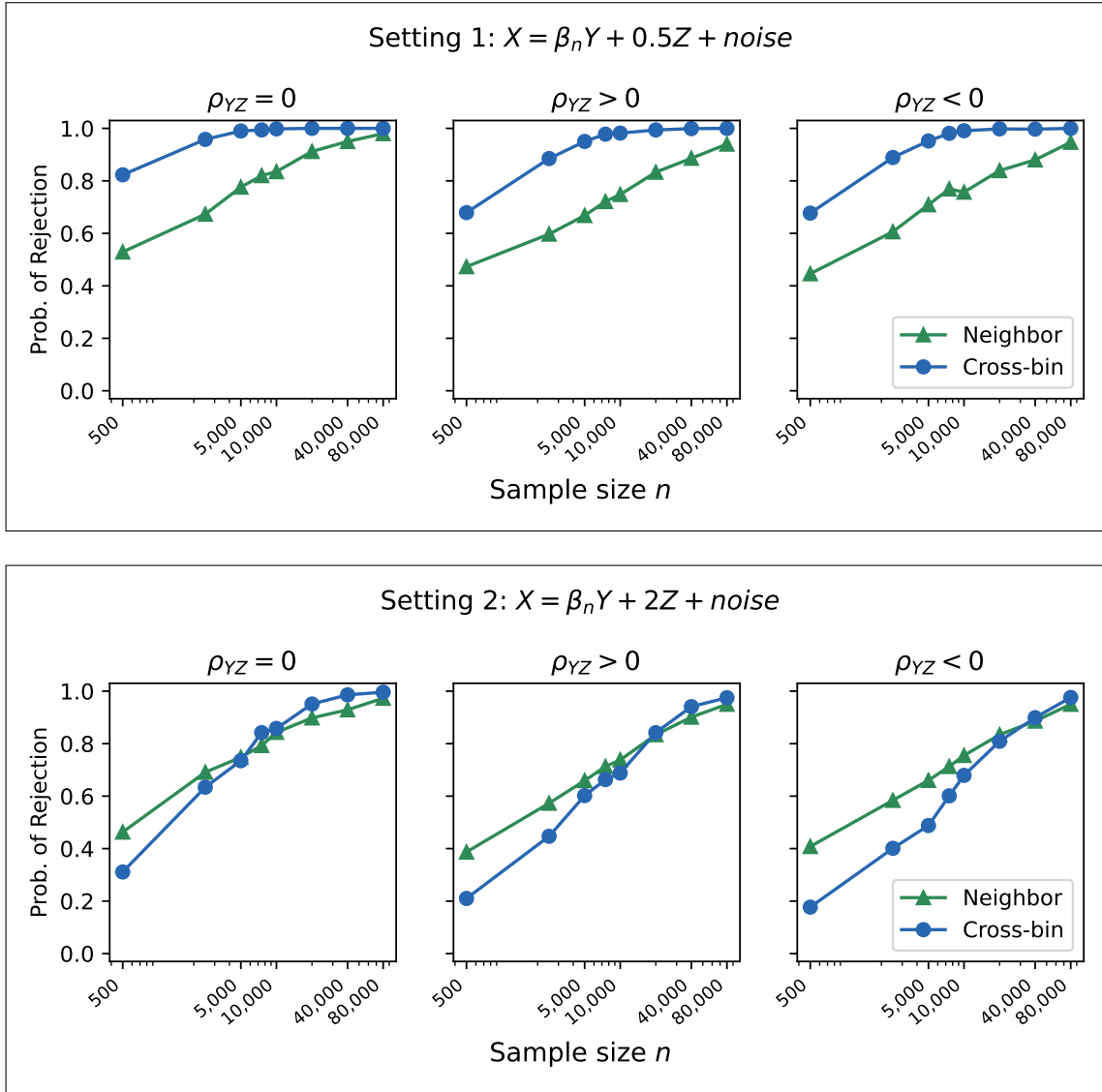


Figure 3: Simulation results demonstrating power for two alternatives at level $\alpha = 0.1$. Each subplot shows the rejection probability, averaged over 10^3 simulation trials, as a function of the sample size n . Columns correspond to different relationships between Y and Z . In each setting, X follows a Gaussian linear model with mean $\beta_n Y + \gamma Z$, where $\beta_n = n^{-1/3}$ and $\gamma = 0.5$ (above) or $\gamma = 2$ (below).

now examine this setting empirically. We sample data from the Gaussian linear model

$$X \mid Y, Z \sim \mathcal{N}(\beta_n Y + \gamma Z, 1),$$

with $\beta_n = n^{-1/3}$. The pair (Y, Z) is drawn from a bivariate Gaussian

$$\mathcal{N}\left(0, \begin{bmatrix} 1 & \rho_{YZ} \\ \rho_{YZ} & 1 \end{bmatrix}\right).$$

Figure 3 shows the power as a function of the sample size n for various choices of γ and ρ_{YZ} . In Setting 1, we set $\gamma = 0.5$, and cross-bin matching uniformly dominates neighbour

matching because it allows us to make many more matches of similar quality. On the other hand, in Setting 2 we set $\gamma = 2$, so the strong dependence of X on Z means the quality of a match (i_ℓ, j_ℓ) degrades much more quickly as the gap $Z_{j_\ell} - Z_{i_\ell}$ increases—that is, for cross-bin matching, where Z_{i_ℓ} and Z_{j_ℓ} may be farther apart than for neighbour matching, this gap may lead to conservativeness that results in a loss of power. However, with sufficiently large sample size, cross-bin matching performs at least as well as the neighbour matching. This is because the bin width decreases as n increases, so the quality of the cross-bin matches rivals that of the neighbour matches (with many more matches). The dependence ρ_{YZ} between Y and Z does not have a major impact on the power of these two methods.

6 Experiment on real data: risk factors for diabetes

In this section, we evaluate the performance of our proposed testing procedure on a real data set in three different experimental settings. The PairSwap-ICI method is implemented with the heuristic neighbour matching or with cross-bin matching (with $K = 50$ bins), and our statistic T takes the form in (11) with linear kernel $\psi(x, x') = x - x'$ and weights $w_\ell = \max\{Y_{i_\ell} - Y_{j_\ell}, 0\}$.

We use a data set⁵ on the incidence of diabetes among the Pima population near Phoenix, Arizona, originally collected by the US National Institute of Diabetes and Digestive and Kidney Diseases. The data set contains 768 observations, and it includes information on whether each of the patient has been diagnosed with diabetes according to World Health Organization standards. Additional variables provide data on the number of pregnancies, plasma glucose concentration, diastolic blood pressure, triceps skinfold thickness, 2-hour serum insulin levels, body mass index (BMI), diabetes pedigree function and age.

It is well-known that the likelihood of developing diabetes increases with age (e.g., the CDC⁶ lists advanced age as one of the risk factors for type 1 and type 2 diabetes). Therefore, if we choose X to represent the incidence of diabetes and Z as the age of the patient, then we would expect X to exhibit stochastic monotonicity with respect to Z (i.e., we expect that Assumption 1 holds, at least approximately). This is supported by the increasing trend we observe in Figure 4. Most of the other variables, such as `BloodPressure`, `BMI`, `Glucose`, and `Pregnancies`, are also considered potential risk factors for diabetes, as visualized in Figure 5—but does this association remain after we control for age? In this experiment, we aim to determine whether these variables remain significant risk factors for diabetes, even after controlling for age.

Experiment 1: marginal independence testing: In our first experiment, we consider six variables: `Pregnancies`, `Glucose`, `BloodPressure`, `SkinThickness`, `Insulin`, and `BMI`, and aim to assess whether each of them is an individual risk factor for diabetes incidence. Specifically, we test the hypothesis $H_0 : X \perp\!\!\!\perp Y$, where Y represents one of the six variables listed above, while X is `Diabetes` (and since we are testing marginal rather than conditional independence, we do not attempt to control for Z , i.e., `Age`). For this purpose, we will be

⁵The data for this experiment were obtained from <https://www.kaggle.com/datasets/uciml/pima-indians-diabetes-database>. Additional data descriptions can be found in Smith et al. (1988).

⁶For more details, refer to the [list of diabetes risk factors](#) from U.S. Centers for Disease Control and Prevention.

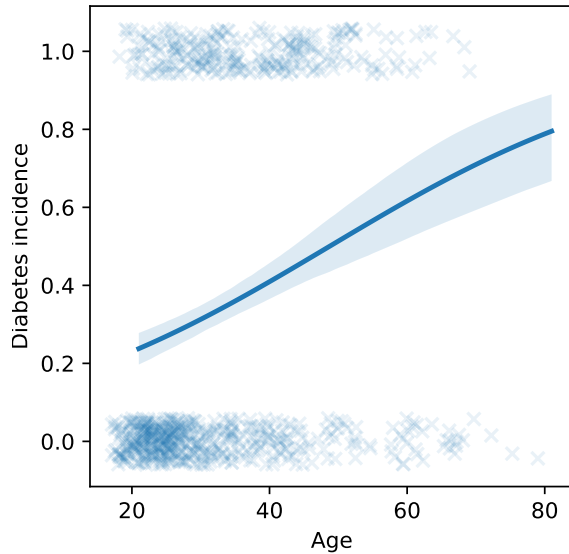


Figure 4: A scatter-plot (jittered for better visibility) of `Age` and `Diabetes Incidence` along with the fitted logistic regression model to demonstrate the stochastic monotonicity between them.

using the permutation test for independence with $T(\mathbf{X}, \mathbf{Y}) = \mathbf{X}^T \mathbf{Y}$, as outlined in Section 1.2.

Experiment 2: conditional independence testing, after controlling for age: Next, for the same set of six choices for Y , we test the hypothesis $H_0^{CI} : X \perp\!\!\!\perp Y \mid Z$, where Z denotes `Age` (and X is `Diabetes` as before). This allows us to identify risk factors for diabetes after controlling for age. As noted earlier, we expect the distribution of $X \mid Z = z$ to be stochastically monotone in z , which supports the application of the `PairSwap-ICI` testing procedure developed in this paper for this purpose.

Experiment 3: conditional independence testing, with synthetic control \tilde{X} : Finally, we consider a semi-synthetic experiment where X is replaced by synthetic observations \tilde{X} , generated from an estimated model for $P_{X|Z}$ that satisfies stochastic monotonicity. We then test the hypothesis $\tilde{\mathcal{H}}_0 : \tilde{X} \perp\!\!\!\perp Y \mid Z$ for the same choices of Y from Experiment 1. Since \tilde{X} is generated solely based on Z , the null hypothesis of conditional independence holds trivially in this synthetic setting. The validity of our procedure should therefore ensure that the p -values generated by `PairSwap-ICI` are (super)uniformly distributed.

Now we give details on how the synthetic feature \tilde{X} is generated. Since X is binary, it suffices to fit an isotonic regression to estimate the conditional mean $\mathbb{E}[X \mid Z]$ and then sample \tilde{X} from the Bernoulli distribution with this fitted conditional mean. Following the theory of [Henzi, Ziegel and Gneiting \(2021, Theorem 1\)](#), this is the best approximation for $P_{X|Z}$ under *continuous ranked probability score (CRPS)*, while respecting the monotonicity constraint.

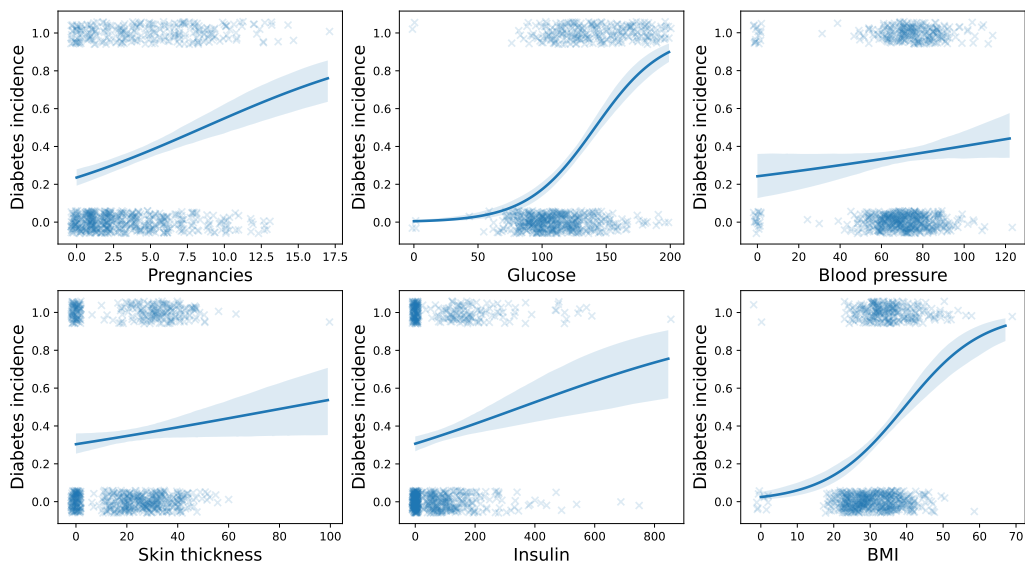


Figure 5: Scatter plots of X (jittered for better visibility) and other feature variables along with the fitted logistic regression models to demonstrate the dependence among these variables and Diabetes Incidence.

		Pregnancies	Glucose	Blood pressure	Skin thickness	Insulin	BMI
Permutation test (testing marginal indep.)		0.001 (0.00)	0.001 (0.00)	0.155(0.003)	0.117(0.002)	0.023 (0.001)	0.001 (0.00)
PairSwap-ICI	neighbour matching	0.433(0.005)	0.002 (0.00)	0.496(0.005)	0.243(0.004)	0.159(0.003)	0.033 (0.001)
	cross-bin matching	0.424(0.004)	0.001 (0.00)	0.499(0.004)	0.123(0.003)	0.224(0.004)	0.001 (0.00)
PairSwap-ICI (with synthetic control)	neighbour matching	0.511(0.005)	0.499(0.005)	0.501(0.005)	0.506(0.005)	0.497(0.005)	0.499(0.005)
	cross-bin matching	0.544(0.005)	0.550(0.005)	0.555(0.005)	0.558(0.005)	0.545(0.005)	0.555(0.005)

Table 1: p -values, averaged over 3000 random sub-samples along with the estimated standard errors (within brackets) for the different tests from different experiments, as outlined in Section 6. The p -values significant at the 0.05 level are marked in bold.

For each experiment, we generate 3,000 random sub-samples of the data, each consisting of half the size of the full data set. We then compute p -values using the permutation test for marginal independence and the PairSwap-ICI test for conditional independence. For experiments involving synthetic control, which require estimating $P_{X|Z}$, the sub-sampled data is further divided into training and test sets, with $\hat{P}_{X|Z}$ being computed in training set. For all the experiments, p -values are computed in the test set. Finally, we report the average p -values from the 3,000 sub-samples, along with the corresponding choice of the Y variable in Table 1.

Results: Under the marginal independence test, four of the six variables are identified as having significant association with `Diabetes`—but, once we test conditional independence with the `PairSwap-ICI` test, only two of these associations are identified as significant. Specifically, `Glucose` and `BMI` both are identified as potential risk factors at the 0.05 level of significance by the marginal independence test, and also by the `PairSwap-ICI` test, even after controlling for `Age`. On the other hand, the variables `Pregnancies` and `Insulin` are significant only under the marginal test; this suggests that, after controlling for `Age`, the data does not provide sufficient evidence to support them as risk factors for `Diabetes`.

Finally, we also note that all the averaged p -value from Experiment 3 with synthetic control \tilde{X} is concentrated around 0.5, for each of the choices of Y . Since (\tilde{X}, Y, Z) satisfy H_0^{ICI} the p -values from Experiment 3 should be roughly uniform (or, if the test is conservative, superuniform), and thus this behaviour is expected as per the result we have established in Theorem 1.

7 Discussion

In this paper, we have developed a nonparametric test of conditional independence assuming only stochastic monotonicity of the conditional distribution $P_{X|Z}$. This nonparametric constraint is natural in many applications, and allows us to circumvent the impossibility of assumption-free conditional independence testing (Shah and Peters, 2020). We have introduced a variety of approaches to constructing a valid test statistic. Our test controls the Type I error in finite samples and has power against an array of alternatives. We close our discussion with some interesting connections to the literature, and potential avenues for future work.

- *Optimal power in general settings.* Theorems 4 and 5 bound the asymptotic power of our test above and below, where the upper and lower bounds differ by the appearance of the constant 2 in the lower bound—as we will see in Appendix C.4.3, this difference can be removed if we assume that the conditional distribution of $Y | Z$ is symmetric, but it remains an open question whether other tests (or, perhaps, the `PairSwap-ICI` but with a different kernel) may be able to avoid this assumption.
- *Avoiding data splitting.* The oracle matching test derived in Section 3 requires modeling the conditional mean and conditional variance of the kernel $\psi(X_i, X_j)$ as a function of Y_i, Y_j, Z_i, Z_j . We proposed to estimate these moments on a hold-out data set. Can we instead perform cross-fitting to improve power and retain finite-sample error control?
- *Alternative methods.* A notable benefit of our stochastic monotonicity assumption is that one can consistently estimate the conditional distribution $P_{X|Z}$ using isotonic distributional regression (Mösching and Dümbgen, 2020; Henzi, Ziegel and Gneiting, 2021). Hence, an alternative approach to testing the restricted null H_0^{ICI} is to first estimate this conditional distribution on one split of the data, and then run a conditional independence test which assumes knowledge of $P_{X|Z}$ (Berrett et al., 2020; Candès et al., 2018). Since we are plugging in the estimated conditional distribution, such tests will only be valid asymptotically. Is there any way to modify such tests to be valid in finite samples?

- *Connection with knockoffs and conditional randomization tests.* Creating synthetic copies of \mathbf{X} via pairwise swaps has resemblance to other conditional independence testing procedures, such as knockoffs and the conditional randomization test (Candès et al., 2018), and the conditional permutation test (Berrett et al., 2020). One difference is that in our method, due to the stochastic ordering assumption, creating the swapped copies of \mathbf{X} is potentially more conservative (i.e., the resulting p -value may be superuniform), since we are not working under the “sharp null”.
- *Alternative shape constraints.* We view stochastic monotonicity as one form of positive dependence for the joint distribution (X, Z) . Are there natural approaches to test conditional independence under other models of dependence, such as likelihood-ratio ordering or total positivity, or under other shape constraints, such as unimodality (Karlin, 1968; Shaked and Shanthikumar, 2007; Mösching and Dümbgen, 2024)?

References

- Albert, M., Laurent, B., Marrel, A. and Meynaoui, A. (2022) Adaptive test of independence based on HSIC measures. *The Annals of Statistics*, **50**, 858–879.
- Azadkia, M. and Chatterjee, S. (2021) A simple measure of conditional dependence. *The Annals of Statistics*, **49**, 3070–3102.
- Barber, R. F., Candès, E. J. and Samworth, R. J. (2020) Robust inference with knockoffs. *The Annals of Statistics*, **48**, 1409–1431.
- Berrett, T. B., Kontoyiannis, I. and Samworth, R. J. (2021) Optimal rates for independence testing via U-statistic permutation tests. *The Annals of Statistics*, **49**, 2457–2490.
- Berrett, T. B. and Samworth, R. J. (2019) Nonparametric independence testing via mutual information. *Biometrika*, **106**, 547–566.
- Berrett, T. B. and Samworth, R. J. (2021) USP: an independence test that improves on Pearson’s chi-squared and the G-test. *Proceedings of the Royal Society A*, **477**, 20210549.
- Berrett, T. B., Wang, Y., Barber, R. F. and Samworth, R. J. (2020) The conditional permutation test for independence while controlling for confounders. *J. R. Stat. Soc. Ser. B. Stat. Methodol.*, **82**, 175–197.
- Candès, E., Fan, Y., Janson, L. and Lv, J. (2018) Panning for gold: ‘model-X’ knockoffs for high dimensional controlled variable selection. *J. R. Stat. Soc. Ser. B. Stat. Methodol.*, **80**, 551–577.
- Davison, A. C. and Hinkley, D. V. (1997) *Bootstrap Methods and their Application*, vol. 1 of *Cambridge Series in Statistical and Probabilistic Mathematics*. Cambridge University Press, Cambridge.
- Duan, R. and Pettie, S. (2014) Linear-time approximation for maximum weight matching. *Journal of the ACM (JACM)*, **61**, 1–23.

- Edmonds, J. (1965) Paths, trees, and flowers. *Canadian Journal of mathematics*, **17**, 449–467.
- Gabow, H. N. (1985) A scaling algorithm for weighted matching on general graphs. In *26th Annual Symposium on Foundations of Computer Science (sfcs 1985)*, 90–100, IEEE.
- Hagberg, A., Swart, P. and Schult, D. (2008) Exploring network structure, dynamics, and function using NetworkX. Tech. rep., Los Alamos National Lab.(LANL), Los Alamos, NM (United States).
- Henzi, A., Ziegel, J. F. and Gneiting, T. (2021) Isotonic distributional regression. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, **83**, 963–993.
- Kalisch, M. and Bühlmann, P. (2007) Estimating high-dimensional directed acyclic graphs with the PC-algorithm. *Journal of Machine Learning Research*, **8**.
- Karlin, S. (1968) *Total positivity. Vol. I*. Stanford University Press, Stanford, CA.
- Kim, I., Balakrishnan, S. and Wasserman, L. (2022) Minimax optimality of permutation tests. *The Annals of Statistics*, **50**, 225–251.
- Kim, I., Neykov, M., Balakrishnan, S. and Wasserman, L. (2022) Local permutation tests for conditional independence. *Ann. Statist.*, **50**, 3388–3414.
- Lei, J. (2020) Convergence and concentration of empirical measures under Wasserstein distance in unbounded functional spaces. *Bernoulli*, **26**, 767 – 798.
- Lundborg, A. R., Kim, I., Shah, R. D. and Samworth, R. J. (2024+) The Projected Covariance Measure for assumption-lean variable significance testing. *The Annals of Statistics*, to appear. *arXiv preprint arXiv:2211.02039*.
- Lundborg, A. R., Shah, R. D. and Peters, J. (2022) Conditional independence testing in Hilbert spaces with applications to functional data analysis. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, **84**, 1821–1850.
- Mösching, A. and Dümbgen, L. (2020) Monotone least squares and isotonic quantiles. *Electron. J. Stat.*, **14**, 24–49.
- Mösching, A. and Dümbgen, L. (2024) Estimation of a likelihood ratio ordered family of distributions. *Stat. Comput.*, **34**, Paper No. 58, 16.
- Neykov, M., Balakrishnan, S. and Wasserman, L. (2021) Minimax optimal conditional independence testing. *Ann. Statist.*, **49**, 2151–2177.
- Niu, Z., Chakraborty, A., Dukes, O. and Katsevich, E. (2024) Reconciling model-X and doubly robust approaches to conditional independence testing. *The Annals of Statistics*, **52**, 895–921.
- O’Mahony, C., Jichi, F., Pavlou, M., Monserrat, L., Anastasakis, A., Rapezzi, C., Biagini, E., Gimeno, J. R., Limongelli, G., McKenna, W. J. et al. (2014) A novel clinical risk prediction model for sudden cardiac death in hypertrophic cardiomyopathy (HCM risk-SCD). *European Heart Journal*, **35**, 2010–2020.

- Pfister, N., Bühlmann, P., Schölkopf, B. and Peters, J. (2018) Kernel-based tests for joint independence. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, **80**, 5–31.
- Phipson, B. and Smyth, G. K. (2010) Permutation p -values should never be zero: calculating exact p -values when permutations are randomly drawn. *Stat. Appl. Genet. Mol. Biol.*, **9**, Art. 39, 14.
- Shah, R. D. and Peters, J. (2020) The hardness of conditional independence testing and the generalised covariance measure. *Ann. Statist.*, **48**, 1514–1538.
- Shaked, M. and Shanthikumar, J. G. (2007) *Stochastic orders*. Springer Series in Statistics, Springer, New York.
- Shevtsova, I. G. (2010) An improvement of convergence rate estimates in the Lyapunov theorem. *Doklady Mathematics*, **82**, 862–864.
- Smith, J. W., Everhart, J. E., Dickson, W., Knowler, W. C. and Johannes, R. S. (1988) Using the ADAP learning algorithm to forecast the onset of diabetes mellitus. In *Proceedings of the annual symposium on computer application in medical care*, 261, American Medical Informatics Association.
- Villani, C. (2009) *Optimal transport: old and new*, vol. 338. Springer.
- Yan, Z., Cai, M., Han, X., Chen, Q. and Lu, H. (2023) The interaction between age and risk factors for diabetes and prediabetes: a community-based cross-sectional study. *Diabetes, Metabolic Syndrome and Obesity*, 85–93.

A Proof of Theorem 2

The proof of this result follows the same structure as the proof of Theorem 1.

Step 1: some deterministic properties of the p -value. Define a function $\hat{p}_M : \mathbb{R}^n \times (\{\pm 1\}^L)^M \rightarrow [0, 1]$ as

$$\hat{p}_M(\mathbf{x}; \mathbf{s}^{(1)}, \dots, \mathbf{s}^{(M)}) = \frac{1 + \sum_{m=1}^M \mathbb{1} \left\{ T(\mathbf{x}^{\mathbf{s}^{(m)}}) \geq T(\mathbf{x}) \right\}}{1 + M}.$$

As in the proof of Theorem 1, this function is monotone nonincreasing in each x_{i_ℓ} , and monotone nondecreasing in each x_{j_ℓ} .

Step 2: compare to the sharp null. Define $\mathbf{X}_\#$ as in the proof of Theorem 1. Following identical arguments as in that proof, we can verify that, for any fixed $\mathbf{s}^{(1)}, \dots, \mathbf{s}^{(M)}$, it holds that

$$\hat{p}_M(\mathbf{X}_\#; \mathbf{s}^{(1)}, \dots, \mathbf{s}^{(M)}) \preceq_{\text{st}} \hat{p}_M(\mathbf{X}; \mathbf{s}^{(1)}, \dots, \mathbf{s}^{(M)})$$

conditional on \mathbf{Y}, \mathbf{Z} . Since $\hat{p}_M = \hat{p}_M(\mathbf{X}; \mathbf{s}^{(1)}, \dots, \mathbf{s}^{(M)})$ by construction, we therefore have

$$\mathbb{P} \left\{ \hat{p}_M \leq \alpha \mid \mathbf{Y}, \mathbf{Z}, \mathbf{s}^{(1)}, \dots, \mathbf{s}^{(M)} \right\} \leq \mathbb{P} \left\{ \hat{p}_M(\mathbf{X}_\#; \mathbf{s}^{(1)}, \dots, \mathbf{s}^{(M)}) \leq \alpha \mid \mathbf{Y}, \mathbf{Z}, \mathbf{s}^{(1)}, \dots, \mathbf{s}^{(M)} \right\}.$$

Marginalizing over the random draw of the swaps, $\mathbf{s}^{(1)}, \dots, \mathbf{s}^{(M)} \stackrel{\text{iid}}{\sim} \text{Unif}(\{\pm 1\}^L)$, we therefore have

$$\mathbb{P} \{ \hat{p}_M \leq \alpha \mid \mathbf{Y}, \mathbf{Z} \} \leq \mathbb{P} \{ \hat{p}_M(\mathbf{X}_\#; \mathbf{s}^{(1)}, \dots, \mathbf{s}^{(M)}) \leq \alpha \mid \mathbf{Y}, \mathbf{Z} \}.$$

Step 3: validity under the sharp null. We now need to verify the validity of the Monte Carlo p -value, under the sharp null. Unlike the first two steps, for this step the arguments are somewhat different than in the proof of Theorem 1.

First, let $\mathbf{s}^{(0)}$ be an additional draw from $\text{Unif}(\{\pm 1\}^L)$, sampled independently from all other random variables. Then it holds that

$$(\mathbf{s}^{(1)}, \dots, \mathbf{s}^{(M)}) \stackrel{\text{d}}{=} (\mathbf{s}^{(0)} \circ \mathbf{s}^{(1)}, \dots, \mathbf{s}^{(0)} \circ \mathbf{s}^{(M)}),$$

where \circ denotes the elementwise product, and so

$$\hat{p}_M(\mathbf{X}_\#; \mathbf{s}^{(1)}, \dots, \mathbf{s}^{(M)}) \stackrel{\text{d}}{=} \hat{p}_M(\mathbf{X}_\#; \mathbf{s}^{(0)} \circ \mathbf{s}^{(1)}, \dots, \mathbf{s}^{(0)} \circ \mathbf{s}^{(M)})$$

conditional on \mathbf{Y}, \mathbf{Z} . Moreover, by construction of the sharp null data $\mathbf{X}_\#$,

$$\mathbf{X}_\# \stackrel{\text{d}}{=} (\mathbf{X}_\#)^{\mathbf{s}^{(0)}}$$

holds conditional on $\mathbf{Y}, \mathbf{Z}, \mathbf{s}^{(0)}, \mathbf{s}^{(1)}, \dots, \mathbf{s}^{(M)}$, and therefore

$$\hat{p}_M(\mathbf{X}_\#; \mathbf{s}^{(0)} \circ \mathbf{s}^{(1)}, \dots, \mathbf{s}^{(0)} \circ \mathbf{s}^{(M)}) \stackrel{\text{d}}{=} \hat{p}_M((\mathbf{X}_\#)^{\mathbf{s}^{(0)}}; \mathbf{s}^{(0)} \circ \mathbf{s}^{(1)}, \dots, \mathbf{s}^{(0)} \circ \mathbf{s}^{(M)})$$

holds conditional on $\mathbf{Y}, \mathbf{Z}, \mathbf{s}^{(0)}, \mathbf{s}^{(1)}, \dots, \mathbf{s}^{(M)}$. Combining all these calculations so far, then, we have

$$\hat{p}_M(\mathbf{X}_\#; \mathbf{s}^{(1)}, \dots, \mathbf{s}^{(M)}) \stackrel{\text{d}}{=} \hat{p}_M((\mathbf{X}_\#)^{\mathbf{s}^{(0)}}; \mathbf{s}^{(0)} \circ \mathbf{s}^{(1)}, \dots, \mathbf{s}^{(0)} \circ \mathbf{s}^{(M)}), \quad (22)$$

conditional on \mathbf{Y}, \mathbf{Z} .

Next we calculate this last p -value: by definition,

$$\begin{aligned} \hat{p}_M((\mathbf{X}_\#)^{\mathbf{s}^{(0)}}; \mathbf{s}^{(0)} \circ \mathbf{s}^{(1)}, \dots, \mathbf{s}^{(0)} \circ \mathbf{s}^{(M)}) &= \frac{1 + \sum_{m=1}^M \mathbb{1} \left\{ T((\mathbf{X}_\#)^{\mathbf{s}^{(m)}}) \geq T((\mathbf{X}_\#)^{\mathbf{s}^{(0)}}) \right\}}{1 + M} \\ &= \frac{\sum_{m=0}^M \mathbb{1} \left\{ T((\mathbf{X}_\#)^{\mathbf{s}^{(m)}}) \geq T((\mathbf{X}_\#)^{\mathbf{s}^{(0)}}) \right\}}{1 + M}, \end{aligned}$$

where the first step holds since, for each $m = 1, \dots, M$,

$$\left[(\mathbf{X}_\#)^{\mathbf{s}^{(0)}} \right]^{\mathbf{s}^{(0)} \circ \mathbf{s}^{(m)}} = (\mathbf{X}_\#)^{\mathbf{s}^{(0)} \circ \mathbf{s}^{(0)} \circ \mathbf{s}^{(m)}} = (\mathbf{X}_\#)^{\mathbf{s}^{(m)}}$$

by definition of the swap operation. In other words, the p -value $\hat{p}_M((\mathbf{X}_\#)^{\mathbf{s}^{(0)}}; \mathbf{s}^{(0)} \circ \mathbf{s}^{(1)}, \dots, \mathbf{s}^{(0)} \circ \mathbf{s}^{(M)})$ is simply comparing the value of the statistic $T((\mathbf{X}_\#)^{\mathbf{s}^{(0)}})$ against the list of $M + 1$ values $T((\mathbf{X}_\#)^{\mathbf{s}^{(0)}}), \dots, T((\mathbf{X}_\#)^{\mathbf{s}^{(M)}})$. We therefore have

$$\mathbb{P} \left\{ \hat{p}_M((\mathbf{X}_\#)^{\mathbf{s}^{(0)}}; \mathbf{s}^{(0)} \circ \mathbf{s}^{(1)}, \dots, \mathbf{s}^{(0)} \circ \mathbf{s}^{(M)}) \leq \alpha \mid \mathbf{X}_\#, \mathbf{Y}, \mathbf{Z} \right\} \leq \alpha,$$

since, conditional on $\mathbf{X}_\#$, \mathbf{Y} , \mathbf{Z} , the sign vectors $\mathbf{s}^{(0)}, \dots, \mathbf{s}^{(M)}$ are i.i.d., and therefore the rank of the statistic $T((\mathbf{X}_\#)^{\mathbf{s}^{(0)}})$ among the list $T((\mathbf{X}_\#)^{\mathbf{s}^{(0)}}), \dots, T((\mathbf{X}_\#)^{\mathbf{s}^{(M)}})$ is uniformly distributed. Marginalizing over $\mathbf{X}_\#$, therefore,

$$\mathbb{P} \left\{ \hat{p}_M \left((\mathbf{X}_\#)^{\mathbf{s}^{(0)}}; \mathbf{s}^{(0)} \circ \mathbf{s}^{(1)}, \dots, \mathbf{s}^{(0)} \circ \mathbf{s}^{(M)} \right) \leq \alpha \mid \mathbf{Y}, \mathbf{Z} \right\} \leq \alpha.$$

Finally, combining this with our earlier calculation (22), we have

$$\mathbb{P} \left\{ \hat{p}_M(\mathbf{X}_\#; \mathbf{s}^{(1)}, \dots, \mathbf{s}^{(M)}) \leq \alpha \mid \mathbf{Y}, \mathbf{Z} \right\} \leq \alpha,$$

which completes the proof.

B Proof of Lemma 3

For $\ell \in [L]$, let

$$u_\ell = \mathbb{E} [\psi(X_{i_\ell}, X_{j_\ell}) \mid \mathbf{Y}, \mathbf{Z}] \quad \text{and} \quad v_\ell = (\text{Var}(\psi(X_{i_\ell}, X_{j_\ell}) \mid \mathbf{Y}, \mathbf{Z}))^{1/2}.$$

Then by Cauchy–Schwarz, for any $(w_1, \dots, w_L) \in [0, \infty)^L$,

$$\begin{aligned} \frac{\sum_{\ell=1}^L w_\ell \mathbb{E} [\psi(X_{i_\ell}, X_{j_\ell}) \mid \mathbf{Y}, \mathbf{Z}]}{\sqrt{\sum_{\ell=1}^L w_\ell^2 \text{Var}(\psi(X_{i_\ell}, X_{j_\ell}) \mid \mathbf{Y}, \mathbf{Z})}} &= \frac{\sum_{\ell=1}^L w_\ell u_\ell}{\sqrt{\sum_{\ell=1}^L w_\ell^2 v_\ell^2}} \leq \frac{\sum_{\ell=1}^L w_\ell \max\{u_\ell, 0\}}{\sqrt{\sum_{\ell=1}^L w_\ell^2 v_\ell^2}} \\ &= \frac{\sum_{\ell=1}^L (w_\ell v_\ell) \cdot \frac{\max\{u_\ell, 0\}}{v_\ell}}{\sqrt{\sum_{\ell=1}^L (w_\ell v_\ell)^2}} \leq \left(\sum_{\ell=1}^L \frac{\max\{u_\ell, 0\}^2}{v_\ell^2} \right)^{1/2}, \end{aligned}$$

with equality if and only if $w_\ell \propto \max\{u_\ell, 0\}/v_\ell^2$ for $\ell \in [L]$.

C Proof of the results from Section 4

In this section, we prove the results presented in Section 4. Throughout this appendix, we assume that the statistic T admits the form in (11) with $\psi(x, x') = x - x'$, and that $\mathcal{Y} = \mathcal{Z} = \mathbb{R}$; the partial ordering \preceq for Z will simply be the usual ordering \leq on \mathbb{R} . The organization of this appendix is as follows.

- We begin in Appendix C.1 by proving finite sample and asymptotic upper and lower bounds on the conditional power of PairSwap-ICI test for any valid matching and weighting scheme, and any statistic T of the form (11) with a shared linear kernel $\psi(x, x') = x - x'$.
- Next, in Appendix C.2, we specialize these results to the oracle matching under two cases: one assuming access to oracle knowledge of μ (i.e., Theorem 4) and another with μ estimated from data (i.e., Theorem 5).
- Then, we shift our attention to the partially linear Gaussian models in (20). In Appendices C.3 and C.4 we prove the asymptotic behavior of conditional power for neighbour matching (Theorem 10) and for cross bin matching (Theorem 11), respectively.

- In Appendix C.5 we prove the corollaries and lemmas from Section 4.
- In Appendix C.6 we introduce an oracle matching, namely isotonic median matching, and discuss a key property of the same, which allows us to prove the results in Appendix C.2.
- Finally, in Appendix C.7 we prove the lemmas stated in Appendices C.1—C.4.

Notation. We write $(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$ as shorthand for independent triples $(X_i, Y_i, Z_i)_{i \in [n]}$. For $L \equiv L_n \in \mathbb{N}$, a matching $\{(i_1, j_1), (i_2, j_2), \dots, (i_L, j_L)\} \in \mathcal{M}_n(\mathbf{Z})$ and a vector $\mathbf{V} = (V_1, \dots, V_n) \in \mathbb{R}^n$, we define $\Delta \mathbf{V} = (\Delta_1 \mathbf{V}, \dots, \Delta_L \mathbf{V}) \in \mathbb{R}^L$ with entries $\Delta_\ell \mathbf{V} := V_{i_\ell} - V_{j_\ell}$. Given any vector \mathbf{v} , we write \mathbf{v}^+ to denote the vector with i th component $v_{i+} = \max\{v_i, 0\}$. We write $\mathbf{a} \circ \mathbf{b}$ for the Hadamard product of vectors \mathbf{a}, \mathbf{b} of the same dimension, with i th component $a_i \cdot b_i$. For $k \geq 1$ and a distribution P_ζ on \mathbb{R} with finite k -th moment, let

$$\rho_k = \left(\mathbb{E}_{\substack{\zeta, \zeta' \sim P_\zeta \\ \zeta \perp \zeta'}} [|\zeta - \zeta'|^k] \right)^{1/k}. \quad (23)$$

In particular, $\rho_2 = \sqrt{2}\sigma$ where σ^2 is variance of P_ζ .

C.1 A general result on power of PairSwap-ICI test

Here, we consider any valid matching and weighting scheme, and state finite-sample lower and upper bounds on the conditional power $\mathbb{P}\{p \leq \alpha \mid \mathbf{Y}, \mathbf{Z}\}$ of PairSwap-ICI test. Further, under the asymptotic regime of Section 4.1, we state asymptotic high-probability upper and lower bounds for the same quantity. We first define several quantities in terms of the weight vector $\mathbf{w} = (w_1, \dots, w_L) \in \mathbb{R}^L$, $\delta > 0$ and $\{\rho_k : k \leq 6\}$, defined in (23): let

$$\begin{aligned} \epsilon_{1,\delta,U} &= \frac{\|\mathbf{w} \circ \Delta \mu(\mathbf{Y}, \mathbf{Z})\|_2^2}{\rho_2^2 \|\mathbf{w}\|_2^2} + \frac{\rho_4^2}{\rho_2^2 \delta^{1/2}} \frac{\|\mathbf{w}\|_\infty}{\|\mathbf{w}\|_2} + \frac{2}{\rho_2 \delta^{1/2}} \frac{\|\mathbf{w} \circ \Delta \mu(\mathbf{Y}, \mathbf{Z})\|_\infty}{\|\mathbf{w}\|_2}, \\ \epsilon_{1,\delta,L} &= \frac{\|\mathbf{w} \circ \Delta \mu(\mathbf{Y}, \mathbf{Z})\|_2^2}{\rho_2^2 \|\mathbf{w}\|_2^2} - \frac{\rho_4^2}{\rho_2^2 \delta^{1/2}} \frac{\|\mathbf{w}\|_\infty}{\|\mathbf{w}\|_2} - \frac{2}{\rho_2 \delta^{1/2}} \frac{\|\mathbf{w} \circ \Delta \mu(\mathbf{Y}, \mathbf{Z})\|_\infty}{\|\mathbf{w}\|_2}, \\ \epsilon_{2,\delta} &= \frac{0.56}{((1 + \epsilon_{1,\delta,L}) \vee 0)^{3/2}} \left\{ \left(\frac{\|\mathbf{w} \circ \Delta \mu(\mathbf{Y}, \mathbf{Z})\|_2}{\rho_2 \|\mathbf{w}\|_2} \right)^{2/3} \cdot \left(\frac{\|\mathbf{w} \circ \Delta \mu(\mathbf{Y}, \mathbf{Z})\|_\infty}{\rho_2 \|\mathbf{w}\|_2} \right)^{1/3} \right. \\ &\quad \left. + \left(\frac{\rho_3^3 \|\mathbf{w}\|_\infty}{\rho_2^3 \|\mathbf{w}\|_2} + \frac{\rho_6^3}{\delta^{1/2} \rho_2^3} \cdot \frac{\|\mathbf{w}\|_\infty^2}{\|\mathbf{w}\|_2^2} \right)^{1/3} \right\}^3, \\ \epsilon_{3,\delta} &= \frac{0.56 \rho_3^3}{\rho_2^3} \cdot \frac{\|\mathbf{w}\|_\infty}{\|\mathbf{w}\|_2} + 5\delta. \end{aligned} \quad (24)$$

Theorem C.12. *Suppose that $(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$ satisfies (16), and that $\mathbf{w} = \mathbf{w}(\mathbf{Y}, \mathbf{Z}) \in \mathbb{R}^n$ is our chosen weight vector, which is assumed to satisfy $\|\mathbf{w}\|_2 > 0$. Write $\Omega_0 := \{\|\mathbf{w} \circ \Delta \mathbf{X}\|_2 > 0\}$. Then*

(i) For any $\delta > 0$ and $\alpha \in (0, 1/2 - \epsilon_{2,\delta}]$,

$$\begin{aligned} & \Phi\left(\frac{\mathbf{w}^T \Delta\mu(\mathbf{Y}, \mathbf{Z})}{\rho_2 \|\mathbf{w}\|_2} - \bar{\Phi}^{-1}((\alpha - \epsilon_{2,\delta}) \vee 0)(1 + \epsilon_{1,\delta,U})^{1/2}\right) - \epsilon_{3,\delta} - \mathbb{P}\{\Omega_0^c \mid \mathbf{Y}, \mathbf{Z}\} \\ & \leq \mathbb{P}\{p \leq \alpha \mid \mathbf{Y}, \mathbf{Z}\} \leq \Phi\left(\frac{\mathbf{w}^T \Delta\mu(\mathbf{Y}, \mathbf{Z})}{\rho_2 \|\mathbf{w}\|_2} - \bar{\Phi}^{-1}(\alpha + \epsilon_{2,\delta}) \cdot ((1 + \epsilon_{1,\delta,L}) \vee 0)^{1/2}\right) + \epsilon_{3,\delta}, \end{aligned}$$

where $\epsilon_{1,\delta,U}, \epsilon_{1,\delta,L}, \epsilon_{2,\delta}$, and $\epsilon_{3,\delta}$ are as in (24).

(ii) Further, suppose that the weights and matching scheme satisfy

Assumption A1. $\|\mathbf{w}\|_\infty \vee (\|\mathbf{w} \circ \Delta\mu(\mathbf{Y}, \mathbf{Z})\|_2^{2/3} \cdot \|\mathbf{w} \circ \Delta\mu(\mathbf{Y}, \mathbf{Z})\|_\infty^{1/3}) = o_P(\|\mathbf{w}\|_2)$.

Assumption A2. $\|\mathbf{w} \circ \Delta\mu(\mathbf{Y}, \mathbf{Z})\|_2 = o_P(\mathbf{w}^T \Delta\mu(\mathbf{Y}, \mathbf{Z}))$.

Then for $\alpha \in (0, 1/2)$,

$$\begin{aligned} \Phi\left(\frac{\mathbf{w}^T \Delta\mu(\mathbf{Y}, \mathbf{Z})}{\sqrt{2}\sigma \|\mathbf{w}\|_2} - \bar{\Phi}^{-1}(\alpha)\right) - o_P(1) & \leq \mathbb{P}\{p \leq \alpha \mid \mathbf{Y}, \mathbf{Z}\} \\ & \leq \Phi\left(\frac{\mathbf{w}^T \Delta\mu(\mathbf{Y}, \mathbf{Z})}{\sqrt{2}\sigma \|\mathbf{w}\|_2} - \bar{\Phi}^{-1}(\alpha)\right) + o_P(1). \end{aligned}$$

Proof. We note that for fixed $n \in \mathbb{N}$, and conditional on $(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$, the quantity $\mathbf{s}^T(\mathbf{w} \circ \Delta\mathbf{X})$ is a sum of independent random variables $\{s_\ell \cdot w_\ell \cdot \Delta_\ell \mathbf{X}\}_{\ell \in [L]}$, and has mean 0 and variance $\|\mathbf{w} \circ \Delta\mathbf{X}\|_2^2$. Hence, writing $U := \frac{\mathbf{1}^T(\mathbf{w} \circ \Delta\mathbf{X})}{\|\mathbf{w} \circ \Delta\mathbf{X}\|_2} \mathbb{1}_{\Omega_0}$, we have by the Berry–Esseen theorem (Shevtsova, 2010, Theorem 1) that on Ω_0 ,

$$|p - \bar{\Phi}(U)| \leq \sup_{x \in \mathbb{R}} \left| \mathbb{P}\left\{ \frac{\mathbf{s}^T(\mathbf{w} \circ \Delta\mathbf{X})}{\|\mathbf{w} \circ \Delta\mathbf{X}\|_2} \geq x \mid \mathbf{X}, \mathbf{Y}, \mathbf{Z} \right\} - \bar{\Phi}(x) \right| \leq 0.56 \cdot \frac{\|\mathbf{w} \circ \Delta\mathbf{X}\|_3^3}{\|\mathbf{w} \circ \Delta\mathbf{X}\|_2^3}.$$

Hence, since $p = 1$ on Ω_0^c , we therefore have

$$\begin{aligned} \mathbb{P}\left[\left\{U \geq \bar{\Phi}^{-1}\left(\left(\alpha - 0.56 \frac{\|\mathbf{w} \circ \Delta\mathbf{X}\|_3^3}{\|\mathbf{w} \circ \Delta\mathbf{X}\|_2^3}\right) \vee 0\right)\right\} \cap \Omega_0 \mid \mathbf{Y}, \mathbf{Z}\right] & \leq \mathbb{P}\{p \leq \alpha \mid \mathbf{Y}, \mathbf{Z}\} \\ & \leq \mathbb{P}\left[\left\{U \geq \bar{\Phi}^{-1}\left(\left(\alpha + 0.56 \frac{\|\mathbf{w} \circ \Delta\mathbf{X}\|_3^3}{\|\mathbf{w} \circ \Delta\mathbf{X}\|_2^3}\right) \wedge 1\right)\right\} \cap \Omega_0 \mid \mathbf{Y}, \mathbf{Z}\right]. \end{aligned}$$

Now $\mathbf{1}^T(\mathbf{w} \circ \Delta\mathbf{X}) = \mathbf{w}^T \Delta\mathbf{X} = \mathbf{w}^T \Delta\mu(\mathbf{Y}, \mathbf{Z}) + \mathbf{w}^T \Delta\boldsymbol{\zeta}$ and $\mathbf{w}^T \Delta\mu(\mathbf{Y}, \mathbf{Z})$ is a measurable function of (\mathbf{Y}, \mathbf{Z}) . On the other hand, conditional on (\mathbf{Y}, \mathbf{Z}) , the quantity $\mathbf{w}^T \Delta\boldsymbol{\zeta}$ is a sum of independent random variables $\{w_\ell \cdot \Delta_\ell \boldsymbol{\zeta}\}_{\ell \in [L]}$, and has mean 0 and variance $\rho_2^2 \|\mathbf{w}\|_2^2$. Another application of the Berry–Esseen theorem then yields that

$$\sup_{x \in \mathbb{R}} \left| \mathbb{P}\left\{ \frac{\mathbf{w}^T \Delta\boldsymbol{\zeta}}{\rho_2 \|\mathbf{w}\|_2} \geq x \mid \mathbf{Y}, \mathbf{Z} \right\} - \bar{\Phi}(x) \right| \leq \frac{0.56 \rho_3^3}{\rho_2^3} \cdot \frac{\|\mathbf{w}\|_3^3}{\|\mathbf{w}\|_2^3} \leq \frac{0.56 \rho_3^3}{\rho_2^3} \cdot \frac{\|\mathbf{w}\|_\infty}{\|\mathbf{w}\|_2}. \quad (25)$$

Thus, recalling the events Ω_2 and Ω_3 from Lemma C.19 we have by that lemma that

$$\begin{aligned}
& \mathbb{P} \left[\left\{ U \geq \bar{\Phi}^{-1} \left(\left(\alpha + 0.56 \frac{\|\mathbf{w} \circ \Delta \mathbf{X}\|_3^3}{\|\mathbf{w} \circ \Delta \mathbf{X}\|_2^3} \right) \wedge 1 \right) \right\} \cap \Omega_0 \mid \mathbf{Y}, \mathbf{Z} \right] \\
& \leq \mathbb{P} \left[\left\{ \frac{\mathbf{w}^T \Delta \boldsymbol{\zeta}}{\rho_2 \|\mathbf{w}\|_2} \geq -\frac{\mathbf{w}^T \Delta \mu(\mathbf{Y}, \mathbf{Z})}{\rho_2 \|\mathbf{w}\|_2} + \frac{\|\mathbf{w} \circ \Delta \mathbf{X}\|_2}{\rho_2 \|\mathbf{w}\|_2} \bar{\Phi}^{-1}(\alpha + \epsilon_{2,\delta}) \right\} \cap \Omega_0 \mid \mathbf{Y}, \mathbf{Z} \right] + \mathbb{P} \{\Omega_3^c\} \\
& \leq \mathbb{P} \left[\left\{ \frac{\mathbf{w}^T \Delta \boldsymbol{\zeta}}{\rho_2 \|\mathbf{w}\|_2} \geq -\frac{\mathbf{w}^T \Delta \mu(\mathbf{Y}, \mathbf{Z})}{\rho_2 \|\mathbf{w}\|_2} + \bar{\Phi}^{-1}(\alpha + \epsilon_{2,\delta}) \left((1 + \epsilon_{1,\delta,L}) \vee 0 \right)^{1/2} \right\} \cap \Omega_0 \mid \mathbf{Y}, \mathbf{Z} \right] \\
& \quad + \mathbb{P} \{\Omega_2^c\} + \mathbb{P} \{\Omega_3^c\} \\
& \leq \Phi \left(\frac{\mathbf{w}^T \Delta \mu(\mathbf{Y}, \mathbf{Z})}{\rho_2 \|\mathbf{w}\|_2} - \bar{\Phi}^{-1}(\alpha + \epsilon_{2,\delta}) \cdot \left((1 + \epsilon_{1,\delta,L}) \vee 0 \right)^{1/2} \right) + \frac{0.56 \rho_3^3}{\rho_2^3} \cdot \frac{\|\mathbf{w}\|_\infty}{\|\mathbf{w}\|_2} + 5\delta \\
& \leq \Phi \left(\frac{\mathbf{w}^T \Delta \mu(\mathbf{Y}, \mathbf{Z})}{\rho_2 \|\mathbf{w}\|_2} - \bar{\Phi}^{-1}(\alpha + \epsilon_{2,\delta}) \cdot \left((1 + \epsilon_{1,\delta,L}) \vee 0 \right)^{1/2} \right) + \epsilon_{3,\delta}.
\end{aligned}$$

Similarly, recalling the events Ω_1, Ω_2 and Ω_3 from Lemma C.19 for the lower bound, and writing $\Omega_0^* = \Omega_0 \cap \Omega_1 \cap \Omega_3$,

$$\begin{aligned}
& \mathbb{P} \left[\left\{ U \geq \bar{\Phi}^{-1} \left(\left(\alpha - 0.56 \frac{\|\mathbf{w} \circ \Delta \mathbf{X}\|_3^3}{\|\mathbf{w} \circ \Delta \mathbf{X}\|_2^3} \right) \vee 0 \right) \right\} \cap \Omega_0 \mid \mathbf{Y}, \mathbf{Z} \right] \\
& \geq \mathbb{P} \left[\left\{ \frac{\mathbf{w}^T \Delta \boldsymbol{\zeta}}{\rho_2 \|\mathbf{w}\|_2} \geq -\frac{\mathbf{w}^T \Delta \mu(\mathbf{Y}, \mathbf{Z})}{\rho_2 \|\mathbf{w}\|_2} + \frac{\|\mathbf{w} \circ \Delta \mathbf{X}\|_2}{\rho_2 \|\mathbf{w}\|_2} \bar{\Phi}^{-1}((\alpha - \epsilon_{2,\delta}) \vee 0) \right\} \cap \Omega_0^* \mid \mathbf{Y}, \mathbf{Z} \right] \\
& \geq \mathbb{P} \left[\left\{ \frac{\mathbf{w}^T \Delta \boldsymbol{\zeta}}{\rho_2 \|\mathbf{w}\|_2} \geq -\frac{\mathbf{w}^T \Delta \mu(\mathbf{Y}, \mathbf{Z})}{\rho_2 \|\mathbf{w}\|_2} + \bar{\Phi}^{-1}((\alpha - \epsilon_{2,\delta}) \vee 0) (1 + \epsilon_{1,\delta,U})^{1/2} \right\} \cap \Omega_0^* \mid \mathbf{Y}, \mathbf{Z} \right] \\
& \geq \Phi \left(\frac{\mathbf{w}^T \Delta \mu(\mathbf{Y}, \mathbf{Z})}{\rho_2 \|\mathbf{w}\|_2} - \bar{\Phi}^{-1}((\alpha - \epsilon_{2,\delta}) \vee 0) (1 + \epsilon_{1,\delta,U})^{1/2} \right) - \frac{0.56 \rho_3^3 \|\mathbf{w}\|_\infty}{\rho_2^3 \|\mathbf{w}\|_2} - \mathbb{P} \{\Omega_0^{*c} \mid \mathbf{Y}, \mathbf{Z}\} \\
& \geq \Phi \left(\frac{\mathbf{w}^T \Delta \mu(\mathbf{Y}, \mathbf{Z})}{\rho_2 \|\mathbf{w}\|_2} - \bar{\Phi}^{-1}((\alpha - \epsilon_{2,\delta}) \vee 0) (1 + \epsilon_{1,\delta,U})^{1/2} \right) - \epsilon_{3,\delta} - \mathbb{P} \{\Omega_0^c \mid \mathbf{Y}, \mathbf{Z}\}.
\end{aligned}$$

This completes the proof of the first part of the result.

For the second part, define

$$\epsilon_{1,\delta} := \frac{\rho_4^2}{\rho_2^2 \delta^{1/2}} \cdot \frac{\|\mathbf{w}\|_\infty}{\|\mathbf{w}\|_2} + \frac{2}{\rho_2 \delta^{1/2}} \cdot \frac{\|\mathbf{w} \circ \Delta \mu(\mathbf{Y}, \mathbf{Z})\|_\infty}{\|\mathbf{w}\|_2} = o_P(1),$$

by Assumption A1, and observe that $(1 + \epsilon_{1,\delta,U})^{1/2} \leq \frac{\|\mathbf{w} \circ \Delta \mu(\mathbf{Y}, \mathbf{Z})\|_2}{\rho_2 \|\mathbf{w}\|_2} + (1 + \epsilon_{1,\delta})^{1/2}$ and $(1 + \epsilon_{1,\delta,L})^{1/2} \geq (1 - \epsilon_{1,\delta})^{1/2}$. It follows by Assumption A1 that $\epsilon_{2,\delta} = o_P(1)$. Moreover,

$$\begin{aligned}
\limsup_{n \rightarrow \infty} \mathbb{P} \{\Omega_0^c \mid \mathbf{Y}, \mathbf{Z}\} & \leq \limsup_{n \rightarrow \infty} \mathbb{P} \{\Omega_2^c \cap \{1 + \epsilon_{1,\delta,L} > 0\} \mid \mathbf{Y}, \mathbf{Z}\} + \limsup_{n \rightarrow \infty} \mathbb{P} \{1 + \epsilon_{1,\delta,L} \leq 0\} \\
& = \limsup_{n \rightarrow \infty} \mathbb{P} \{\Omega_2^c \mid \mathbf{Y}, \mathbf{Z}\} \leq 2\delta.
\end{aligned}$$

Now, $\epsilon_{3,\delta} - 5\delta = o_P(1)$ by Assumption **A1**. Thus, by Assumption **A2** and part (i),

$$\begin{aligned} & \mathbb{P}\{p \leq \alpha \mid \mathbf{Y}, \mathbf{Z}\} \\ & \geq \Phi\left(\frac{\mathbf{w}^T \Delta\mu(\mathbf{Y}, \mathbf{Z})}{\rho_2 \|\mathbf{w}\|_2} - \bar{\Phi}^{-1}(\alpha - o_P(1)) \left(\frac{\|\mathbf{w} \circ \Delta\mu(\mathbf{Y}, \mathbf{Z})\|_2}{\rho_2 \|\mathbf{w}\|_2} + (1 + o_P(1))\right)\right) - o_P(1) - 7\delta \\ & = \Phi\left(\left(\frac{\mathbf{w}^T \Delta\mu(\mathbf{Y}, \mathbf{Z})}{\rho_2 \|\mathbf{w}\|_2} - \bar{\Phi}^{-1}(\alpha) + o_P(1)\right)(1 + o_P(1))\right) - o_P(1) - 7\delta, \end{aligned}$$

and since $\delta > 0$ was arbitrary, the desired asymptotic lower bound holds by noting that $\rho_2 = \sqrt{2}\sigma$. The asymptotic upper bound follows by a very similar (in fact, slightly more straightforward) argument. \square

C.2 Proof of Theorems 4 and 5

We first consider the more practical setting where μ is estimated from data, and prove Theorem 5 using the general upper and lower bounds on the conditional power of the PairSwap-ICI test from Theorem C.12. Theorem 4 will then follow as a special case of Theorem 5.

Before turning to the proofs, we observe from Lemma 3 that under the model class (16) and with the linear kernel $\psi(x, x') = x - x'$, the oracle weights in (13) satisfy $w_\ell^* \propto (\mu(Y_{i_\ell}, Z_{i_\ell}) - \mu(Y_{j_\ell}, Z_{j_\ell}))_+$ for $\ell \in [L]$. Since the p -value of our test is independent of the scales of weights, we may take $\mathbf{w}^* = \Delta\mu^+(\mathbf{Y}, \mathbf{Z})$. When we do not have access to oracle knowledge of μ , we assume that we are able to learn estimates $\hat{\mu}$ and $\hat{\sigma}$ of μ and σ respectively from a prior dataset D_{prior} , independent of $(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$. With these estimates in place, the plug-in weights satisfy $\hat{w}_\ell \propto (\hat{\mu}(Y_{i_\ell}, Z_{i_\ell}) - \hat{\mu}(Y_{j_\ell}, Z_{j_\ell}))_+$, i.e., we may take $\hat{\mathbf{w}} = \Delta\hat{\mu}^+(\mathbf{Y}, \mathbf{Z})$.

Given $(\mathbf{Y}, \mathbf{Z}) \in \mathbb{R}^n \times \mathbb{R}^n$, we define

$$\hat{\mu}_{\text{ISO}}(\mathbf{Y}, \mathbf{Z}) = \operatorname{argmin}_{g \in \mathcal{C}_{\text{ISO}}} \|\mu(\mathbf{Y}, \mathbf{Z}) - g(\mathbf{Z})\|_2, \quad (26)$$

the empirical isotonic Euclidean projection of μ onto \mathcal{C}_{ISO} .

C.2.1 Proof of Theorem 5

Our proof is split into three steps. First we establish a key property of the oracle matching M^* and use this to deduce relative error properties of $\hat{\mu}$. These in turn enable us to show that Assumptions **A1** and **A2** of Theorem C.12 are satisfied by \hat{M} and weight vector $\mathbf{w} = \Delta\hat{\mu}^+(\mathbf{Y}, \mathbf{Z})$. Finally we apply part (ii) of Theorem C.12 to conclude the proof.

We write

$$\begin{aligned} \text{Err}_2(\hat{\mu}, \mu) &= \frac{\|\Delta\hat{\mu}(\mathbf{Y}, \mathbf{Z}) - \Delta\mu(\mathbf{Y}, \mathbf{Z})\|_2}{\|\Delta\mu^+(\mathbf{Y}, \mathbf{Z})\|_2}, \\ \text{Err}_\infty(\hat{\mu}, \mu) &= \frac{\|\Delta\hat{\mu}(\mathbf{Y}, \mathbf{Z}) - \Delta\mu(\mathbf{Y}, \mathbf{Z})\|_\infty^4}{\|\Delta\mu^+(\mathbf{Y}, \mathbf{Z})\|_2}. \end{aligned}$$

Step 1: a key property of oracle matching and its implications. Our notation $\Delta\mu(\mathbf{Y}, \mathbf{Z})$ and $\Delta\hat{\mu}(\mathbf{Y}, \mathbf{Z})$ suppresses the dependence of the matching M for which it is computed. Since we will need both matchings M^* and \hat{M} in this step, we make the dependence on the matching explicit by writing $\Delta\hat{\mu}(\mathbf{Y}, \mathbf{Z}; M)$, where $M \in \mathcal{M}_n(\mathbf{Z})$.

By Lemma C.13, Theorem C.18 and the definition of $M^* \in \mathcal{M}_n(\mathbf{Z})$ from (14), we have

$$\widehat{\text{ISS}}_n \leq \|\Delta\mu^+(\mathbf{Y}, \mathbf{Z}; M^*)\|_2 = \sup_{M \in \mathcal{M}_n(\mathbf{Z})} \|\Delta\mu^+(\mathbf{Y}, \mathbf{Z}; M)\|_2 \leq \sqrt{2} \widehat{\text{ISS}}_n. \quad (27)$$

Now, for any $M \in \mathcal{M}_n(\mathbf{Z})$,

$$\begin{aligned} \left| \|\Delta\hat{\mu}^+(\mathbf{Y}, \mathbf{Z}; M)\|_2 - \|\Delta\mu^+(\mathbf{Y}, \mathbf{Z}; M)\|_2 \right| &\leq \|\Delta\hat{\mu}^+(\mathbf{Y}, \mathbf{Z}; M) - \Delta\mu^+(\mathbf{Y}, \mathbf{Z}; M)\|_2 \\ &\leq \|\Delta\hat{\mu}(\mathbf{Y}, \mathbf{Z}; M) - \Delta\mu(\mathbf{Y}, \mathbf{Z}; M)\|_2 \\ &\leq \sqrt{2} \|\hat{\mu}(\mathbf{Y}, \mathbf{Z}) - \mu(\mathbf{Y}, \mathbf{Z})\|_2 = o_P(\widehat{\text{ISS}}_n), \end{aligned} \quad (28)$$

where the final step is from the hypothesis in the statement of the result. Further by (15),

$$\|\Delta\hat{\mu}^+(\mathbf{Y}, \mathbf{Z}; \hat{M})\|_2 \geq \|\Delta\hat{\mu}^+(\mathbf{Y}, \mathbf{Z}; M^*)\|_2 \geq \|\Delta\mu^+(\mathbf{Y}, \mathbf{Z}; M^*)\|_2 - o_P(\widehat{\text{ISS}}_n) \geq \widehat{\text{ISS}}_n(1 - o_P(1)). \quad (29)$$

Similarly, by (28),

$$\|\Delta\mu^+(\mathbf{Y}, \mathbf{Z}; \hat{M})\|_2 \geq \|\Delta\hat{\mu}^+(\mathbf{Y}, \mathbf{Z}; \hat{M})\|_2 - o_P(\widehat{\text{ISS}}_n) \geq \widehat{\text{ISS}}_n(1 - o_P(1)). \quad (30)$$

We deduce that for plug-in matching,

$$\max \left\{ \frac{\|\mu(\mathbf{Y}, \mathbf{Z})\|_\infty \vee \|\mu(\mathbf{Y}, \mathbf{Z})\|_\infty^4}{\|\Delta\mu^+(\mathbf{Y}, \mathbf{Z})\|_2}, \text{Err}_2(\hat{\mu}, \mu), \text{Err}_\infty(\hat{\mu}, \mu) \right\} = o_P(1). \quad (31)$$

Step 2: establishing Assumptions A1 and A2. Henceforth we work with the matching \hat{M} . We have by (27) that

$$\begin{aligned} &\|\Delta\hat{\mu}^+(\mathbf{Y}, \mathbf{Z}) \circ \Delta\mu(\mathbf{Y}, \mathbf{Z})\|_2 \\ &\leq \|\Delta\mu^+(\mathbf{Y}, \mathbf{Z}) \circ \Delta\mu(\mathbf{Y}, \mathbf{Z})\|_2 + \|(\Delta\hat{\mu}^+(\mathbf{Y}, \mathbf{Z}) - \Delta\mu^+(\mathbf{Y}, \mathbf{Z})) \circ \Delta\mu(\mathbf{Y}, \mathbf{Z})\|_2 \\ &\leq \|\Delta\mu^+(\mathbf{Y}, \mathbf{Z})\|_4^2 + \|\Delta\mu(\mathbf{Y}, \mathbf{Z})\|_\infty \cdot \|\Delta\hat{\mu}(\mathbf{Y}, \mathbf{Z}) - \Delta\mu(\mathbf{Y}, \mathbf{Z})\|_2 \\ &\leq 2\|\mu(\mathbf{Y}, \mathbf{Z})\|_\infty \|\Delta\mu^+(\mathbf{Y}, \mathbf{Z})\|_2 \cdot (1 + \text{Err}_2(\hat{\mu}, \mu)) = o_P(\widehat{\text{ISS}}_n^2). \end{aligned}$$

Moreover, by the Cauchy–Schwarz inequality, (27) and (28),

$$\begin{aligned} \left| \Delta\hat{\mu}^+(\mathbf{Y}, \mathbf{Z})^T \Delta\mu(\mathbf{Y}, \mathbf{Z}) - \|\Delta\hat{\mu}^+(\mathbf{Y}, \mathbf{Z})\|_2^2 \right| &= |(\Delta\hat{\mu}(\mathbf{Y}, \mathbf{Z}) - \Delta\mu(\mathbf{Y}, \mathbf{Z}))^T \Delta\hat{\mu}^+(\mathbf{Y}, \mathbf{Z})| \\ &\leq \|\Delta\hat{\mu}^+(\mathbf{Y}, \mathbf{Z})\|_2 \cdot \|\Delta\mu^+(\mathbf{Y}, \mathbf{Z})\|_2 \cdot \text{Err}_2(\hat{\mu}, \mu) \\ &= o_P(\widehat{\text{ISS}}_n^2). \end{aligned} \quad (32)$$

Hence, by (31), Assumption A2 is satisfied by \hat{M} . Moreover, (31) also yields that

$$\begin{aligned} \frac{\|\Delta\hat{\mu}^+(\mathbf{Y}, \mathbf{Z})\|_\infty}{\|\Delta\hat{\mu}^+(\mathbf{Y}, \mathbf{Z})\|_2} &\leq \frac{\|\Delta\mu^+(\mathbf{Y}, \mathbf{Z})\|_\infty + \|\Delta\hat{\mu}(\mathbf{Y}, \mathbf{Z}) - \Delta\mu(\mathbf{Y}, \mathbf{Z})\|_2}{\|\Delta\mu^+(\mathbf{Y}, \mathbf{Z})\|_2 - \|\Delta\hat{\mu}(\mathbf{Y}, \mathbf{Z}) - \Delta\mu(\mathbf{Y}, \mathbf{Z})\|_2} \\ &\leq \frac{2(\|\mu(\mathbf{Y}, \mathbf{Z})\|_\infty / \|\Delta\mu^+(\mathbf{Y}, \mathbf{Z})\|_2) + \text{Err}_2(\hat{\mu}, \mu)}{1 - \text{Err}_2(\hat{\mu}, \mu)} = o_P(1). \end{aligned}$$

Now, we focus on the other term that appears in Assumption **A1**. By (31) and (27),

$$\begin{aligned}
& \|\Delta\hat{\mu}^+(\mathbf{Y}, \mathbf{Z}) \circ \Delta\mu(\mathbf{Y}, \mathbf{Z})\|_2^{2/3} \\
& \leq \left(\|\Delta\mu^+(\mathbf{Y}, \mathbf{Z}) \circ \Delta\mu^+(\mathbf{Y}, \mathbf{Z})\|_2 + \|(\Delta\hat{\mu}^+(\mathbf{Y}, \mathbf{Z}) - \Delta\mu^+(\mathbf{Y}, \mathbf{Z})) \circ \Delta\mu(\mathbf{Y}, \mathbf{Z})\|_2 \right)^{2/3} \\
& \leq \left(\|\Delta\mu^+(\mathbf{Y}, \mathbf{Z})\|_4^2 + \text{Err}_2(\hat{\mu}, \mu) \cdot \|\Delta\mu^+(\mathbf{Y}, \mathbf{Z})\|_2 \cdot \|\Delta\mu(\mathbf{Y}, \mathbf{Z})\|_\infty \right)^{2/3} \\
& \leq \|\Delta\mu(\mathbf{Y}, \mathbf{Z})\|_\infty^{2/3} \cdot \|\Delta\mu^+(\mathbf{Y}, \mathbf{Z})\|_2^{2/3} \cdot (1 + o_P(1)) \\
& \leq \|\Delta\mu(\mathbf{Y}, \mathbf{Z})\|_\infty^{2/3} \cdot 2^{1/3} \widehat{\text{ISS}}_n^{2/3} \cdot (1 + o_P(1)).
\end{aligned}$$

Moreover, again by (31),

$$\begin{aligned}
& \|\Delta\hat{\mu}^+(\mathbf{Y}, \mathbf{Z}) \circ \Delta\mu(\mathbf{Y}, \mathbf{Z})\|_\infty^{1/3} \\
& \leq \|\Delta\mu^+(\mathbf{Y}, \mathbf{Z})\|_\infty^{2/3} + \|\Delta\hat{\mu}(\mathbf{Y}, \mathbf{Z}) - \Delta\mu(\mathbf{Y}, \mathbf{Z})\|_\infty^{1/3} \|\Delta\mu(\mathbf{Y}, \mathbf{Z})\|_\infty^{1/3} \\
& \leq \frac{\|\Delta\mu(\mathbf{Y}, \mathbf{Z})\|_2^{1/3}}{\|\Delta\mu(\mathbf{Y}, \mathbf{Z})\|_\infty^{2/3}} \left(\frac{2^{4/3} \|\mu(\mathbf{Y}, \mathbf{Z})\|_\infty^{4/3}}{\|\Delta\mu(\mathbf{Y}, \mathbf{Z})\|_2^{1/3}} + \text{Err}_\infty^{1/12}(\hat{\mu}, \mu) \frac{2\|\mu(\mathbf{Y}, \mathbf{Z})\|_\infty}{\|\Delta\mu^+(\mathbf{Y}, \mathbf{Z})\|_2^{1/4}} \right) \\
& = o_P \left(\frac{\widehat{\text{ISS}}_n^{1/3}}{\|\Delta\mu(\mathbf{Y}, \mathbf{Z})\|_\infty^{2/3}} \right).
\end{aligned}$$

Hence, by (29) we see that Assumption **A1** is satisfied by plug-in matching.

Step 3: applying Theorem C.12. By (28), (29) and (32),

$$\left| \frac{\Delta\hat{\mu}^+(\mathbf{Y}, \mathbf{Z})^T \Delta\mu(\mathbf{Y}, \mathbf{Z})}{\|\Delta\hat{\mu}^+(\mathbf{Y}, \mathbf{Z})\|_2} - \|\Delta\mu^+(\mathbf{Y}, \mathbf{Z})\|_2 \right| = o_P(\widehat{\text{ISS}}_n).$$

Hence, by (27) and (30),

$$\widehat{\text{ISS}}_n (1 - o_P(1)) \leq \frac{\Delta\hat{\mu}^+(\mathbf{Y}, \mathbf{Z})^T \Delta\mu(\mathbf{Y}, \mathbf{Z})}{\|\Delta\hat{\mu}^+(\mathbf{Y}, \mathbf{Z})\|_2} \leq \sqrt{2} \cdot \widehat{\text{ISS}}_n (1 + o_P(1)).$$

Finally the result follows from part (ii) of Theorem C.12. \square

Lemma C.13. For any $M \in \mathcal{M}_n(\mathbf{Z})$, we have $\|\Delta\mu^+(\mathbf{Y}, \mathbf{Z})\|_2 \leq \sqrt{2} \widehat{\text{ISS}}_n$.

Proof. By definition, $\widehat{\text{ISS}}_n = \|\mu(\mathbf{Y}, \mathbf{Z}) - \hat{\mu}_{\text{ISO}}(\mathbf{Z})\|_2$ where $\hat{\mu}_{\text{ISO}}$ is as defined in (26). Recall that $(\Delta\mu^+(\mathbf{Y}, \mathbf{Z}))_\ell = (\mu(Y_{i_\ell}, Z_{i_\ell}) - \mu(Y_{j_\ell}, Z_{j_\ell}))_+$ for $\ell \in [L]$. We claim that

$$\left(\mu(Y_{i_\ell}, Z_{i_\ell}) - \mu(Y_{j_\ell}, Z_{j_\ell}) \right)_+ \leq |\mu(Y_{i_\ell}, Z_{i_\ell}) - \hat{\mu}_{\text{ISO}}(Z_{i_\ell})| + |\hat{\mu}_{\text{ISO}}(Z_{j_\ell}) - \mu(Y_{j_\ell}, Z_{j_\ell})| \quad (33)$$

for every $\ell \in [L]$. To see this, we may assume that $\mu(Y_{i_\ell}, Z_{i_\ell}) > \mu(Y_{j_\ell}, Z_{j_\ell})$ since otherwise the left-hand side is zero. Now, if $\hat{\mu}_{\text{ISO}}(Z_{i_\ell}) \geq \mu(Y_{j_\ell}, Z_{j_\ell})$, then

$$\begin{aligned}
\mu(Y_{i_\ell}, Z_{i_\ell}) - \mu(Y_{j_\ell}, Z_{j_\ell}) & \leq |\mu(Y_{i_\ell}, Z_{i_\ell}) - \hat{\mu}_{\text{ISO}}(Z_{i_\ell})| + \hat{\mu}_{\text{ISO}}(Z_{i_\ell}) - \mu(Y_{j_\ell}, Z_{j_\ell}) \\
& \leq |\mu(Y_{i_\ell}, Z_{i_\ell}) - \hat{\mu}_{\text{ISO}}(Z_{i_\ell})| + \hat{\mu}_{\text{ISO}}(Z_{j_\ell}) - \mu(Y_{j_\ell}, Z_{j_\ell}).
\end{aligned}$$

Otherwise, if $\hat{\mu}_{\text{ISO}}(Z_{i_\ell}) < \mu(Y_{j_\ell}, Z_{j_\ell})$, then

$$\begin{aligned} \mu(Y_{i_\ell}, Z_{i_\ell}) - \mu(Y_{j_\ell}, Z_{j_\ell}) &< \mu(Y_{i_\ell}, Z_{i_\ell}) - \hat{\mu}_{\text{ISO}}(Z_{i_\ell}) \\ &\leq \mu(Y_{i_\ell}, Z_{i_\ell}) - \hat{\mu}_{\text{ISO}}(Z_{i_\ell}) + |\hat{\mu}_{\text{ISO}}(Z_{j_\ell}) - \mu(Y_{j_\ell}, Z_{j_\ell})|. \end{aligned}$$

This proves the claim (33), and it follows that

$$\begin{aligned} \|\Delta\mu^+(\mathbf{Y}, \mathbf{Z})\|_2^2 &\leq 2 \sum_{\ell=1}^L \{\mu(Y_{i_\ell}, Z_{i_\ell}) - \hat{\mu}_{\text{ISO}}(Z_{i_\ell})\}^2 + 2 \sum_{\ell=1}^L \{\mu(Y_{j_\ell}, Z_{j_\ell}) - \hat{\mu}_{\text{ISO}}(Z_{j_\ell})\}^2 \\ &\leq 2 \sum_{i=1}^n \{\mu(Y_i, Z_i) - \hat{\mu}_{\text{ISO}}(Z_i)\}^2 = 2\|\text{Res}_n(\mathbf{Y}, \mathbf{Z})\|_2^2, \end{aligned}$$

which proves the result. \square

C.2.2 Proof of Theorem 4

Given access to oracle knowledge of μ , we observe that $\hat{\mu} = \mu$ satisfies the assumptions of Theorem 5. Hence, Theorem 4 follows as an immediate corollary of Theorem 5. \square

C.3 Proof of Theorem 10

Let us define π as any permutation of $[n]$ for which

$$Z_{\pi(1)} \leq Z_{\pi(2)} \leq \dots \leq Z_{\pi(n)},$$

and note that the collection of matched pairs for neighbour matching is given by $M := \{(\pi(2\ell - 1), \pi(2\ell))\}_{\ell \in [n/2]}$. By part (ii) of Theorem C.12, with $\mathbf{w} = \Delta\mathbf{Y}^+$ the dominant term in the upper and lower bound on conditional power reduces to

$$\frac{\Delta\mathbf{Y}^{+T} \Delta\mu(\mathbf{Y}, \mathbf{Z})}{\rho_2 \|\Delta\mathbf{Y}^+\|_2}.$$

Since $\rho_2 = \sqrt{2}\sigma$, by the same result it suffices to prove a concentration of the dominant term as

$$\frac{\Delta\mathbf{Y}^{+T} \Delta\mu(\mathbf{Y}, \mathbf{Z})}{\|\Delta\mathbf{Y}^+\|_2} = \sqrt{n/2} \cdot \beta_n \cdot \{\mathbb{E}[\text{Var}(Y | Z)]\}^{1/2} + o_P(1),$$

and that the Assumptions A1 and A2 are satisfied.

Step 1: concentration of the dominant term. Under the partially linear Gaussian models (20),

$$\Delta\mathbf{Y}^{+T} \Delta\mu(\mathbf{Y}, \mathbf{Z}) = \Delta\mathbf{Y}^{+T} \Delta\mu_0(\mathbf{Z}) + \beta_n \|\Delta\mathbf{Y}^+\|_2^2,$$

where the first term is negative since $\mu_0 \in \mathcal{C}_{\text{ISO}}$. In fact, we can say more. Firstly,

$$|\Delta\mathbf{Y}^{+T} \Delta\mu_0(\mathbf{Z})| \leq \|\Delta\mathbf{Y}^+ \circ \Delta\mu_0(\mathbf{Z})\|_1 \leq \|\Delta\mathbf{Y}^+\|_\infty \|\Delta\mu_0(\mathbf{Z})\|_1.$$

Since $\mu_0 \in \mathcal{C}_{\text{ISO}}$, we further have that

$$\begin{aligned} \|\Delta\mu_0(\mathbf{Z})\|_1 &= \sum_{\ell=1}^{\lfloor n/2 \rfloor} (\mu_0(Z_{\pi(2\ell)}) - \mu_0(Z_{\pi(2\ell-1)})) \leq \sum_{i=2}^n (\mu_0(Z_{\pi(i)}) - \mu_0(Z_{\pi(i-1)})) \\ &= \mu_0(Z_{\pi(n)}) - \mu_0(Z_{\pi(1)}) \leq 2\|\mu_0(\mathbf{Z})\|_\infty. \end{aligned} \quad (34)$$

Moreover by Lemma C.20, we have that

$$\left| \frac{1}{n/2} \|\Delta\mathbf{Y}^+\|_2^2 - \mathbb{E}[\text{Var}(Y | Z)] \right| = o_P(1),$$

which also implies that $\|\Delta\mathbf{Y}^+\|_2 \geq \sqrt{n/2} \cdot \{\mathbb{E}[\text{Var}(Y | Z)]\}^{1/2} - o_P(\sqrt{n})$. Since $Y \in [-1, 1]$ and $\mu_0(\mathbf{Z})$ is sub-Gaussian, it follows that

$$\begin{aligned} &\left| \frac{\Delta\mathbf{Y}^{+T} \Delta\mu(\mathbf{Y}, \mathbf{Z})}{\|\Delta\mathbf{Y}^+\|_2} - \sqrt{n/2} \cdot \beta_n \cdot \{\mathbb{E}[\text{Var}(Y | Z)]\}^{1/2} \right| \\ &\leq \left| \frac{\Delta\mathbf{Y}^{+T} \Delta\mu(\mathbf{Y}, \mathbf{Z})}{\|\Delta\mathbf{Y}^+\|_2} - \beta_n \|\Delta\mathbf{Y}^+\|_2 \right| + \left| \beta_n \|\Delta\mathbf{Y}^+\|_2 - \sqrt{n/2} \cdot \beta_n \cdot \{\mathbb{E}[\text{Var}(Y | Z)]\}^{1/2} \right| \\ &\leq 4 \frac{\|\mu_0(\mathbf{Z})\|_\infty}{\|\Delta\mathbf{Y}^+\|_2} + o_P(\beta_n \sqrt{n}) = o_P(1), \end{aligned}$$

where the last equality follows by recalling that $\mu_0(Z)$ is sub-Gaussian and $\beta_n \gtrsim n^{-1/2}$.

Step 2: establishing Assumptions A1 and A2. Since $Y \in [-1, 1]$,

$$\|\Delta\mathbf{Y}^+ \circ \Delta\mu(\mathbf{Y}, \mathbf{Z})\|_\infty \leq \|\Delta\mathbf{Y}^+ \circ \Delta\mu_0(\mathbf{Z})\|_\infty + \beta_n \|\Delta\mathbf{Y}^+\|_\infty^2 \leq 4(\|\mu_0(\mathbf{Z})\|_\infty + \beta_n).$$

Moreover, by an argument, similar to that in (34), we have that

$$\begin{aligned} \|\Delta\mathbf{Y}^+ \circ \Delta\mu(\mathbf{Y}, \mathbf{Z})\|_2 &\leq \|\Delta\mathbf{Y}^+ \circ \Delta\mu_0(\mathbf{Z})\|_2 + \beta_n \|\Delta\mathbf{Y}^+\|_4^2 \\ &\leq \|\Delta\mathbf{Y}^+\|_\infty \|\Delta\mu_0(\mathbf{Z})\|_\infty \|\Delta\mu_0(\mathbf{Z})\|_1 + \beta_n \|\Delta\mathbf{Y}^+\|_\infty \|\Delta\mathbf{Y}^+\|_2 \\ &\leq \|\Delta\mathbf{Y}^+\|_\infty \left(\|\Delta\mu_0(\mathbf{Z})\|_\infty \sum_{i=2}^n (\mu_0(Z_{\pi(i)}) - \mu_0(Z_{\pi(i-1)})) + \beta_n \|\Delta\mathbf{Y}^+\|_2 \right) \\ &\leq 2(4\|\mu_0(\mathbf{Z})\|_\infty^2 + \beta_n \|\Delta\mathbf{Y}^+\|_2). \end{aligned}$$

Since $\mu_0(Z)$ is sub-Gaussian and that $\|\Delta\mathbf{Y}^+\|_2 \geq \sqrt{n/2} \cdot \{\mathbb{E}[\text{Var}(Y | Z)]\}^{1/2} - o_P(\sqrt{n})$, Assumption A1 is satisfied by the neighbour matching. Similarly, by (34) and recalling that $\beta_n \gtrsim 1/\sqrt{n}$, it follows that

$$\frac{\|\Delta\mathbf{Y}^+ \circ \Delta\mu(\mathbf{Y}, \mathbf{Z})\|_2}{\Delta\mathbf{Y}^{+T} \Delta\mu(\mathbf{Y}, \mathbf{Z})} \leq \frac{2(4\|\mu_0(\mathbf{Z})\|_\infty^2 + \beta_n \|\Delta\mathbf{Y}^+\|_2)}{-4\|\mu_0(\mathbf{Z})\|_\infty + \beta_n \|\Delta\mathbf{Y}^+\|_2^2} = o_P(1).$$

Hence Assumption A2 holds too, which completes the proof. \square

C.4 Power of cross bin matching

In this section, we study power of the `PairSwap-ICI` test for cross-bin matching. We will start with stating and proving a concentration result for the conditional power of cross-bin matching in Theorem C.14. This subsequently leads to the proof of Theorem 11 in Appendix C.4.2. Before diving into the main results, we establish some necessary notations.

Definition 1. For any distribution P , we define the deviation of P as

$$\text{Dev}(P) := \mathbb{E}_{q \sim \text{Unif}[0,1]} [(F^{-1}(1-q) - F^{-1}(q))^2], \quad (35)$$

where F^{-1} is the generalized inverse CDF of P .

C.4.1 A general version of Theorem 11: asymptotic conditional power of cross bin matching

Theorem C.14. Under the model class (20) suppose $\beta_n \gtrsim 1/\sqrt{n}$. Additionally, also assume that

- (i) $\mu_0(Z)$ is a sub-Gaussian random variable, and the function μ_0 satisfies the Lipschitz property

(36)

- (ii) there exists a constant $L_W \geq 0$ such that for any $z_1, z_2 \in \mathbb{R}$,

$$d_{W_1}(P_{Y|Z}(\cdot | z_1), P_{Y|Z}(\cdot | z_2)) \leq L_W |P_Z(z_1) - P_Z(z_2)|, \quad (37)$$

where $d_{W_1}(\cdot, \cdot)$ is the 1-Wasserstein distance.

Then, the conditional power of cross-bin matching (Algorithm 2), implemented with kernel $\psi(x, x') = x - x'$ and weights $w_\ell = \max\{Y_{i_\ell} - Y_{j_\ell}, 0\}$, satisfies

$$|\mathbb{P}\{p \leq \alpha \mid \mathbf{Y}, \mathbf{Z}\} - \Phi\left(\sqrt{n/2} \cdot (\beta_n/\sigma) \cdot (\mathbb{E}[\text{Dev}(P_{Y|Z})])^{1/2} - \bar{\Phi}^{-1}(\alpha)\right)| = o_P(1).$$

Proof. The proof follows the same key steps as from the Appendix C.3. By part (ii) of Theorem C.12 and with replacing $\mathbf{w} = \Delta \mathbf{Y}^+$, the dominant term in the upper and lower bound on conditional power reduces to

$$\frac{\Delta \mathbf{Y}^{+T} \Delta \mu(\mathbf{Y}, \mathbf{Z})}{\rho_2 \|\Delta \mathbf{Y}^+\|_2},$$

Since $\rho_2 = \sqrt{2}\sigma$, it suffices to prove a concentration of the dominant term as

$$\frac{\Delta \mathbf{Y}^{+T} \Delta \mu(\mathbf{Y}, \mathbf{Z})}{\|\Delta \mathbf{Y}^+\|_2} = \sqrt{n/2} \cdot \beta_n \cdot \{\mathbb{E}[\text{Dev}(P_{Y|Z})]\}^{1/2} + o_P(1),$$

and that the Assumptions **A1** and **A2** are satisfied.

Step 1: concentration of the dominant term. Under the partially linear Gaussian model (20),

$$\Delta \mathbf{Y}^{+T} \Delta \mu(\mathbf{Y}, \mathbf{Z}) = \Delta \mathbf{Y}^{+T} \Delta \mu_0(\mathbf{Z}) + \beta \|\Delta \mathbf{Y}^+\|_2^2,$$

where for the first term, it holds that

$$|\Delta \mathbf{Y}^{+T} \Delta \mu_0(\mathbf{Z})| = \sum_{\ell=1}^{L_n} |Y_{i_{\ell,n}} - Y_{j_{\ell,n}}| |\mu_0(Z_{j_{\ell,n}}) - \mu_0(Z_{i_{\ell,n}})| \leq n \|\Delta \mathbf{Y}^+\|_{\infty} \|\Delta \mu_0(\mathbf{Z})\|_{\infty}. \quad (38)$$

Next, we analyze the large sample behaviour of $\|\Delta \mathbf{Y}^+\|_2^2$. Conditioned on \mathbf{Z} , $\|\Delta \mathbf{Y}^+\|_2^2$ is a sum of independent and uniformly bounded terms. Hence, by the law of large numbers,

$$\left| \frac{1}{n/2} \|\Delta \mathbf{Y}^+\|_2^2 - \mathbb{E} \left[\frac{1}{n/2} \|\Delta \mathbf{Y}^+\|_2^2 \mid \mathbf{Z} \right] \right| = o_P(1).$$

Finally, by Lemma C.21 we have that

$$\left| \frac{1}{n/2} \|\Delta \mathbf{Y}^+\|_2^2 - \mathbb{E} [\text{Dev}(P_{Y|Z})] \right| = o_P(1). \quad (39)$$

Since μ_0 is L_{μ} - Lipschitz, we note that $\|\Delta \mu_0(\mathbf{Z})\|_{\infty} = O_P(1/\sqrt{n})$. Combining behaviour of both terms, we have that

$$\begin{aligned} \frac{\Delta \mathbf{Y}^{+T} \Delta \mu(\mathbf{Y}, \mathbf{Z})}{\rho_2 \|\Delta \mathbf{Y}^+\|_2} &= \frac{-O_P(\sqrt{n}) + \beta \|\Delta \mathbf{Y}^+\|_2^2}{\rho_2 \|\Delta \mathbf{Y}^+\|_2} = (\beta/\rho_2) \cdot \|\Delta \mathbf{Y}^+\|_2 - o_P(1) \\ &= (\beta/\rho_2) \left(\sqrt{(n/2)} \cdot \sqrt{\mathbb{E} [\text{Dev}(P_{Y|Z})]} + o_P(\sqrt{n}) \right) - o_P(1). \end{aligned}$$

Since $\rho_2 = \sqrt{2}\sigma$, the result follows by Theorem C.12 as long as we can show that Assumption A1 and A2 are satisfied.

Step 2: establishing Assumptions A1 and A2. Since $Y \in [-1, 1]$,

$$\|\Delta \mathbf{Y}^+ \circ \Delta \mu(\mathbf{Y}, \mathbf{Z})\|_{\infty} \leq \|\Delta \mathbf{Y}^+ \circ \Delta \mu_0(\mathbf{Z})\|_{\infty} + \beta_n \|\Delta \mathbf{Y}^+\|_{\infty}^2 \leq 4(\|\mu_0(\mathbf{Z})\|_{\infty} + \beta_n).$$

Similarly, we also have that

$$\begin{aligned} \|\Delta \mathbf{Y}^+ \circ \Delta \mu(\mathbf{Y}, \mathbf{Z})\|_2 &\leq \|\Delta \mathbf{Y}^+ \circ \Delta \mu_0(\mathbf{Z})\|_2 + \beta_n \|\Delta \mathbf{Y}^+\|_4^2 \\ &\leq \|\Delta \mathbf{Y}^+\|_2 (\|\Delta \mu_0(\mathbf{Z})\|_{\infty} + \beta_n \|\Delta \mathbf{Y}^+\|_{\infty}) \\ &\leq 2\|\Delta \mathbf{Y}^+\|_2 (\|\mu_0(\mathbf{Z})\|_{\infty} + \beta_n). \end{aligned}$$

Since $\mu_0(Z)$ is sub-Gaussian and that $\|\Delta \mathbf{Y}^+\|_2 \geq \sqrt{n/2} \cdot \{\mathbb{E} [\text{Dev}(P_{Y|Z})]\}^{1/2} - o_P(\sqrt{n})$, Assumption A1 is satisfied by the cross-bin matching. Moreover, by (38) and recalling that $\beta_n \gtrsim 1/\sqrt{n}$, it follows that

$$\frac{\|\Delta \mathbf{Y}^+ \circ \Delta \mu(\mathbf{Y}, \mathbf{Z})\|_2}{\Delta \mathbf{Y}^{+T} \Delta \mu(\mathbf{Y}, \mathbf{Z})} \leq \frac{2\|\Delta \mathbf{Y}^+\|_2 (\|\mu_0(\mathbf{Z})\|_{\infty} + \beta_n)}{-n\|\Delta \mathbf{Y}^+\|_{\infty} \|\Delta \mu_0(\mathbf{Z})\|_{\infty} + \beta_n \|\Delta \mathbf{Y}^+\|_2^2} = o_P(1).$$

Hence Assumption A2 holds too, which completes the proof. \square

C.4.2 Proof of Theorem 11

Observe that for any $q \neq 0.5$,

$$(Q_P(1 - q) - Q_P(0.5)) \cdot (Q_P(0.5) - Q_P(q)) \geq 0.$$

Therefore, for any distribution P , we have that

$$\begin{aligned} \text{Dev}(P) &= \mathbb{E}_{q \sim \text{Unif}[0,1]} \left[[Q_P(1 - q) - Q_P(0.5) + Q_P(0.5) - Q_P(q)]^2 \right] \\ &\geq \mathbb{E}_{q \sim \text{Unif}[0,1]} \left[[Q_P(1 - q) - Q_P(0.5)]^2 + [Q_P(0.5) - Q_P(q)]^2 \right] \\ &= 2 \mathbb{E}_{q \sim \text{Unif}[0,1]} \left[(Q_P(q) - Q_P(0.5))^2 \right] \\ &= 2 \mathbb{E}_{X \sim P} \left[(X - \text{Median}(P))^2 \right] \geq 2 \text{Var}_{X \sim P}(X). \end{aligned}$$

Hence, it follows that

$$\mathbb{E} [\text{Dev}(P_{Y|Z})] \geq \sqrt{2} \mathbb{E} [\text{Var}(Y | Z)]$$

The result now follows from Theorem C.14. \square

C.4.3 Matching upper and lower bounds on conditional power

In (18), the asymptotic upper and lower bounds on conditional power for oracle matching match up to a factor of $\sqrt{2}$. In this section, we try to answer the natural question: can we close out this gap? We note that in Appendix C.2.1, this gap arises from (27), i.e.,

$$\widehat{\text{ISS}}_n \leq \|\Delta\mu^+(\mathbf{Y}, \mathbf{Z})\|_2 \leq \sqrt{2} \widehat{\text{ISS}}_n,$$

where $\|\Delta\mu^+(\mathbf{Y}, \mathbf{Z})\|_2$ is computed for the oracle matching.

Without additional assumptions on the model, it is unclear whether for oracle matching one can ensure $\|\Delta\mu^+(\mathbf{Y}, \mathbf{Z})\|_2 \approx \sqrt{2} \widehat{\text{ISS}}_n$. However, the proof of Lemma C.13 suggests that it is possible if oracle matching satisfies

$$Z_{i_\ell} \approx Z_{j_\ell} \text{ so that } \widehat{\mu}_{\text{ISO}}(Z_{i_\ell}) \approx \widehat{\mu}_{\text{ISO}}(Z_{j_\ell}), \text{ and } \mu(Y_{i_\ell}, Z_{i_\ell}) \approx -\mu(Y_{j_\ell}, Z_{j_\ell}).$$

However, we need to assume symmetry of $\mu(Y, Z)$ conditional on Z to have $\mu(Y_{i_\ell}, Z_{i_\ell}) \approx -\mu(Y_{j_\ell}, Z_{j_\ell})$. We investigate this in greater details for the special case of partial linear Gaussian model (20). With this additional structure, and the restriction $Z_{i_\ell} \approx Z_{j_\ell}$, we only need to ensure $Y_{i_\ell} \approx -Y_{j_\ell}$.

We show that under symmetry of the conditional distribution of $P_{Y|Z}$, even cross-bin matching can attain the upper bound on asymptotic conditional power of oracle matching in (18) and thus, oracle matching will attain it too.

Corollary C.15. *Under the setting of Theorem C.14, if $P_{Y|Z}$ is symmetric almost surely, then the conditional power of cross-bin matching satisfies*

$$\left| \mathbb{P} \{p \leq \alpha \mid \mathbf{Y}, \mathbf{Z}\} - \Phi \left(\sqrt{n} \beta_n \left\{ \frac{\mathbb{E} [\text{Var}(Y | Z)]}{\sigma^2} \right\}^{1/2} - \bar{\Phi}^{-1}(\alpha) \right) \right| = o_P(1).$$

Proof. If $P_{Y|Z=z}$ is symmetric, then

$$Q_{P_{Y|Z=z}}(1-q) - Q_{P_{Y|Z=z}}(q) = 2(Q_{P_{Y|Z=z}}(1-q) - Q_{P_{Y|Z=z}}(0.5)).$$

Consequently,

$$\text{Dev}(P_{Y|Z}) = 4 \mathbb{E}_{q \sim \text{Unif}[0,0.5]} \left[(Q_P(0.5) - Q_P(q))^2 \right] = 4 \mathbb{E}_{X \sim P_{Y|Z}} \left[(X - \text{Median}(P_{Y|Z}))^2 \right]$$

almost surely. Moreover, the median and mean coincide and hence, $\mathbb{E} [\text{Dev}(P_{Y|Z})] = 4 \text{Var}(Y | Z)$. The result now follows by Theorem C.14. \square

Since by Lemma 9, under the model class (20), $\sqrt{n}\beta_n (\mathbb{E} [\text{Var}(Y | Z)])^{1/2} (1 + o_P(1)) \leq \text{ISS}_n$, the conditional power for cross-bin matching further satisfies

$$\left| \mathbb{P} \{p \leq \alpha \mid \mathbf{Y}, \mathbf{Z}\} - \Phi \left(\frac{\widehat{\text{ISS}}_n}{\sigma} - \bar{\Phi}^{-1}(\alpha) \right) \right| = o_P(1),$$

as required.

C.5 Proof of propositions, corollaries and lemmas from Section 4

C.5.1 Proof of Proposition 6

Fix a distribution $Q_{X,Y,Z} \in H_0^{\text{ICI}}$. For any test function

$$\phi : (\mathcal{X} \times \mathbb{R} \times \mathbb{R})^n \rightarrow [0, 1] \text{ such that } \sup_{P \in H_0^{\text{ICI}}} \mathbb{E}_P [\phi(\mathbf{X}, \mathbf{Y}, \mathbf{Z})] \leq \alpha,$$

by definition of total-variation distance it holds that

$$\begin{aligned} \mathbb{E}_{P_{X,Y,Z}} [\phi(\mathbf{X}, \mathbf{Y}, \mathbf{Z})] &\leq \mathbb{E}_{Q_{X,Y,Z}} [\phi(\mathbf{X}, \mathbf{Y}, \mathbf{Z})] + d_{\text{TV}}(P_{X,Y,Z}^n, Q_{X,Y,Z}^n) \\ &\leq \alpha + d_{\text{TV}}(P_{X,Y,Z}^n, Q_{X,Y,Z}^n). \end{aligned}$$

The result now follows since the last inequality holds for any $Q_{X,Y,Z} \in H_0^{\text{ICI}}$.

C.5.2 Proof of Corollary 7

By Proposition 6, it is enough to argue that

$$\inf_{Q_{X,Y,Z} \in H_0^{\text{ICI}}} d_{\text{TV}}(P_{X,Y,Z}^n, Q_{X,Y,Z}^n) \leq \frac{\text{ISS}_n}{2\sigma}.$$

Consider any $\mu' \in \mathcal{C}_{\text{ISO}}$. Similar to (16), define a model $Q_{X,Y,Z}$ as

$$\mathbf{X} = \mu'(\mathbf{Z}) + \boldsymbol{\zeta}, \quad (Y_1, Z_1), \dots, (Y_n, Z_n) \stackrel{\text{iid}}{\sim} P_{Y,Z}, \quad \zeta_1, \dots, \zeta_n \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2).$$

By definition, $Q_{X,Y,Z} \in H_0^{\text{ICI}}$ and

$$\begin{aligned} d_{\text{TV}}(P_{X,Y,Z}, Q_{X,Y,Z}) &\leq \mathbb{E}_{P_{Y,Z}} [d_{\text{TV}}(\mathcal{N}(\mu(\mathbf{Y}, \mathbf{Z}), \sigma^2 I_n), \mathcal{N}(\mu'(\mathbf{Z}), \sigma^2 I_n) \mid \mathbf{Y}, \mathbf{Z})] \\ &= \frac{\mathbb{E}_{P_{Y,Z}} [\|\mu(\mathbf{Y}, \mathbf{Z}) - \mu'(\mathbf{Z})\|_2]}{2\sigma} \end{aligned}$$

Since this is true for any $\mu' \in \mathcal{C}_{\text{ISO}}$, then,

$$d_{\text{TV}}(P_{X,Y,Z}, Q_{X,Y,Z}) \leq \inf_{\mu' \in \mathcal{C}_{\text{ISO}}} \frac{\mathbb{E}_{P_{Y,Z}} [\|\mu(\mathbf{Y}, \mathbf{Z}) - \mu'(\mathbf{Z})\|_2]}{2\sigma} = \frac{\text{ISS}_n}{2\sigma}.$$

This proves the result. \square

C.5.3 Proof of Corollary 8

Similar to the proof of Corollary 7, it is enough to argue that

$$\inf_{Q_{X,Y,Z} \in H_0^{\text{ICI}}} d_{\text{TV}}(P_{X,Y,Z}^n, Q_{X,Y,Z}^n) \leq \left(\frac{1}{\epsilon(1-\epsilon)} \right)^{1/2} \text{ISS}_n.$$

Consider any $\mu' \in \mathcal{C}_{\text{ISO}}$, and define the model $Q_{X,Y,Z} \in H_0^{\text{ICI}}$ given by

$$X_i \sim \text{Ber}(\mu'(Z_i)), \quad (Y_1, Z_1), \dots, (Y_n, Z_n) \stackrel{\text{iid}}{\sim} P_{Y,Z}.$$

We note that $d_{\text{TV}}(P_{X,Y,Z}^n, Q_{X,Y,Z}^n)$ can be computed as

$$\mathbb{E}_{\mathbf{Y}, \mathbf{Z}} [d_{\text{TV}}(\text{Ber}(\mu(Y_1, Z_1)) \times \dots \times \text{Ber}(\mu(Y_n, Z_n)), \text{Ber}(\mu'(Z_1)) \times \dots \times \text{Ber}(\mu'(Z_n)))],$$

where the inner total variation term can be upper bounded using the Hellinger distance as

$$\begin{aligned} & d_{\text{TV}}(\text{Ber}(\mu(Y_1, Z_1)) \times \dots \times \text{Ber}(\mu(Y_n, Z_n)), \text{Ber}(\mu'(Z_1)) \times \dots \times \text{Ber}(\mu'(Z_n))) \\ & \leq \sqrt{2} \cdot H(\text{Ber}(\mu(Y_1, Z_1)) \times \dots \times \text{Ber}(\mu(Y_n, Z_n)), \text{Ber}(\mu'(Z_1)) \times \dots \times \text{Ber}(\mu'(Z_n))). \end{aligned}$$

Further, $H^2(P_1 \times \dots \times P_k, Q_1 \times \dots \times Q_k) \leq \sum_{i=1}^k H^2(P_i, Q_i)$ and thus by Lemma D.25, the inner total variation can be further upper bounded by

$$\left(\sum_{i=1}^n \frac{(\mu(Y_i, Z_i) - \mu'(Z_i))^2}{2\mu(Y_i, Z_i)(1 - \mu(Y_i, Z_i))} \right)^{1/2}.$$

Finally, since $\mu(Y, Z) \in (\epsilon, 1 - \epsilon)$ almost surely, the aforementioned term is upper bounded by $\frac{1}{\epsilon(1-\epsilon)} \cdot \|\mu(\mathbf{Y}, \mathbf{Z}) - \mu'(\mathbf{Z})\|_2^2$. Since this is true for any $\mu' \in \mathcal{C}_{\text{ISO}}$, we have

$$\begin{aligned} & d_{\text{TV}}(\text{Ber}(\mu(Y_1, Z_1)) \times \dots \times \text{Ber}(\mu(Y_n, Z_n)), \text{Ber}(\mu'(Z_1)) \times \dots \times \text{Ber}(\mu'(Z_n))) \\ & \leq \left(\inf_{\mu' \in \mathcal{C}_{\text{ISO}}} \frac{1}{\epsilon(1-\epsilon)} \cdot \|\mu(\mathbf{Y}, \mathbf{Z}) - \mu'(\mathbf{Z})\|_2^2 \right)^{1/2} = \left(\frac{1}{\epsilon(1-\epsilon)} \cdot \text{ISS}_n^2 \right)^{1/2}, \end{aligned}$$

which concludes the proof. \square

C.5.4 Proof of Lemma 9

Firstly, under the model class (20) with $\mathbb{E}[Y] = 0$, we have that

$$\begin{aligned} \text{ISS}_n & \stackrel{(1)}{\leq} \mathbb{E}_{P_{Y,Z}} [\|\mu(\mathbf{Y}, \mathbf{Z}) - \mu_0(\mathbf{Z})\|_2] \\ & = \mathbb{E}_{P_{Y,Z}} [\|\beta_n \mathbf{Y}\|_2] \\ & \stackrel{(2)}{\leq} \beta_n \left(\mathbb{E}_{P_{Y,Z}} \left[\sum_i Y_i^2 \right] \right)^{1/2} = \sqrt{n} \beta_n (\text{Var}(Y))^{1/2}, \end{aligned}$$

where the inequality (1) holds since $\mu_0 \in \mathcal{C}_{\text{ISO}}$ and (2) holds by Jensen's inequality. Next, we note that under the model class (20),

$$\begin{aligned} \text{ISS}_n &= \inf_{g \in \mathcal{C}_{\text{ISO}}} \mathbb{E}_{P_{Y,Z}} [\|\mu(\mathbf{Y}, \mathbf{Z}) - g(\mathbf{Z})\|_2] \\ &\geq \inf_{g: \mathbb{R} \rightarrow \mathbb{R}} \mathbb{E}_{P_{Y,Z}} [\|\mu(\mathbf{Y}, \mathbf{Z}) - g(\mathbf{Z})\|_2] \\ &= \mathbb{E}_{P_{Y,Z}} \left[\left\| \mu(\mathbf{Y}, \mathbf{Z}) - \mathbb{E}[\mu(\mathbf{Y}, \mathbf{Z}) \mid \mathbf{Z}] \right\|_2 \right] \\ &= \mathbb{E}_{P_{Y,Z}} \left[\left\| \beta_n (\mathbf{Y} - \mathbb{E}[\mathbf{Y} \mid \mathbf{Z}]) \right\|_2 \right]. \end{aligned}$$

Finally, since Y is a bounded random variable, we have

$$\frac{1}{n} \|\mathbf{Y} - \mathbb{E}[\mathbf{Y} \mid \mathbf{Z}]\|_2^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - E[Y \mid Z_i])^2 = \mathbb{E}[\text{Var}(Y \mid Z)] (1 + o_P(1)).$$

Thus, we have $\text{ISS}_n \geq \sqrt{n} \beta_n (\mathbb{E}[\text{Var}(Y \mid Z)])^{1/2} (1 + o_P(1))$. \square

C.6 An oracle matching—Isotonic median matching

A key step in Appendix C.2.1 is demonstrating (27) for a matching in $\mathcal{M}_n(\mathbf{Z})$. The purpose of this section is to show that in the case where \mathbb{R} is a totally ordered set, this is achievable via a method that we refer to as isotonic median matching (IMM). Since \mathbb{R} is totally ordered, we write \leq instead of \leq , and assume that $Z_1 \leq \dots \leq Z_n$. Given $a_1, \dots, a_m \in \mathbb{R}$, we define their median by

$$\text{Med}(a_1, \dots, a_m) = \begin{cases} a_{(k+1)} & \text{if } m = 2k + 1 \text{ for some } k \in \mathbb{N}_0, \\ \frac{a_{(k)} + a_{(k+1)}}{2} & \text{if } m = 2k \text{ for some } k \in \mathbb{N}, \end{cases}$$

where $a_{(1)} \leq \dots \leq a_{(m)}$ denote the order statistics of a_1, \dots, a_m . Now given (\mathbf{Y}, \mathbf{Z}) and any function $\gamma : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$, we consider the function $\tilde{\gamma}_{\text{ISO}} : \mathbb{R} \rightarrow \mathbb{R}$ where

$$\tilde{\gamma}_{\text{ISO}}(z) := \max_{j \in [n]: Z_j \leq z} \min_{z_r \in \mathbb{R}: z_r \geq z} \text{Med}(\{\gamma(Y_\ell, Z_\ell) : Z_j \leq Z_\ell \leq z_r\}). \quad (40)$$

We claim that $\tilde{\gamma}_{\text{ISO}} \in \mathcal{C}_{\text{ISO}}$. To see this, observe that for $z \leq z'$,

$$\begin{aligned} \tilde{\gamma}_{\text{ISO}}(z) &= \max_{j \in [n]: Z_j \leq z} \min_{z_r \in \mathbb{R}: z_r \geq z} \text{Med}(\{\gamma(Y_\ell, Z_\ell) : Z_j \leq Z_\ell \leq z_r\}) \\ &\leq \max_{j \in [n]: Z_j \leq z} \min_{z_r \in \mathbb{R}: z_r \geq z'} \text{Med}(\{\gamma(Y_\ell, Z_\ell) : Z_j \leq Z_\ell \leq z_r\}) \\ &\leq \max_{j \in [n]: Z_j \leq z'} \min_{z_r \in \mathbb{R}: z_r \geq z'} \text{Med}(\{\gamma(Y_\ell, Z_\ell) : Z_j \leq Z_\ell \leq z_r\}) = \tilde{\gamma}_{\text{ISO}}(z') \end{aligned}$$

which proves the claim. Further, we note that $\tilde{\gamma}_{\text{ISO}}$ is a piecewise constant function in \mathcal{C}_{ISO} , i.e., there exist $m \in [n + 1]$, integers $1 = n_1 < n_2 < n_3 < \dots < n_m = n + 1$ and real numbers $r_1 < \dots < r_{m-1}$ such that

$$\tilde{\gamma}_{\text{ISO}}(z) = \sum_{i=1}^{m-1} r_i \mathbb{1} \{Z_{n_i} \leq z < Z_{n_{i+1}}\}, \quad (41)$$

with the convention that $Z_{n_{m+1}} = \infty$. In fact, we can say more:

Lemma C.16. *In the representation (41), we have $r_i = \text{Med}(\{\gamma(Y_\ell, Z_\ell)\}_{n_i \leq \ell \leq n_{i+1}-1})$ for $i \in [m-1]$. Further, for any $i \in [m-1]$ and any integer $k \in \{n_i, n_i+1, \dots, n_{i+1}-1\}$,*

$$\text{Med}(\{\gamma(Y_\ell, Z_\ell)\}_{n_i \leq \ell \leq k}) \geq r_i. \quad (42)$$

Proof. Fix $i \in [m-1]$ and let $z_0 \in \mathbb{R}$ be such that $Z_{n_{i-1}} < z_0 < Z_{n_i}$. For any $j \in [n]$ with $Z_j \leq z_0$, we have that

$$\min_{z_r \in \mathbb{R}: z_r \geq z_0} \text{Med}(\{\gamma(Y_\ell, Z_\ell) : Z_j \leq Z_\ell \leq z_r\}) = \min_{z_r \in \mathbb{R}: z_r \geq Z_{n_i}} \text{Med}(\{\gamma(Y_\ell, Z_\ell) : Z_j \leq Z_\ell \leq z_r\}).$$

Since $\tilde{\gamma}_{\text{ISO}}(z_0) = r_{i-1} < r_i = \tilde{\gamma}_{\text{ISO}}(Z_{n_i})$, it follows that

$$\begin{aligned} \tilde{\gamma}_{\text{ISO}}(Z_{n_i}) &= \max \left\{ \min_{z_r \in \mathbb{R}: z_r \geq Z_{n_i}} \text{Med}(\{\gamma(Y_\ell, Z_\ell) : Z_{n_i} \leq Z_\ell \leq z_r\}), \tilde{\gamma}_{\text{ISO}}(z_0) \right\} \\ &= \min_{k \in [n]: k \geq n_i} \text{Med}(\{\gamma(Y_\ell, Z_\ell) : n_i \leq \ell \leq k\}). \end{aligned} \quad (43)$$

Let $\hat{k} \in [n]$ denote the largest index k at which the minimum in (43) is attained. Since $\tilde{\gamma}_{\text{ISO}}(Z_{n_i}) = \tilde{\gamma}_{\text{ISO}}(Z_{n_{i+1}}) = \dots = \tilde{\gamma}_{\text{ISO}}(Z_{n_{i+1}-1})$, we have $\hat{k} \geq n_{i+1} - 1$. Now, it suffices to show that

$$\text{Med}(\{\gamma(Y_\ell, Z_\ell) : n_i \leq \ell \leq \hat{k}\}) = \text{Med}(\{\gamma(Y_\ell, Z_\ell) : n_i \leq \ell \leq n_{i+1} - 1\}). \quad (44)$$

If $i = m-1$, then $\hat{k} = n_m - 1 = n$, so we may assume that $i \in [m-2]$. Suppose for a contradiction that (44) is not true, so $\hat{k} \geq n_{i+1}$. Now by (43), we have that

$$\begin{aligned} \tilde{\gamma}_{\text{ISO}}(Z_{n_{i+1}}) &= \min_{k \in [n]: k \geq n_{i+1}} \text{Med}(\{\gamma(Y_\ell, Z_\ell) : n_{i+1} \leq \ell \leq k\}) \\ &\leq \text{Med}(\{\gamma(Y_\ell, Z_\ell) : n_{i+1} \leq \ell \leq \hat{k}\}) \end{aligned}$$

But since we have assumed that (44) is not true, we have

$$\text{Med}(\{\gamma(Y_\ell, Z_\ell) : n_i \leq \ell \leq \hat{k}\}) < \text{Med}(\{\gamma(Y_\ell, Z_\ell) : n_i \leq \ell \leq n_{i+1} - 1\}),$$

and thus by Lemma D.26, we further have

$$\tilde{\gamma}_{\text{ISO}}(Z_{n_{i+1}}) \leq \text{Med}(\{\gamma(Y_\ell, Z_\ell) : n_i \leq \ell \leq \hat{k}\}).$$

However by (41), we also have

$$\tilde{\gamma}_{\text{ISO}}(Z_{n_{i+1}}) > \tilde{\gamma}_{\text{ISO}}(Z_{n_i}) = \text{Med}(\{\gamma(Y_\ell, Z_\ell) : n_i \leq \ell \leq \hat{k}\}),$$

which is a contradiction. This establishes (44) and completes the proof of the first part of the result.

For the final part, fix any $i \in [m-1]$ and $k_0 \in \{n_i, n_i+1, \dots, n_{i+1}-1\}$. Then by (43),

$$\begin{aligned} r_i = \tilde{\gamma}_{\text{ISO}}(Z_{n_i}) &= \min_{k \in [n]: k \geq n_i} \text{Med}(\{\gamma(Y_\ell, Z_\ell) : n_i \leq \ell \leq k\}) \\ &\leq \text{Med}(\{\gamma(Y_\ell, Z_\ell) : n_i \leq \ell \leq k_0\}), \end{aligned}$$

as required. \square

We are now in a position to define the isotonic median matching. For $i \in [m - 1]$, define

$$\begin{aligned} P_i &:= \{t \in \{n_i, n_i + 1, \dots, n_{i+1} - 1\} : \gamma(Y_t, Z_t) > r_i\}, \\ N_i &:= \{t \in \{n_i, n_i + 1, \dots, n_{i+1} - 1\} : \gamma(Y_t, Z_t) < r_i\}. \end{aligned}$$

Let us order the indices in P_i as $t_{i,1}^+ < \dots < t_{i,n_i^+}^+$ and the indices in N_i as $t_{i,1}^- < \dots < t_{i,n_i^-}^-$. Defining $L_i := n_i^+ \wedge n_i^-$, we consider the collection of ordered pairs

$$\mathcal{C}_i := \{(t_{i,1}^+, t_{i,1}^-), \dots, (t_{i,L_i}^+, t_{i,L_i}^-)\}.$$

Finally, $\widetilde{M}(\gamma) := \bigcup_{i=1}^{m-1} \mathcal{C}_i$ defines the isotonic median matching.

Lemma C.17. *Suppose that the elements of $\{\gamma(Y_\ell, Z_\ell) : \ell \in [n]\}$ are all distinct. Then for any $i \in [m - 1]$, we have $n_i^+ = n_i^-$ and $\gamma(Y_\ell, Z_\ell) = \widetilde{\gamma}_{\text{ISO}}(Z_\ell)$ for all $\ell \in \{n_i, n_i + 1, \dots, n_{i+1} - 1\} \setminus (P_i \cup N_i)$. Moreover, $\widetilde{M}(\gamma) \in \mathcal{M}_n(\mathbf{Z})$.*

Proof. Fix $i \in [m - 1]$. When $n_{i+1} - n_i$ is even, we have by the first part of Lemma C.16 that $n_i^+ = n_i^- = (n_{i+1} - n_i)/2$ and $\{n_i, n_i + 1, \dots, n_{i+1} - 1\} \setminus (P_i \cup N_i) = \emptyset$. On the other hand, when $n_{i+1} - n_i$ is odd, we have $n_i^+ = n_i^- = (n_{i+1} - n_i - 1)/2$ and $\{n_i, n_i + 1, \dots, n_{i+1} - 1\} \setminus (P_i \cup N_i) =: \{\ell_0\}$ is a singleton set. Moreover, $\gamma(Y_{\ell_0}, Z_{\ell_0}) = \widetilde{\gamma}_{\text{ISO}}(Z_{\ell_0})$ which proves the first two claims.

Since $Z_1 \leq \dots \leq Z_n$, in order to prove that $\widetilde{M}(\gamma)$ is a valid matching, it suffices to show that $t_{i,\ell}^+ < t_{i,\ell}^-$ for any $i \in [m - 1]$ and $\ell \in [L_i]$. To see this, noting that $t_{i,\ell}^+ \neq t_{i,\ell}^-$, suppose for a contradiction that $t_{i,\ell_0}^+ > t_{i,\ell_0}^-$ for some $i \in [m - 1]$ and some minimal $\ell_0 \in [L_i]$. Since the elements of $\{\gamma(Y_\ell, Z_\ell) : \ell \in [n]\}$ are all distinct, we have $|\{t \in P_i : t \leq t_{i,\ell_0}^-\}| = \ell_0 - 1$ and $|\{t \in N_i : t \leq t_{i,\ell_0}^-\}| = \ell_0$, so

$$\text{Med}(\{\gamma(Y_\ell, Z_\ell)\}_{n_i \leq \ell \leq t_{i,\ell_0}^-}) < r_i,$$

which contradicts (42). □

The following key property of IMM ensures that when the `PairSwap-ICI` test is run with the oracle matching and oracle choice of weights from (13), it has valid Type I error control by Theorem 1 and moreover satisfies the power guarantees of Corollary ??.

Theorem C.18. *Given $(\mathbf{Y}, \mathbf{Z}) \in \mathbb{R}^n \times \mathbb{R}^n$, there exists $\gamma : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ for which, with the corresponding isotonic median matching $\widetilde{M}(\gamma) \in \mathcal{M}_n(\mathbf{Z})$, we have*

$$\|\Delta\mu^+(\mathbf{Y}, \mathbf{Z})\|_2 \geq \widehat{\text{ISS}}_n.$$

Proof. Fix $\epsilon > 0$. We can find $\gamma : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ such that the coordinates of $\gamma(\mathbf{Y}, \mathbf{Z})$ are all distinct and $\|\mu(\mathbf{Y}, \mathbf{Z}) - \gamma(\mathbf{Y}, \mathbf{Z})\|_2 < \epsilon$. Consider the IMM $\widetilde{M}(\gamma) := \{(i_\ell, j_\ell)\}_{\ell \in \bar{L}}$, where $\widetilde{M}(\gamma) \in \mathcal{M}_n(\mathbf{Z})$ by Lemma C.17. Furthermore by Lemma C.17 and (41), we have for any

$\ell \in [\tilde{L}]$ that $\gamma(Y_{i_\ell}, Z_{i_\ell}) > \tilde{\gamma}_{\text{ISO}}(Z_{i_\ell}) = \tilde{\gamma}_{\text{ISO}}(Z_{j_\ell}) > \gamma(Y_{j_\ell}, Z_{j_\ell})$. Hence,

$$\begin{aligned} \|\Delta\gamma^+(\mathbf{Y}, \mathbf{Z})\|_2^2 &= \sum_{\ell=1}^{\tilde{L}} (\gamma(Y_{i_\ell}, Z_{i_\ell}) - \gamma(Y_{j_\ell}, Z_{j_\ell}))^2 \\ &= \sum_{\ell=1}^{\tilde{L}} (\gamma(Y_{i_\ell}, Z_{i_\ell}) - \tilde{\gamma}_{\text{ISO}}(Z_{i_\ell}) + \tilde{\gamma}_{\text{ISO}}(Z_{j_\ell}) - \gamma(Y_{j_\ell}, Z_{j_\ell}))^2 \\ &\geq \sum_{\ell=1}^{\tilde{L}} \{(\gamma(Y_{i_\ell}, Z_{i_\ell}) - \tilde{\gamma}_{\text{ISO}}(Z_{i_\ell}))^2 + (\tilde{\gamma}_{\text{ISO}}(Z_{j_\ell}) - \gamma(Y_{j_\ell}, Z_{j_\ell}))^2\} \\ &= \sum_{i=1}^n (\gamma(Y_i, Z_i) - \tilde{\gamma}_{\text{ISO}}(Z_i))^2 = \|\gamma(\mathbf{Y}, \mathbf{Z}) - \tilde{\gamma}_{\text{ISO}}(\mathbf{Z})\|_2^2, \end{aligned}$$

where the penultimate equality holds because $\gamma(Y_i, Z_i) = \tilde{\gamma}_{\text{ISO}}(Z_i)$ for $i \in [n] \setminus \{i_1, j_1, \dots, i_{\tilde{L}}, j_{\tilde{L}}\}$, by Lemma C.17. Next, we observe that

$$\begin{aligned} \|\Delta\gamma^+(\mathbf{Y}, \mathbf{Z})\|_2 - \|\Delta\mu^+(\mathbf{Y}, \mathbf{Z})\|_2 &\leq \|\Delta\gamma^+(\mathbf{Y}, \mathbf{Z}) - \Delta\mu^+(\mathbf{Y}, \mathbf{Z})\|_2 \\ &\leq \|\Delta\gamma(\mathbf{Y}, \mathbf{Z}) - \Delta\mu(\mathbf{Y}, \mathbf{Z})\|_2 \\ &\leq \sqrt{2}\|\gamma(\mathbf{Y}, \mathbf{Z}) - \mu(\mathbf{Y}, \mathbf{Z})\|_2 \leq \sqrt{2}\epsilon, \end{aligned}$$

and that by (17),

$$\widehat{\text{ISS}}_n \leq \|\mu(\mathbf{Y}, \mathbf{Z}) - \tilde{\gamma}_{\text{ISO}}(\mathbf{Z})\|_2 \leq \|\gamma(\mathbf{Y}, \mathbf{Z}) - \tilde{\gamma}_{\text{ISO}}(\mathbf{Z})\|_2 + \epsilon.$$

Therefore,

$$\|\Delta\mu^+(\mathbf{Y}, \mathbf{Z})\|_2 + \sqrt{2}\epsilon \geq \|\Delta\gamma^+(\mathbf{Y}, \mathbf{Z})\|_2 \geq \|\gamma(\mathbf{Y}, \mathbf{Z}) - \tilde{\gamma}_{\text{ISO}}(\mathbf{Z})\|_2 \geq \widehat{\text{ISS}}_n - \epsilon.$$

Since $\epsilon > 0$ was arbitrary, the result follows. \square

C.7 Proof of lemmas from Appendices C.1—C.4

Lemma C.19. *In the setting of Theorem C.12, let*

$$\begin{aligned} \Omega_1 &:= \{\|\mathbf{w} \circ \Delta\mathbf{X}\|_2^2 \leq \rho_2^2 \|\mathbf{w}\|_2^2 \cdot (1 + \epsilon_{1,\delta,U})\}, \\ \Omega_2 &:= \{\|\mathbf{w} \circ \Delta\mathbf{X}\|_2^2 \geq \rho_2^2 \|\mathbf{w}\|_2^2 \cdot \max\{1 + \epsilon_{1,\delta,L}, 0\}\}, \\ \Omega_3 &:= \{0.56\|\mathbf{w} \circ \Delta\mathbf{X}\|_3^3 \leq \|\mathbf{w} \circ \Delta\mathbf{X}\|_2^3 \cdot \epsilon_{2,\delta}\}, \end{aligned}$$

where $\epsilon_{1,\delta,U}, \epsilon_{1,\delta,L}, \epsilon_{2,\delta}$ are as defined in (24). Then for any $\delta \in [0, 1]$, we have $\mathbb{P}\{\Omega_1 \mid \mathbf{Y}, \mathbf{Z}\} \geq 1 - 2\delta$, $\mathbb{P}\{\Omega_2 \mid \mathbf{Y}, \mathbf{Z}\} \geq 1 - 2\delta$ and $\mathbb{P}\{\Omega_3 \mid \mathbf{Y}, \mathbf{Z}\} \geq 1 - 2\delta$.

Proof. We have

$$\|\mathbf{w} \circ \Delta\mathbf{X}\|_2^2 = \|\mathbf{w} \circ \Delta\mu(\mathbf{Y}, \mathbf{Z})\|_2^2 + \|\mathbf{w} \circ \Delta\zeta\|_2^2 + 2 \cdot (\mathbf{w} \circ \Delta\mu(\mathbf{Y}, \mathbf{Z}))^T (\mathbf{w} \circ \Delta\zeta).$$

Now, conditional on (\mathbf{Y}, \mathbf{Z}) , both $\|\mathbf{w} \circ \Delta\zeta\|_2^2$, and $(\mathbf{w} \circ \Delta\mu(\mathbf{Y}, \mathbf{Z}))^T(\mathbf{w} \circ \Delta\zeta)$ are weighted sums of independent and identically distributed random variables. Hence,

$$\begin{aligned}\mathbb{E} [\|\mathbf{w} \circ \Delta\zeta\|_2^2 \mid \mathbf{Y}, \mathbf{Z}] &= \rho_2^2 \sum_{\ell=1}^L w_\ell^2 = \rho_2^2 \|\mathbf{w}\|_2^2, \\ \mathbb{E} \left[(\mathbf{w} \circ \Delta\mu(\mathbf{Y}, \mathbf{Z}))^T (\mathbf{w} \circ \Delta\zeta) \mid \mathbf{Y}, \mathbf{Z} \right] &= 0 \\ \text{Var} (\|\mathbf{w} \circ \Delta\zeta\|_2^2 \mid \mathbf{Y}, \mathbf{Z}) &= \sum_{\ell=1}^L w_\ell^4 \text{Var} ((\Delta_\ell \zeta)^2) \\ &\leq \sum_{\ell=1}^L w_\ell^4 \mathbb{E} [(\Delta_\ell \zeta)^4] = \rho_4^4 \|\mathbf{w}\|_4^4, \\ \text{Var} \left((\mathbf{w} \circ \Delta\mu(\mathbf{Y}, \mathbf{Z}))^T (\mathbf{w} \circ \Delta\zeta) \mid \mathbf{Y}, \mathbf{Z} \right) &= \sum_{\ell=1}^L w_\ell^4 (\Delta_\ell \mu(\mathbf{Y}, \mathbf{Z}))^2 \mathbb{E} [(\Delta_\ell \zeta)^2] \\ &= \rho_2^2 \|\mathbf{w}^2 \circ \Delta\mu(\mathbf{Y}, \mathbf{Z})\|_2^2.\end{aligned}$$

By Chebychev's inequality, for any $\delta \in [0, \infty)$,

$$\begin{aligned}\mathbb{P} \left\{ \left| \|\mathbf{w} \circ \Delta\zeta\|_2^2 - \rho_2^2 \|\mathbf{w}\|_2^2 \right| \geq \frac{\rho_4^2}{\delta^{1/2}} \|\mathbf{w}\|_4^2 \mid \mathbf{Y}, \mathbf{Z} \right\} &\leq \delta, \\ \mathbb{P} \left\{ \left| (\mathbf{w} \circ \Delta\mu(\mathbf{Y}, \mathbf{Z}))^T (\mathbf{w} \circ \Delta\zeta) \right| \geq \frac{\rho_2}{\delta^{1/2}} \|\mathbf{w}^2 \circ \Delta\mu(\mathbf{Y}, \mathbf{Z})\|_2 \mid \mathbf{Y}, \mathbf{Z} \right\} &\leq \delta.\end{aligned}$$

Hence,

$$\begin{aligned}1 - 2\delta &\leq \mathbb{P} \left\{ \|\mathbf{w} \circ \Delta\mathbf{X}\|_2^2 \leq \|\mathbf{w} \circ \Delta\mu(\mathbf{Y}, \mathbf{Z})\|_2^2 + \rho_2^2 \|\mathbf{w}\|_2^2 + \frac{\rho_4^2}{\delta^{1/2}} \|\mathbf{w}\|_4^2 \right. \\ &\quad \left. + \frac{2\rho_2}{\delta^{1/2}} \|\mathbf{w}^2 \circ \Delta\mu(\mathbf{Y}, \mathbf{Z})\|_2 \mid \mathbf{Y}, \mathbf{Z} \right\} \\ &\leq \mathbb{P} \{ \Omega_1 \mid \mathbf{Y}, \mathbf{Z} \},\end{aligned}$$

where, in the final step, we have used the facts that $\|\mathbf{w}\|_4^2 \leq \|\mathbf{w}\|_\infty \cdot \|\mathbf{w}\|_2$ and $\|\mathbf{w}^2 \circ \Delta\mu(\mathbf{Y}, \mathbf{Z})\|_2 \leq \|\mathbf{w} \circ \Delta\mu(\mathbf{Y}, \mathbf{Z})\|_\infty \cdot \|\mathbf{w}\|_2$.

The lower bound on $\mathbb{P} \{ \Omega_2 \mid \mathbf{Y}, \mathbf{Z} \}$ follows by a very similar argument.

For the final bound, observe that $\|\mathbf{w} \circ \Delta\zeta\|_3^3$ is also a weighted sum of independent and identically distributed random variables, and has (conditional) mean and variance given by

$$\begin{aligned}\mathbb{E} [\|\mathbf{w} \circ \Delta\zeta\|_3^3 \mid \mathbf{Y}, \mathbf{Z}] &= \sum_{\ell=1}^L |w_\ell|^3 \rho_3^3 = \rho_3^3 \|\mathbf{w}\|_3^3, \\ \text{Var} (\|\mathbf{w} \circ \Delta\zeta\|_3^3 \mid \mathbf{Y}, \mathbf{Z}) &= \sum_{\ell=1}^L w_\ell^6 \text{Var} (|\Delta_\ell \zeta|^3) \leq \sum_{\ell=1}^L w_\ell^6 \mathbb{E} [(\Delta_\ell \zeta)^6] = \|\Delta\mathbf{w}\|_6^6 \rho_6^6.\end{aligned}$$

Thus, by the triangle inequality and Chebychev's inequality,

$$\begin{aligned}
1 - \delta &\leq \mathbb{P} \left\{ \|\mathbf{w} \circ \Delta \mathbf{X}\|_3 \leq \|\mathbf{w} \circ \Delta \mu(\mathbf{Y}, \mathbf{Z})\|_3 + \left(\rho_3^3 \|\mathbf{w}\|_3^3 + \frac{\rho_6^3}{\delta^{1/2}} \cdot \|\mathbf{w}\|_6^3 \right)^{1/3} \mid \mathbf{Y}, \mathbf{Z} \right\} \\
&\leq \mathbb{P} \left\{ \|\mathbf{w} \circ \Delta \mathbf{X}\|_3 \leq \rho_2 \|\mathbf{w}\|_2 \left\{ \left(\frac{\|\mathbf{w} \circ \Delta \mu(\mathbf{Y}, \mathbf{Z})\|_2}{\rho_2 \|\mathbf{w}\|_2} \right)^{2/3} \left(\frac{\|\mathbf{w} \circ \Delta \mu(\mathbf{Y}, \mathbf{Z})\|_\infty}{\rho_2 \|\mathbf{w}\|_2} \right)^{1/3} \right. \right. \\
&\quad \left. \left. + \left(\frac{\rho_3^3 \|\mathbf{w}\|_\infty}{\rho_2^3 \|\mathbf{w}\|_2} + \frac{\rho_6^3}{\rho_2^3 \delta^{1/2}} \cdot \frac{\|\mathbf{w}\|_\infty^2}{\|\mathbf{w}\|_2^2} \right)^{1/3} \right\} \mid \mathbf{Y}, \mathbf{Z} \right\} \\
&\leq \mathbb{P} \{ \Omega_2 \cap \Omega_3 \mid \mathbf{Y}, \mathbf{Z} \} + \mathbb{P} \{ \Omega_2^c \mid \mathbf{Y}, \mathbf{Z} \} \leq \mathbb{P} \{ \Omega_3 \mid \mathbf{Y}, \mathbf{Z} \} + 2\delta,
\end{aligned}$$

as required. \square

Lemma C.20. *Suppose we have i.i.d. samples $(Y_1, Z_1), \dots, (Y_n, Z_n)$ from some distribution $P_{Y,Z} \in [-1, 1] \times \mathbb{R}$, and π is any permutation of $[n]$ such that $Z_{\pi(1)} \leq Z_{\pi(2)} \leq \dots \leq Z_{\pi(n)}$ are ordered. Then,*

$$\left| \frac{1}{n/2} \sum_{i=1}^{\lfloor n/2 \rfloor} (Y_{\pi(2i-1)} - Y_{\pi(2i)})_+^2 - \mathbb{E}[\text{Var}(Y \mid Z)] \right| = o_P(1).$$

Proof. Consider any joint distribution $P_{Y,Z}$ on $[-1, 1] \times \mathbb{R}$. We write P_Z and $P_{Y|Z}(\cdot \mid z)$ for the distributions of Z and Y given $Z = z$, and we write F_Z^{-1} and $F_{Y|Z}^{-1}(\cdot \mid z)$ to denote the generalized inverse for P_Z and $P_{Y|Z}(\cdot \mid z)$ respectively. Without loss of generality we can also assume that π is the identity permutation, i.e., $Z_1 \leq Z_2 \leq \dots \leq Z_n$ are ordered and Y_i is the Y value corresponding to Z_i for any $i \in [n]$.

Step 1: concentration, conditional on \mathbf{Z} . Conditioned on $\mathbf{Z} = (Z_1, Z_2, \dots, Z_n)$, $\{(Y_{2i-1} - Y_{2i})\}_{i \in [n/2]}$ is a collection of independent and uniformly bounded random variables. Thus, by weak law of large numbers,

$$\left| \frac{1}{n/2} \sum_{i=1}^{\lfloor n/2 \rfloor} (Y_{2i-1} - Y_{2i})_+^2 - \frac{1}{n/2} \sum_{i=1}^{\lfloor n/2 \rfloor} \mathbb{E}[(Y_{2i-1} - Y_{2i})_+^2 \mid \mathbf{Z}] \right| = o_P(1).$$

Hence, it suffices to show that

$$\left| \frac{1}{n/2} \sum_{i=1}^{\lfloor n/2 \rfloor} \mathbb{E}[(Y_{2i-1} - Y_{2i})_+^2 \mid \mathbf{Z}] - \mathbb{E}[\text{Var}(Y \mid Z)] \right| = o_P(1),$$

Step 2: constructing a coupling. Consider $\{(U_{2i-1}, U_{2i})\}_{i \in [n/2]} \stackrel{\text{iid}}{\sim} \text{Uniform}[0, 1]$. Without loss of generality, we may take

$$Y_{2i-1} = F_{Y|Z}^{-1}(U_{2i-1} \mid Z_{2i-1}), \quad Y_{2i} = F_{Y|Z}^{-1}(U_{2i} \mid Z_{2i}).$$

We also define for $i \in [n/2]$

$$Y'_{2i-1} = F_{Y|Z}^{-1}(U_{2i-1} \mid Z_{2i}), \quad Y'_{2i} = F_{Y|Z}^{-1}(U_{2i} \mid Z_{2i-1}).$$

Conditioned on \mathbf{Z} , $Y'_{2i-1} \stackrel{d}{=} Y_{2i}$ and $Y'_{2i} \stackrel{d}{=} Y_{2i-1}$. Since $t \rightarrow t_+^2$ is 4-Lipschitz on $[-2, 2]$, it follows that

$$\begin{aligned} & \left| 2(Y_{2i-1} - Y_{2i})_+^2 - [(Y'_{2i-1} - Y_{2i})_+^2 + (Y_{2i-1} - Y'_{2i})_+^2] \right| \\ & \leq |(Y_{2i-1} - Y_{2i})_+^2 - (Y'_{2i-1} - Y_{2i})_+^2| + |(Y_{2i-1} - Y_{2i})_+^2 - (Y_{2i-1} - Y'_{2i})_+^2| \\ & \leq 4(|Y_{2i-1} - Y'_{2i-1}| + |Y_{2i} - Y'_{2i}|). \end{aligned}$$

Therefore, we have

$$\begin{aligned} & \left| \frac{1}{n/2} \sum_{i=1}^{\lfloor n/2 \rfloor} \mathbb{E} [(Y_{2i-1} - Y_{2i})_+^2 \mid \mathbf{Z}] - \frac{1}{n} \sum_{i=1}^{\lfloor n/2 \rfloor} \mathbb{E} [(Y'_{2i-1} - Y_{2i})_+^2 + (Y_{2i-1} - Y'_{2i})_+^2 \mid \mathbf{Z}] \right| \\ & \leq \frac{4}{n} \sum_{i=1}^n \mathbb{E} [|Y_{2i-1} - Y'_{2i-1}| + |Y_{2i} - Y'_{2i}| \mid \mathbf{Z}]. \end{aligned}$$

Moreover, by construction

$$\mathbb{E} [|Y_{2i-1} - Y'_{2i-1}| \mid \mathbf{Z}] = \mathbb{E} [|Y_{2i} - Y'_{2i}| \mid \mathbf{Z}] = d_{W_1}(P_{Y|Z}(\cdot \mid Z_{2i-1}), P_{Y|Z}(\cdot \mid Z_{2i})),$$

where $d_{W_1}(\cdot, \cdot)$ denotes the 1-Wasserstein distance. We also write $\sigma_z^2 := \text{Var}(Y \mid Z = z)$ and note that

$$\mathbb{E} [(Y_{2i-1} - Y'_{2i-1})_+^2 \mid \mathbf{Z}] = \sigma_{Z_{2i-1}}^2 \quad \text{and} \quad \mathbb{E} [(Y_{2i} - Y'_{2i-1})_+^2 \mid \mathbf{Z}] = \sigma_{Z_{2i}}^2.$$

Hence, it follows that

$$\begin{aligned} & \left| \frac{1}{n/2} \sum_{i=1}^{\lfloor n/2 \rfloor} \mathbb{E} [(Y_{2i-1} - Y_{2i})_+^2 \mid \mathbf{Z}] - \frac{1}{n} \sum_{i=1}^{\lfloor n/2 \rfloor} (\sigma_{Z_{2i-1}}^2 + \sigma_{Z_{2i}}^2) \right| \\ & \leq \frac{8}{n} \sum_{i=1}^{\lfloor n/2 \rfloor} d_{W_1}(P_{Y|Z}(\cdot \mid Z_{2i-1}), P_{Y|Z}(\cdot \mid Z_{2i})). \end{aligned}$$

Step 3: another concentration step. Since $Y \in [-1, 1]$,

$$\left| \frac{1}{n} \sum_{i=1}^{\lfloor n/2 \rfloor} (\sigma_{Z_{2i-1}}^2 + \sigma_{Z_{2i}}^2) - \frac{1}{n} \sum_{i=1}^n \sigma_{Z_i}^2 \right| \leq \frac{4}{n}.$$

Furthermore, $\sigma^2 Z_1, \dots, \sigma^2 Z_n$ are independent and uniformly bounded random variables. Thus, by weak law of large numbers,

$$\left| \frac{1}{n} \sum_{i=1}^n \sigma_{Z_i}^2 - \mathbb{E} [\text{Var}(Y \mid Z)] \right| = o_P(1).$$

Hence, so far we have established that

$$\left| \frac{1}{n/2} \sum_{i=1}^{\lfloor n/2 \rfloor} \mathbb{E} [(Y_{2i-1} - Y_{2i})_+^2 \mid \mathbf{Z}] - \mathbb{E} [\text{Var}(Y \mid Z)] \right| \leq \frac{8}{n} \sum_{i=1}^{\lfloor n/2 \rfloor} d_{W_1}(P_{Y|Z}(\cdot \mid Z_{2i-1}), P_{Y|Z}(\cdot \mid Z_{2i})) + o_P(1).$$

Step 4: bounding the Wasserstein distance by a TV distance. Fix $N \in \mathbb{N}$, and define the random variable $\tilde{Y} := (1/N)\lfloor NY \rfloor$, which essentially amounts to rounding off Y to the closest point to the left of Y on the grid $(\mathbb{N})/N$.

For any $i \in \lfloor n/2 \rfloor$,

$$\begin{aligned} & \left| d_{W_1}(P_{Y|Z}(\cdot | Z_{2i-1}), P_{Y|Z}(\cdot | Z_{2i})) - d_{W_1}(P_{\tilde{Y}|Z}(\cdot | Z_{2i-1}), P_{\tilde{Y}|Z}(\cdot | Z_{2i})) \right| \\ & \leq d_{W_1}(P_{Y|Z}(\cdot | Z_{2i-1}), P_{\tilde{Y}|Z}(\cdot | Z_{2i-1})) + d_{W_1}(P_{Y|Z}(\cdot | Z_{2i}), P_{\tilde{Y}|Z}(\cdot | Z_{2i})) \leq 2/N, \end{aligned}$$

where the last inequality follows by noting that $|Y - \tilde{Y}| \leq 1/N$. Furthermore, $\tilde{Y} \in [-1, 1]$ and thus, by Villani (2009, Theorem 6.15)

$$\begin{aligned} d_{W_1}(P_{\tilde{Y}|Z}(\cdot | Z_{2i-1}), P_{\tilde{Y}|Z}(\cdot | Z_{2i})) & \leq 2 \, d_{TV}(P_{\tilde{Y}|Z}(\cdot | Z_{2i-1}), P_{\tilde{Y}|Z}(\cdot | Z_{2i})) \\ & \leq 2 \sum_{k=-N}^N |P_{\tilde{Y}|Z}(k/N | Z_{2i-1}) - P_{\tilde{Y}|Z}(k/N | Z_{2i})|. \end{aligned}$$

Writing $f_k(Z_1, \dots, Z_n) := \frac{2}{n-1} \sum_{i=1}^{n-1} |P_{\tilde{Y}|Z}(k/N | Z_i) - P_{\tilde{Y}|Z}(k/N | Z_{i+1})|$ for any $k \in [-N, N]$, we have that

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^{\lfloor n/2 \rfloor} d_{W_1}(P_{Y|Z}(\cdot | Z_{2i-1}), P_{Y|Z}(\cdot | Z_{2i})) \\ & \leq \frac{1}{N} + \frac{2}{n} \sum_{i=1}^{\lfloor n/2 \rfloor} \sum_{k=-N}^N |P_{\tilde{Y}|Z}(k/N | Z_{2i-1}) - P_{\tilde{Y}|Z}(k/N | Z_{2i})| \\ & \leq \frac{1}{N} + \frac{n-1}{n} \sum_{k=-N}^N f_k(Z_1, \dots, Z_n). \end{aligned}$$

Step 5: controlling the total-variation term. Fix $K \in [-N, N]$. Observe that if we resample a single Z_i (and write (Z'_1, \dots, Z'_n) to denote this new sample), the perturbation in f_k is given by

$$|f_k(Z_1, \dots, Z_n) - f_k(Z'_1, \dots, Z'_n)| \leq \frac{2}{n} \cdot 2 = \frac{4}{n}.$$

This follows by noting that the resampled Z'_i can alter at most two of the summands.

Now, let $\tilde{Z}_1, \tilde{Z}_2, \dots$ be an infinite sequence of i.i.d. copies of Z . For each $n \geq 2$ and each $1 \leq i \leq n$, let $\tilde{Z}_{n,i}$ be the Euclidean-nearest neighbour of \tilde{Z}_i among $\{\tilde{Z}_j : 1 \leq j \leq n, j \neq i \text{ and } \tilde{Z}_j \geq \tilde{Z}_i\}$. Now, observe that for any $i \in [n]$,

$$\mathbb{E} \left[|P_{\tilde{Y}|Z}(k/N | Z_i) - P_{\tilde{Y}|Z}(k/N | Z_{i+1})| \right] = \mathbb{E}_{P_Z} \left[|P_{\tilde{Y}|Z}(k/N | \tilde{Z}_1) - P_{\tilde{Y}|Z}(k/N | \tilde{Z}_{n,1})| \right].$$

Since Z_1, \dots, Z_n are independent, by McDiarmid's inequality, for any $\delta > 0$

$$|f_k(Z_1, \dots, Z_n) - 2\mathbb{E}_{P_Z} \left[|P_{\tilde{Y}|Z}(k/N | \tilde{Z}_1) - P_{\tilde{Y}|Z}(k/N | \tilde{Z}_{n,1})| \right]| = o_P(1). \quad (45)$$

Now, we claim that

$$|P_{\tilde{Y}|Z}(k/N | \tilde{Z}_1) - P_{\tilde{Y}|Z}(k/N | \tilde{Z}_{n,1})| = o_P(1). \quad (46)$$

To prove this claim, fix an $\epsilon > 0$ and a $\delta > 0$. By Lusin's theorem, there exists a compactly supported continuous function g_k such that

$$\mathbb{P}_{\tilde{Z} \sim P_Z} \left\{ P_{\tilde{Y}|Z}(k/N | \tilde{Z}) \neq g_k(\tilde{Z}) \right\} < \epsilon.$$

Now, $\mathbb{P} \left\{ |P_{\tilde{Y}|Z}(k/N | \tilde{Z}_1) - P_{\tilde{Y}|Z}(k/N | \tilde{Z}_{n,1})| \geq \delta \right\}$ is upper bounded by

$$\mathbb{P} \left\{ |g_k(\tilde{Z}_1) - g_k(\tilde{Z}_{n,1})| \geq \delta \right\} + \mathbb{P} \left\{ P_{\tilde{Y}|Z}(k/N | \tilde{Z}_1) \neq g_k(\tilde{Z}_1) \right\} + \mathbb{P} \left\{ P_{\tilde{Y}|Z}(k/N | \tilde{Z}_{n,1}) \neq g_k(\tilde{Z}_{n,1}) \right\}.$$

By construction, $\mathbb{P} \left\{ P_{\tilde{Y}|Z}(k/N | \tilde{Z}_1) \neq g_k(\tilde{Z}_1) \right\} \leq \epsilon$ and by Lemma D.27, there exists an universal constant C_0 such that

$$\mathbb{P} \left\{ P_{\tilde{Y}|Z}(k/N | \tilde{Z}_{n,1}) \neq g_k(\tilde{Z}_{n,1}) \right\} \leq C_0 \mathbb{P} \left\{ P_{\tilde{Y}|Z}(k/N | \tilde{Z}_1) \neq g_k(\tilde{Z}_1) \right\} \leq C_0 \epsilon.$$

Moreover, by continuity of g_k and by Lemma D.27,

$$\lim_{n \rightarrow \infty} \mathbb{P} \left\{ |g_k(Z) - g_k(Z_{RN})| \geq \delta \right\} = 0.$$

Thus, we have that

$$\lim_{n \rightarrow \infty} \mathbb{P} \left\{ |P_{\tilde{Y}|Z}(k/N | Z) - P_{\tilde{Y}|Z}(k/N | Z_{RN})| \geq \delta \right\} \leq \epsilon(1 + C_0).$$

Since ϵ, δ was arbitrary, (46) holds.

Since the sequence of random variables in (46) is uniformly bounded, it further implies that

$$\mathbb{E}_{P_Z} \left[|P_{\tilde{Y}|Z}(k/N | \tilde{Z}_1) - P_{\tilde{Y}|Z}(k/N | \tilde{Z}_{n,1})| \right] = o_P(1).$$

Since the conclusion above holds for any $k \in [-N, N]$, we finally have that

$$\frac{1}{n} \sum_{i=1}^{\lfloor n/2 \rfloor} d_{W_1}(P_{Y|Z}(\cdot | Z_{2i-1}), P_{Y|Z}(\cdot | Z_{2i})) \leq \frac{1}{N} + \frac{n-1}{n} \sum_{k=-N}^N f_k(Z_1, \dots, Z_n) \leq \frac{1}{N} + o_P(1).$$

But, N was arbitrary to start with. Thus, it follows that

$$\frac{1}{n} \sum_{i=1}^{\lfloor n/2 \rfloor} d_{W_1}(P_{Y|Z}(\cdot | Z_{2i-1}), P_{Y|Z}(\cdot | Z_{2i})) = o_P(1),$$

which concludes the proof. \square

Lemma C.21. *For cross bin matching under any distribution $P_{Y,Z}$, it holds that*

$$\left| \frac{1}{n/2} \|\Delta \mathbf{Y}^+\|_2^2 - \mathbb{E} [\text{Dev}(P_{Y|Z})] \right| \leq C(L_W \vee 1) \left(\frac{1}{\sqrt{K}} + \frac{1}{\sqrt{n/K}} + \sqrt{\frac{\log(4/\delta)}{n}} \right)$$

where C is a universal constant.

Proof. Let $P_{(Y,Z)}$ be a joint distribution on $(Y, Z) \in \mathbb{R} \times \mathbb{R}$. Write P_Z and $P_{Y|Z}$ for the marginal and conditional distributions. Define also F_Z^{-1} and $F_{Y|Z}^{-1}(\cdot | z)$ as the generalized inverse CDF for P_Z and $P_{Y|Z}(\cdot | z)$.

We define some notation. Let $A_k \subseteq [n]$ denote the “ k th bin”—that is, if we sort Z values from smallest to largest, A_k corresponds to the indices appearing in positions $(k-1)m+1, \dots, m$ in the sorted list. If we write $Y_{k,(1)} \geq \dots \geq Y_{k,(m)}$ as the sorted values of Y in the k th bin $A_k = \{(k-1)m+1, \dots, km\}$, then the cross-bin matching returns

$$\widehat{\text{Dev}}_{\text{cross-bin}} = \sum_{k=1}^{K-1} \sum_{i=1}^{\lfloor m/2 \rfloor} (Y_{k,(i)} - Y_{k+1,(m+1-i)})_+^2.$$

Step 1: rewriting with a Lipschitz function. Define a function $f : [-1, 1]^m \times [-1, 1]^m \rightarrow \mathbb{R}$ as

$$f(y, y') = \sum_{i=1}^{\lfloor m/2 \rfloor} (y_{(i)} - y'_{(m+1-i)})_+^2$$

where $y_{(1)} \geq \dots \geq y_{(m)}$ and $y'_{(1)} \geq \dots \geq y'_{(m)}$ denote the sorted values of y and of y' . Note that, by construction,

$$\widehat{\text{Dev}}_{\text{cross-bin}} = \sum_{k=1}^{K-1} f(Y_{A_k}, Y_{A_{k+1}}).$$

This function satisfies several key properties:

Lemma C.22. *The function f satisfies the Lipschitz property*

$$|f(y, y') - f(\tilde{y}, \tilde{y}')| \leq 4\|y - \tilde{y}\|_1 + 4\|y' - \tilde{y}'\|_1.$$

Lemma C.23. *Let Q be any distribution on \mathbb{R} , and let $A = (A_1, \dots, A_m) \sim Q^m$ and $B = (B_1, \dots, B_m) \sim Q^m$. Then*

$$\left| \frac{1}{m/2} \mathbb{E}[f(A, B)] - \text{Dev}(Q) \right| \leq \frac{C}{\sqrt{m}}$$

for a universal constant C .

Step 2: concentration. First, observe that if we resample a single Y_i value (and write (Y'_1, \dots, Y'_n) to denote this new sample), the perturbation in $\widehat{\text{Dev}}_{\text{cross-bin}}$ is given by

$$\left| \sum_{k=1}^{K-1} f(Y_{A_k}, Y_{A_{k+1}}) - \sum_{k=1}^{K-1} f(Y'_{A_k}, Y'_{A_{k+1}}) \right| \leq \sum_{k=1}^{K-1} \left(4\|Y_{A_k} - Y'_{A_k}\|_1 + 4\|Y_{A_{k+1}} - Y'_{A_{k+1}}\|_1 \right) \leq 16$$

where the last step holds since, for each k , $\|Y_{A_k} - Y'_{A_k}\|_1 \leq 2 \cdot \mathbb{1}_{i \in A_k}$ and similarly $\|Y_{A_{k+1}} - Y'_{A_{k+1}}\|_1 \leq 2 \cdot \mathbb{1}_{i \in A_{k+1}}$. Therefore by McDiarmid’s inequality, since the Y_i ’s are independent conditional on Z_1, \dots, Z_n , with probability $\geq 1 - \delta/2$ it holds that

$$\left| \widehat{\text{Dev}}_{\text{cross-bin}} - \mathbb{E} \left[\widehat{\text{Dev}}_{\text{cross-bin}} \mid Z_1, \dots, Z_n \right] \right| \leq 8\sqrt{2n \log(4/\delta)}.$$

Step 3: an equivalent way to sample the data. Now let $V_1, \dots, V_n, U_1, \dots, U_n \stackrel{\text{iid}}{\sim} \text{Unif}[0, 1]$. First we define $Z_i = F_Z^{-1}(V_i)$, so that we have $Z_1, \dots, Z_n \stackrel{\text{iid}}{\sim} P_Z$. Then we define

$$Y_i = F_{Y|Z}^{-1}(U_i | Z_i)$$

for each i —this is equivalent to sampling $(Y_i, Z_i) \stackrel{\text{iid}}{\sim} P_{(Y,Z)}$.

Step 4: replace the Y values with oracle samples. Now let $0 < u_1^* \leq \dots \leq u_{K-1}^* < 1$ be any values and let $z_k^* = F_Z^{-1}(u_k^*)$. For each $i \in A_k$, define

$$\tilde{Y}_i = F_{Y|Z}^{-1}(U_i | z_k^*)$$

and for each $i \in A_{k+1}$, define

$$\tilde{Y}'_i = F_{Y|Z}^{-1}(U_i | z_k^*).$$

Now define

$$\widetilde{\text{Dev}}_{\text{cross-bin}} = \sum_{k=1}^{K-1} f(\tilde{Y}_{A_k}, \tilde{Y}'_{A_{k+1}}),$$

so that

$$\left| \widetilde{\text{Dev}}_{\text{cross-bin}} - \widehat{\text{Dev}}_{\text{cross-bin}} \right| \leq \sum_{k=1}^{K-1} \left| f(Y_{A_k}, Y_{A_{k+1}}) - f(\tilde{Y}_{A_k}, \tilde{Y}'_{A_{k+1}}) \right|.$$

By Lemma C.22, for each k ,

$$\left| f(Y_{A_k}, Y_{A_{k+1}}) - f(\tilde{Y}_{A_k}, \tilde{Y}'_{A_{k+1}}) \right| \leq 4\|Y_{A_k} - \tilde{Y}_{A_k}\|_1 + 4\|Y_{A_{k+1}} - \tilde{Y}'_{A_{k+1}}\|_1.$$

Taking expected values, then,

$$\begin{aligned} \mathbb{E} \left[\left| f(Y_{A_k}, Y_{A_{k+1}}) - f(\tilde{Y}_{A_k}, \tilde{Y}'_{A_{k+1}}) \right| \mid Z_1, \dots, Z_n \right] \\ \leq \mathbb{E} \left[4\|Y_{A_k} - \tilde{Y}_{A_k}\|_1 + 4\|Y_{A_{k+1}} - \tilde{Y}'_{A_{k+1}}\|_1 \mid Z_1, \dots, Z_n \right]. \end{aligned}$$

Now we calculate an upper bound. For any $k = 1, \dots, K-1$ and any $i \in A_k$,

$$\mathbb{E} \left[|Y_i - \tilde{Y}_i| \mid Z_1, \dots, Z_n \right] = \text{d}_W(P_{Y|Z}(\cdot | Z_i), P_{Y|Z}(\cdot | z_k^*)),$$

by definition of the 1-Wasserstein distance, and by construction of Y_i and \tilde{Y}_i . Similarly for $i \in A_{k+1}$,

$$\mathbb{E} \left[|Y_i - \tilde{Y}'_i| \mid Z_1, \dots, Z_n \right] = \text{d}_W(P_{Y|Z}(\cdot | Z_i), P_{Y|Z}(\cdot | z_k^*)).$$

And, since we have made a smoothness assumption on the conditional distributions, since $Z_i = F_Z^{-1}(V_i)$ and $z_k^* = F_Z^{-1}(u_k^*)$, for each i and each k we have

$$\text{d}_W(P_{Y|Z}(\cdot | Z_i), P_{Y|Z}(\cdot | z_k^*)) \leq L_W |V_i - u_k^*|.$$

Therefore,

$$\mathbb{E} \left[\left| f(Y_{A_k}, Y_{A_{k+1}}) - f(\tilde{Y}_{A_k}, \tilde{Y}'_{A_{k+1}}) \right| \mid Z_1, \dots, Z_n \right] \leq 4L_W \sum_{i \in A_k \cup A_{k+1}} |V_i - u_k^*|.$$

We therefore have

$$\begin{aligned} & \left| \mathbb{E} \left[\widehat{\text{Dev}}_{\text{cross-bin}} \mid Z_1, \dots, Z_n \right] - \mathbb{E} \left[\widetilde{\text{Dev}}_{\text{cross-bin}} \mid Z_1, \dots, Z_n \right] \right| \\ & \leq \mathbb{E} \left[\left| \widetilde{\text{Dev}}_{\text{cross-bin}} - \widehat{\text{Dev}}_{\text{cross-bin}} \right| \mid Z_1, \dots, Z_n \right] \leq 4L_W \sum_{k=1}^{K-1} \sum_{i \in A_k \cup A_{k+1}} |V_i - u_k^*|. \end{aligned}$$

Next, by Lemma C.23, since \tilde{Y}_{A_k} and $\tilde{Y}'_{A_{k+1}}$ each represent m i.i.d. draws from $P_{Y|Z}(\cdot \mid z_k^*)$ for each k , we have

$$\left| \frac{1}{m/2} \mathbb{E} \left[f(\tilde{Y}_{A_k}, \tilde{Y}'_{A_{k+1}}) \mid Z_1, \dots, Z_n \right] - \text{Dev}(P_{Y|Z}(\cdot \mid z_k^*)) \right| \leq \frac{C}{\sqrt{m}}.$$

Therefore,

$$\left| \mathbb{E} \left[\widetilde{\text{Dev}}_{\text{cross-bin}} \mid Z_1, \dots, Z_n \right] - \frac{m}{2} \sum_{k=1}^{K-1} \text{Dev}(P_{Y|Z}(\cdot \mid z_k^*)) \right| \leq \frac{m}{2} \cdot (K-1) \cdot \frac{C}{\sqrt{m}}.$$

Combined with the calculations above, then,

$$\begin{aligned} & \left| \mathbb{E} \left[\widehat{\text{Dev}}_{\text{cross-bin}} \mid Z_1, \dots, Z_n \right] - \frac{m}{2} \sum_{k=1}^{K-1} \text{Dev}(P_{Y|Z}(\cdot \mid z_k^*)) \right| \\ & \leq 4L_W \sum_{k=1}^{K-1} \sum_{i \in A_k \cup A_{k+1}} |V_i - u_k^*| + \frac{m}{2} \cdot (K-1) \cdot \frac{C}{\sqrt{m}}. \end{aligned}$$

Step 5: another Wasserstein distance. By definition of the bins A_k (i.e., these bins reflect the ordering of the V_i 's), we can rewrite

$$\sum_{k=1}^{K-1} \sum_{i \in A_k} |V_i - u_k^*| = \sum_{k=1}^{K-1} \sum_{i=(k-1)m+1}^{km} |V_{(i)} - u_k^*|.$$

Since $u_1^* \leq \dots \leq u_{K-1}^*$, by definition of the 1-Wasserstein distance, we have

$$\frac{1}{m(K-1)} \sum_{k=1}^{K-1} \sum_{i=(k-1)m+1}^{km} |V_{(i)} - u_k^*| = d_W(\widehat{P}_V^{(1)}, \widehat{P}_{u^*}),$$

where $\widehat{P}_V^{(1)}$ is the empirical distribution of $V_{(1)}, \dots, V_{((K-1)m)}$, and \widehat{P}_{u^*} is the empirical distribution of u_1^*, \dots, u_{K-1}^* . By the triangle inequality,

$$\begin{aligned} d_W(\widehat{P}_V^{(1)}, \widehat{P}_{u^*}) & \leq d_W(\widehat{P}_V, \text{Unif}[0, 1]) + d_W(\widehat{P}_{u^*}, \text{Unif}[0, 1]) + d_W(\widehat{P}_V^{(1)}, \widehat{P}_V) \\ & \leq d_W(\widehat{P}_V, \text{Unif}[0, 1]) + d_W(\widehat{P}_{u^*}, \text{Unif}[0, 1]) + \frac{2m}{n}, \end{aligned}$$

where \widehat{P}_V is the empirical distribution of V_1, \dots, V_n (note that the last inequality holds since $V_{(1)}, \dots, V_{((K-1)m)}$ is a subset of the list V_1, \dots, V_n , containing all but $n - (K-1)m \leq 2m$ many entries). Therefore,

$$\sum_{k=1}^{K-1} \sum_{i \in A_k} |V_i - u_k^*| \leq m(K-1) \left(d_W(\widehat{P}_V, \text{Unif}[0, 1]) + d_W(\widehat{P}_{u^*}, \text{Unif}[0, 1]) \right) + 2m.$$

By an identical argument the same bound holds for $\sum_{k=1}^{K-1} \sum_{i \in A_{k+1}} |V_i - u_k^*|$. Therefore, combining with the results of the previous step,

$$\begin{aligned} & \left| \mathbb{E} \left[\widehat{\text{Dev}}_{\text{cross-bin}} \mid Z_1, \dots, Z_n \right] - \frac{m}{2} \sum_{k=1}^{K-1} \text{Dev}(P_{Y|Z}(\cdot \mid z_k^*)) \right| \\ & \leq 4L_W m(K-1) \left(d_W(\widehat{P}_V, \text{Unif}[0, 1]) + d_W(\widehat{P}_{u^*}, \text{Unif}[0, 1]) \right) + 8L_W m + \frac{m}{2} \cdot (K-1) \cdot \frac{C}{\sqrt{m}}. \end{aligned}$$

Step 6: averaging over the reference values. The above calculations hold for any fixed $0 < u_1^* \leq \dots \leq u_{K-1}^* < 1$. Now we will make these random: let $U_1^*, \dots, U_{K-1}^* \stackrel{\text{iid}}{\sim} \text{Unif}[0, 1]$ and let $U_{(1)}^* \leq \dots \leq U_{(K-1)}^*$ be the order statistics. And, let $Z_k^* = F_Z^{-1}(U_k^*)$, with $Z_{(k)}^* = F_Z^{-1}(U_{(k)}^*)$ being the order statistics. Since $Z_k^* \sim P_Z$ for each k , we have

$$\mathbb{E} \left[\sum_{k=1}^{K-1} \text{Dev}(P_{Y|Z}(\cdot \mid Z_{(k)}^*)) \right] = \mathbb{E} \left[\sum_{k=1}^{K-1} \text{Dev}(P_{Y|Z}(\cdot \mid Z_k^*)) \right] = (K-1) \mathbb{E} [\text{Dev}(P_{Y|Z})].$$

Combining with the previous step, and applying Jensen's inequality when taking expectation over the distribution of the Z_k^* 's (while conditioning on Z_1, \dots, Z_n), we then have

$$\begin{aligned} & \left| \mathbb{E} \left[\widehat{\text{Dev}}_{\text{cross-bin}} \mid Z_1, \dots, Z_n \right] - \frac{m}{2} \cdot (K-1) \mathbb{E} [\text{Dev}(P_{Y|Z})] \right| \\ & \leq 4L_W m(K-1) d_W(\widehat{P}_V, \text{Unif}[0, 1]) + 8L_W m + \frac{m}{2} \cdot (K-1) \cdot \frac{C}{\sqrt{m}} \\ & \quad + 4L_W n \mathbb{E} \left[d_W(\widehat{P}_{U^*}, \text{Unif}[0, 1]) \right], \end{aligned}$$

where \widehat{P}_{U^*} is the empirical distribution of U_1^*, \dots, U_{K-1}^* . Since $n - 2m \leq m(K-1) \leq n$, and $\mathbb{E} [\text{Dev}(P_{Y|Z})] \in [0, 4]$, then,

$$\begin{aligned} & \left| \mathbb{E} \left[\widehat{\text{Dev}}_{\text{cross-bin}} \mid Z_1, \dots, Z_n \right] - \frac{n}{2} \mathbb{E} [\text{Dev}(P_{Y|Z})] \right| \\ & \leq 4L_W m(K-1) d_W(\widehat{P}_V, \text{Unif}[0, 1]) + (8L_W + 4)m + \frac{m}{2} \cdot (K-1) \cdot \frac{C}{\sqrt{m}} \\ & \quad + 4L_W m(K-1) \mathbb{E} \left[d_W(\widehat{P}_{U^*}, \text{Unif}[0, 1]) \right]. \end{aligned}$$

By bounds on the empirical 1-Wasserstein distance from [Lei \(2020, Corollary 5.2.\)](#), since $V_i \stackrel{\text{iid}}{\sim} \text{Unif}[0, 1]$ and $U_k^* \stackrel{\text{iid}}{\sim} \text{Unif}[0, 1]$, we have

$$\mathbb{E} \left[d_W(\widehat{P}_V, \text{Unif}[0, 1]) \right] \leq C' n^{-1/2}, \quad \mathbb{E} \left[d_W(\widehat{P}_{U^*}, \text{Unif}[0, 1]) \right] \leq C' (K-1)^{-1/2}$$

for a universal constant C' , and so

$$\begin{aligned} & \left| \mathbb{E} \left[\widehat{\text{Dev}}_{\text{cross-bin}} \mid Z_1, \dots, Z_n \right] - \frac{n}{2} \mathbb{E} [\text{Dev}(P_{Y|Z})] \right| \\ & \leq 4L_W m(K-1) \left(d_W(\widehat{P}_V, \text{Unif}[0, 1]) - \mathbb{E} \left[d_W(\widehat{P}_V, \text{Unif}[0, 1]) \right] \right) \\ & \quad + 4CL_W \sqrt{n} + (8L_W + 4)m + \frac{m}{2} \cdot (K-1) \cdot \frac{C}{\sqrt{m}} + 4C' L_W m \sqrt{K-1}. \end{aligned}$$

Step 7: another concentration step. Since V_1, \dots, V_n are i.i.d., and the quantity $d_W(\widehat{P}_V, \text{Unif}[0, 1])$ can change value by at most $\frac{2}{n}$ if we resample one value V_i , by McDiarmid's inequality, with probability $\geq 1 - \delta/2$,

$$\left| d_W(\widehat{P}_V, \text{Unif}[0, 1]) - \mathbb{E} \left[d_W(\widehat{P}_V, \text{Unif}[0, 1]) \right] \right| \leq \sqrt{\frac{2 \log(4/\delta)}{n}}.$$

Step 8: combining everything. From all the steps above, with probability $\geq 1 - \delta$,

$$\begin{aligned} & \left| \widehat{\text{Dev}}_{\text{cross-bin}} - \frac{n}{2} \mathbb{E} [\text{Dev}(P_{Y|Z})] \right| \\ & \leq 8\sqrt{2n \log(4/\delta)} + 4L_W m(K-1) \cdot \sqrt{\frac{2 \log(4/\delta)}{n}} \\ & \quad + 4CL_W \sqrt{n} + (8L_W + 4)m + \frac{m}{2} \cdot (K-1) \cdot \frac{C}{\sqrt{m}} + 4C' L_W m \sqrt{K-1}. \end{aligned}$$

Letting C now denote a different universal constant, we have the final result □

D Auxiliary lemmas

Lemma D.24. *It holds that deterministically, for any $\alpha \in [0, 1]$,*

$$\frac{1}{2^L} \sum_{\mathbf{s} \in \{\pm 1\}^L} \mathbb{1} \{p(\mathbf{x}^{\mathbf{s}}) \leq \alpha\} \leq \alpha,$$

where $p(\cdot)$ is as defined in (7).

Proof. Consider a bijection $\sigma : [2^L] \rightarrow \{-1, 1\}^L$ such that $T(\mathbf{x}^{\sigma(1)}) \geq \dots \geq T(\mathbf{x}^{\sigma(2^L)})$, and let $r \in \{0, 1, \dots, 2^L\}$ be such that $\alpha \in [r/2^L, (r+1)/2^L)$. Then, since $\sum_{k=1}^{2^L} \mathbb{1} \{T(\mathbf{x}^{\sigma(k)}) \geq T(\mathbf{x}^{\sigma(j)})\} \in [2^L]$ for each $j \in [2^L]$, we have deterministically that

$$\begin{aligned} \frac{1}{2^L} \sum_{\mathbf{s} \in \{\pm 1\}^L} \mathbb{1} \{p(\mathbf{x}^{\mathbf{s}}) \leq \alpha\} &= \frac{1}{2^L} \sum_{j=1}^{2^L} \mathbb{1} \left\{ \frac{1}{2^L} \sum_{k=1}^{2^L} \mathbb{1} \{T(\mathbf{x}^{\sigma(k)}) \geq T(\mathbf{x}^{\sigma(j)})\} \leq \alpha \right\} \\ &= \frac{1}{2^L} \sum_{j=1}^{2^L} \mathbb{1} \left\{ \sum_{k=1}^{2^L} \mathbb{1} \{T(\mathbf{x}^{\sigma(k)}) \geq T(\mathbf{x}^{\sigma(j)})\} \leq r \right\} \\ &\leq \frac{1}{2^L} \sum_{j=1}^{2^L} \mathbb{1} \left\{ \sum_{k=1}^{2^L} \mathbb{1} \{k \leq j\} \leq r \right\} = \frac{1}{2^L} \sum_{j=1}^{2^L} \mathbb{1} \{j \leq r\} = \frac{r}{2^L} \leq \alpha. \end{aligned}$$

□

Lemma D.25. For any $p, q \in (0, 1)$, it holds that

$$H^2(\text{Ber}(p), \text{Ber}(q)) \leq \min \left\{ \frac{(p-q)^2}{2p(1-p)}, \frac{(p-q)^2}{2q(1-q)} \right\}.$$

Proof. We note that for any $p, q \in (0, 1)$,

$$\begin{aligned} H^2(\text{Ber}(p), \text{Ber}(q)) &= \frac{1}{2} \left[(\sqrt{p} - \sqrt{q})^2 + (\sqrt{1-p} - \sqrt{1-q})^2 \right] \\ &= \frac{1}{2} \left[\frac{(p-q)^2}{(\sqrt{p} + \sqrt{q})^2} + \frac{(p-q)^2}{(\sqrt{1-p} + \sqrt{1-q})^2} \right] \\ &\leq \frac{1}{2} \left[\frac{(p-q)^2}{p+q} + \frac{(p-q)^2}{(1-p) + (1-q)} \right] \\ &\leq \frac{1}{2} \left[\frac{(p-q)^2}{p} + \frac{(p-q)^2}{(1-p)} \right] \leq \frac{(p-q)^2}{2p(1-p)}. \end{aligned}$$

By symmetry, it is also bounded by $(p-q)^2/(2q(1-q))$, which proves the lemma. □

Lemma D.26. If S_1 and S_2 are finite subsets of \mathbb{R} with $\text{Med}(S_1 \cup S_2) < \text{Med}(S_1)$, then $\text{Med}(S_1 \cup S_2) \geq \text{Med}(S_2)$.

Proof. Let $S_1 = \{a_1, \dots, a_m\}$ with $a_1 \leq \dots \leq a_m$, $S_2 = \{b_1, \dots, b_n\}$ with $b_1 \leq \dots \leq b_n$, and $S_1 \cup S_2 = \{c_1, \dots, c_{m+n}\}$ with $c_1 \leq \dots \leq c_{m+n}$.

We first claim that $\text{Med}(S_2) < \text{Med}(S_1)$. To see this, observe that

$$\begin{aligned} |\{k \in [m+n] : c_k \leq \text{Med}(S_1 \cup S_2)\}| &\leq |\{k \in [m+n] : c_k < \text{Med}(S_1)\}| \\ &= |\{i \in [m] : a_i < \text{Med}(S_1)\}| + |\{j \in [n] : b_j < \text{Med}(S_1)\}|. \end{aligned}$$

Recalling that $|\{j \in [n] : b_j < \text{Med}(S_2)\}| \leq n/2$, there are four cases to consider:

Case 1: m and n both are odd. In this case,

$$|\{k \in [m+n] : c_k \leq \text{Med}(S_1 \cup S_2)\}| \geq (m+n)/2, \quad |\{i \in [m] : a_i < \text{Med}(S_1)\}| \leq (m-1)/2,$$

and hence $|\{j \in [n] : b_j < \text{Med}(S_1)\}| \geq (n+1)/2$, which proves the claim.

Case 2: m is odd, and n is even. In this case,

$$|\{k \in [m+n] : c_k \leq \text{Med}(S_1 \cup S_2)\}| \geq (m+n+1)/2, \quad |\{i \in [m] : a_i < \text{Med}(S_1)\}| \leq (m-1)/2,$$

and hence $|\{j \in [n] : b_j < \text{Med}(S_1)\}| \geq (n+2)/2$, which proves the claim.

Case 3: m is even, and n is odd. In this case,

$$|\{k \in [m+n] : c_k \leq \text{Med}(S_1 \cup S_2)\}| \geq (m+n+1)/2, \quad |\{i \in [m] : a_i < \text{Med}(S_1)\}| \leq m/2,$$

and hence $|\{j \in [n] : b_j < \text{Med}(S_1)\}| \geq (n+1)/2$, which proves the claim.

Case 4: m and n both are even. By definition,

$$\begin{aligned}\text{Med}(S_1) &= \frac{a_{m/2} + a_{m/2+1}}{2}, & \text{Med}(S_2) &= \frac{b_{n/2} + b_{n/2+1}}{2} \\ \text{and } \text{Med}(S_1 \cup S_2) &= \frac{c_{m/2+n/2} + c_{m/2+n/2+1}}{2}.\end{aligned}$$

Now, if $a_{m/2} = a_{m/2+1}$, then

$|\{k \in [m+n] : c_k \leq \text{Med}(S_1 \cup S_2)\}| \geq (m+n)/2$, $|\{i \in [m] : a_i < \text{Med}(S_1)\}| \leq m/2 - 1$, and thus, $|\{j \in [n] : b_j < \text{Med}(S_1)\}| \geq n/2 + 1$, which proves the claim.

Otherwise, when $a_{m/2} < a_{m/2+1}$

$$|\{k \in [m+n] : c_k \leq \text{Med}(S_1 \cup S_2)\}| \geq (m+n)/2, \quad |\{i \in [m] : a_i < \text{Med}(S_1)\}| = m/2,$$

and hence, $|\{j \in [n] : b_j < \text{Med}(S_1)\}| \geq n/2$ and $b_{n/2} < \text{Med}(S_1)$. Now again, there are four cases to consider.

- (i) If $b_{n/2} \leq a_{m/2} < a_{m/2+1} \leq b_{n/2+1}$, then $\text{Med}(S_1 \cup S_2) = \frac{a_{m/2} + a_{m/2+1}}{2} = \text{Med}(S_1)$ which is a contradiction.
- (ii) If $a_{m/2} < b_{n/2} \leq a_{m/2+1} \leq b_{n/2+1}$, then $\text{Med}(S_1 \cup S_2) = \frac{b_{n/2} + a_{m/2+1}}{2} > \frac{a_{m/2} + a_{m/2+1}}{2} = \text{Med}(S_1)$ which is a contradiction.
- (iii) If $a_{m/2} < b_{n/2} \leq b_{n/2+1} \leq a_{m/2+1}$, then $\text{Med}(S_1) > \text{Med}(S_1 \cup S_2) = \frac{b_{n/2} + b_{n/2+1}}{2} = \text{Med}(S_2)$.
- (iv) If $b_{n/2} \leq a_{m/2} \leq b_{n/2+1} < a_{m/2+1}$, then $\text{Med}(S_1) = \frac{a_{m/2} + a_{m/2+1}}{2} > \frac{b_{n/2} + b_{n/2+1}}{2} = \text{Med}(S_2)$.

This establishes the claim that $\text{Med}(S_2) < \text{Med}(S_1)$. If we had $\text{Med}(S_1 \cup S_2) < \text{Med}(S_2)$, then by the same argument as above, we could conclude that $\text{Med}(S_1) < \text{Med}(S_2)$ which is a contradiction. This proves the result. \square

Proof of Lemma C.22. First, for each $i = 1, \dots, \lfloor m/2 \rfloor$, since $t \mapsto (t)_+^2$ is 4-Lipschitz on $t \in [-2, 2]$,

$$|(y_{(i)} - y'_{(m+1-i)})_+^2 - (\tilde{y}_{(i)} - \tilde{y}'_{(m+1-i)})_+^2| \leq 4|y_{(i)} - \tilde{y}_{(i)}| + 4|y'_{(m+1-i)} - \tilde{y}'_{(m+1-i)}|.$$

Therefore,

$$|f(y, y') - f(\tilde{y}, \tilde{y}')| \leq 4 \sum_{i=1}^{\lfloor m/2 \rfloor} |y_{(i)} - \tilde{y}_{(i)}| + 4 \sum_{i=1}^{\lfloor m/2 \rfloor} |y'_{(m+1-i)} - \tilde{y}'_{(m+1-i)}|.$$

Writing $y_{(0)} = (y_{(1)}, \dots, y_{(m)})$ as the sorted vector, and same for $y'_{(0)}$, $\tilde{y}_{(0)}$, and $\tilde{y}'_{(0)}$, then,

$$|f(y, y') - f(\tilde{y}, \tilde{y}')| \leq 4\|y_{(0)} - \tilde{y}_{(0)}\|_1 + 4\|y'_{(0)} - \tilde{y}'_{(0)}\|_1.$$

Finally, it holds that

$$\|y_{(0)} - \tilde{y}_{(0)}\|_1 \leq \|y - \tilde{y}\|_1$$

by an ℓ_1 version of the rearrangement inequality, and similarly for y', \tilde{y}' . This completes the proof. \square

Proof of Lemma C.23. First let $A_{(1)} \geq \dots \geq A_{(m)}$ and $B_{(1)} \geq \dots \geq B_{(m)}$ be the order statistics, and write $A_{(0)} = (A_{(1)}, \dots, A_{(m)})$ and $B_{(0)} = (B_{(1)}, \dots, B_{(m)})$ for the sorted vectors. Then

$$f(A, B) = f(A_{(0)}, B_{(0)})$$

since $f(y, y')$ is invariant to permutations of y , and invariant to permutations of y' .

Next, let F_Q^{-1} be a generalized CDF for Q . Let $U_1, \dots, U_m \stackrel{\text{iid}}{\sim} \text{Unif}[0, 1]$, so that we can equivalently define $A_i = F_Q^{-1}(U_i)$ for $i = 1, \dots, m$. Define $A'_i = F_Q^{-1}(1 - U_i)$, and let A'_0 denote the sorted vector as before. Then

$$|f(A, B) - f(A, A')| = |f(A_{(0)}, B_{(0)}) - f(A_{(0)}, A'_{(0)})| \leq 4\|A'_0 - B_{(0)}\|_1$$

by Lemma C.22, and so applying Jensen's inequality,

$$|\mathbb{E}[f(A_{(0)}, B_{(0)})] - \mathbb{E}[f(A_{(0)}, A'_{(0)})]| \leq 4\mathbb{E}[\|A'_0 - B_{(0)}\|_1].$$

Next, let $U_{(1)} \leq \dots \leq U_{(m)}$ be the sorted values. Then we have $A_{(i)} = F_Q^{-1}(U_{(m+1-i)})$ and $A'_{(i)} = F_Q^{-1}(U_{(i)})$ for each i . Then by construction,

$$f(A, A') = \sum_{i=1}^{\lfloor m/2 \rfloor} (A_{(i)} - A'_{(m+1-i)})_+^2 = \sum_{i=1}^{\lfloor m/2 \rfloor} (F_Q^{-1}(U_{(m+1-i)}) - F_Q^{-1}(1 - U_{(m+1-i)}))_+^2.$$

Therefore,

$$\mathbb{E}[f(A, A')] = \mathbb{E}\left[\sum_{i=1}^{\lfloor m/2 \rfloor} (F_Q^{-1}(U_{(m+1-i)}) - F_Q^{-1}(1 - U_{(m+1-i)}))_+^2\right].$$

And by symmetry,

$$\mathbb{E}[f(A, A')] = \mathbb{E}[f(A', A)] = \mathbb{E}\left[\sum_{i=1}^{\lfloor m/2 \rfloor} (F_Q^{-1}(1 - U_{(i)}) - F_Q^{-1}(U_{(i)}))_+^2\right].$$

Therefore,

$$\mathbb{E}[f(A, A')] = \frac{1}{2}\mathbb{E}\left[\sum_{i=1}^m (F_Q^{-1}(1 - U_{(i)}) - F_Q^{-1}(U_{(i)}))^2 \cdot \mathbb{1}_{\mathcal{E}_i}\right]$$

where \mathcal{E}_i is the event that $i < \frac{m}{2}$ and $U_{(i)} < \frac{1}{2}$, or, $i > \frac{m}{2}$ and $U_{(i)} > \frac{1}{2}$. Therefore, since each term in the sum is bounded by 1,

$$\left|\mathbb{E}[f(A, A')] - \frac{1}{2}\mathbb{E}\left[\sum_{i=1}^m (F_Q^{-1}(1 - U_{(i)}) - F_Q^{-1}(U_{(i)}))^2\right]\right| \leq \frac{1}{2}\mathbb{E}\left[\sum_{i=1}^m \mathbb{1}_{\mathcal{E}_i^c}\right].$$

We can rewrite this as

$$\left|\mathbb{E}[f(A, A')] - \frac{1}{2}\mathbb{E}\left[\sum_{i=1}^m (F_Q^{-1}(1 - U_i) - F_Q^{-1}(U_i))^2\right]\right| \leq \frac{1}{2}\mathbb{E}\left[\sum_{i=1}^m \mathbb{1}_{\mathcal{E}_i^c}\right].$$

By definition, for each i ,

$$\mathbb{E} [(F_Q^{-1}(U_i) - F_Q^{-1}(1 - U_i))^2] = \text{Dev}(Q).$$

Therefore,

$$\left| \mathbb{E} [f(A, A')] - \frac{m}{2} \text{Dev}(Q) \right| \leq \frac{1}{2} \mathbb{E} \left[\sum_{i=1}^m \mathbb{1}_{\mathcal{E}_i^c} \right] \leq C'' \sqrt{m},$$

where the last step holds for a universal constant C'' by properties of the Binomial distribution, since

$$\sum_{i=1}^m \mathbb{1}_{\mathcal{E}_i^c} \leq 1 + \left| \sum_{i=1}^m \mathbb{1}_{U_i < 1/2} - \frac{m}{2} \right|$$

and $\sum_{i=1}^m \mathbb{1}_{U_i < 1/2} \sim \text{Binomial}(m, \frac{1}{2})$.

Returning to the initial calculations, then, we have

$$\left| \mathbb{E} [f(A, B)] - \frac{m}{2} \text{Dev}(Q) \right| \leq 4 \mathbb{E} [\|A'_0 - B_0\|_1] + C'' \sqrt{m}.$$

Next we turn to bounding $\|A'_0 - B_0\|_1$. Write $\widehat{P}_{A'}$ and \widehat{P}_B as the empirical distributions of A' and of B , respectively. Then

$$\|A'_0 - B_0\|_1 = m d_W(\widehat{P}_{A'}, \widehat{P}_B) \leq m d_W(\widehat{P}_{A'}, \text{Unif}[0, 1]) + m d_W(\widehat{P}_B, \text{Unif}[0, 1]).$$

By bounds on the empirical 1-Wasserstein distance from [Lei \(2020, Corollary 5.2.\)](#) we have

$$\mathbb{E} [d_W(\widehat{P}_{A'}, \text{Unif}[0, 1])] \leq C' m^{-1/2}$$

for a universal constant C' , and same for \widehat{P}_B . Therefore,

$$\left| \mathbb{E} [f(A, B)] - \frac{m}{2} \text{Dev}(Q) \right| \leq 8C' \sqrt{m} + C'' \sqrt{m},$$

which completes the proof. \square

Lemma D.27. *Let Z_1, Z_2, \dots be an infinite sequence of i.i.d. samples from P_Z . For each $n \geq 2$, let $Z_{n,i}$ be the Euclidean-nearest neighbour of Z_i among $\{Z_j : 1 \leq j \leq n, j \neq i \text{ and } Z_j \geq \tilde{Z}_i\}$. Then,*

(i) $Z_1 - Z_{n,1} \xrightarrow{P} 0$, and

(ii) there exists an universal constant C_0 such that for any measurable function $f : \mathbb{R} \rightarrow [0, \infty)$ and any n ,

$$\mathbb{E} [f(Z_{n,1})] \leq C_0 \mathbb{E} [f(Z_1)].$$

Proof. Let S be the support of P_Z , i.e., for any $z \in S$, any open interval containing z has strictly positive measure. Consequently, $\mathbb{P}_{P_Z} \{Z \in S\} = 1$. Take any $\epsilon > 0$, and consider the closed interval $[Z_1, Z_1 + \epsilon]$. Note that,

$$\mathbb{P} \{|Z_1 - Z_{n,1}| \geq \epsilon\} = \mathbb{E} [\mathbb{P} \{|Z_1 - Z_{n,1}| \geq \epsilon \mid Z_1\}] \leq \mathbb{E} \left[(1 - P_Z([Z_1, Z_1 + \epsilon]))^{n-1} \right]$$

Since $Z_1 \in S$ almost surely, $P_Z([Z_1, Z_1 + \epsilon)) > 0$ almost surely. Hence,

$$\lim_{n \rightarrow \infty} \mathbb{P}\{|Z_1 - Z_{n,1}| \geq \epsilon\} = 0,$$

which proves the first part of the result.

a calculation⁷ similar to [Azadkia and Chatterjee \(2021, Lemma 11.5\)](#), □

⁷Note that, in Lemma 11.5 of [Azadkia and Chatterjee \(2021\)](#) they work with nearest Euclidean neighbour, while here we deal with nearest neighbour to the right. However, the argument follows in a very similar fashion.