

The conditional permutation test for independence while controlling for confounders

Thomas B. Berrett*, Yi Wang[†], Rina Foygel Barber[†], Richard J. Samworth*

August 22, 2019

Abstract

We propose a general new method, the *conditional permutation test*, for testing the conditional independence of variables X and Y given a potentially high-dimensional random vector Z that may contain confounding factors. The proposed test permutes entries of X non-uniformly, so as to respect the existing dependence between X and Z and thus account for the presence of these confounders. Like the conditional randomization test of Candès et al. [7], our test relies on the availability of an approximation to the distribution of $X|Z$ —while Candès et al. [7]’s test uses this estimate to draw new X values, for our test we use this approximation to design an appropriate non-uniform distribution on permutations of the X values already seen in the true data. We provide an efficient Markov Chain Monte Carlo sampler for the implementation of our method, and establish bounds on the Type I error in terms of the error in the approximation of the conditional distribution of $X|Z$, finding that, for the worst case test statistic, the inflation in Type I error of the conditional permutation test is no larger than that of the conditional randomization test. We validate these theoretical results with experiments on simulated data and on the Capital Bikeshare data set.

1 Introduction

Independence is a central notion in statistical model building, as well as being a foundational concept for much of statistical theory. Originating with Francis Galton’s work on correlation at the end of the 19th century [24], many measures of dependence have been proposed, including mutual information, the Hilbert–Schmidt independence criterion, and distance covariance [8, 13, 30]; see also [15] for an overview. Simultaneously, a great deal of research effort has gone into developing several different tests

*Statistical Laboratory, University of Cambridge

[†]Department of Statistics, University of Chicago

of independence, for example based on ranks, kernel methods, copulas, and nearest neighbours [32, 17, 16, 6]. Permutation tests are particularly attractive due to their simplicity and their ability to control the Type I error (i.e., the false positive rate) without any distributional assumptions.

In practice, it is often conditional independence that is in fact of primary interest [9]. For instance, in generalized linear models for a response $Y \in \mathbb{R}$ regressed on a high-dimensional feature vector $(X, Z) = (X, Z^1, \dots, Z^p) \in \mathbb{R}^{p+1}$, the regression coefficient on feature X is zero if and only if Y and X are conditionally independent given the remaining p features, $Z = (Z^1, \dots, Z^p)$. In this paper, we will study the general problem of testing $X \perp\!\!\!\perp Y|Z$.¹ We are typically interested in the setting where X and Y are one-dimensional while Z is a high-dimensional set of confounding variables that we would like to control for, but our results are not specific to this setting.

Within standard parametric regression models, conditional independence tests are well-developed; unfortunately, however, they fail to control Type I error under model misspecification. In fact, the very recent work of Shah and Peters [22] has shown that, without placing some assumptions on the joint distribution of (X, Y, Z) , conditional testing is effectively impossible—when (X, Y, Z) is continuously distributed, they prove that there is no conditional independence test that both (1) controls Type I error over any null distribution (i.e., any distribution of (X, Y, Z) with $X \perp\!\!\!\perp Y|Z$), and (2) has better than random power against even one alternative hypothesis.

Our work seeks to complement this fundamental result of Shah and Peters [22] by demonstrating that, given some additional knowledge, namely an approximation to the conditional distribution of X given Z , one can in fact derive conditional independence tests that are approximately valid in finite samples, and that have non-trivial power.

1.1 Summary of contributions

In this paper, we introduce a new method, called the conditional permutation test (CPT), which is inspired by the conditional randomization test (CRT) of Candès et al. [7]. The CPT modifies the standard permutation test by using available distributional information to account correctly for the confounding variables Z , which leads to a non-uniform distribution over the set of possible permutations π on the n observations in our data set, and restores Type I error control.

Implementing the CPT is a challenging problem since we are sampling from a highly non-uniform distribution over the space of $n!$ permutations, but we propose a Monte Carlo sampler that yields an efficient implementation of the test. We additionally develop theoretical results examining the robustness of both the CPT and the CRT to slight errors in modeling assumptions, proving that Type I error is only slightly inflated

¹In the regression literature, it is more common to use the notation of regressing Y on (X^1, \dots, X^p) , and testing whether the coefficient on feature X^j is zero after controlling for the remaining features $X^{-j} = (X^1, \dots, X^{j-1}, X^{j+1}, \dots, X^p)$; this X^j and X^{-j} correspond to our X and Z , respectively.

in both tests when our available distributional information is only approximately correct. In fact, in the worst case, Type I error is always *less* inflated for the new CPT method as compared to the CRT. Our empirical results verify the greater robustness of the CPT, while maintaining comparable power in a range of scenarios.

2 Background

In this section, we briefly summarize several existing approaches to the problem of testing for dependence between X and Y in the presence of confounding variables. Before beginning, it will be helpful to define some brief notation. Throughout, we will assume that the data consists of i.i.d. data points $(X_i, Y_i, Z_i) \in \mathcal{X} \times \mathcal{Y} \times \mathcal{Z}$ for $i = 1, \dots, n$, and will write $\mathbf{X} = (X_1, \dots, X_n)$, $\mathbf{Y} = (Y_1, \dots, Y_n)$, and $\mathbf{Z} = (Z_1, \dots, Z_n)$.

2.1 Permutation tests

One key reason why handling conditional independence in nonparametric contexts is so challenging, is that the permutation approaches that are so effective for testing unconditional independence, $X \perp\!\!\!\perp Y$, cannot be directly applied when we seek to test conditional independence, $X \perp\!\!\!\perp Y|Z$. This is because it may be the case that the null hypothesis $H_0 : X \perp\!\!\!\perp Y|Z$ is true, but X and Y are highly marginally dependent due to correlation induced via each variable’s dependence on Z . Under this null, if we sample a permutation π of $\{1, \dots, n\}$ uniformly at random, then the permuted data set $(X_{\pi(1)}, Y_1), \dots, (X_{\pi(n)}, Y_n)$ may have a very different distribution from the original data set $(X_1, Y_1), \dots, (X_n, Y_n)$, due to the confounding effect of Z .

In certain settings, in particular where Z is categorical, there is a simple and well-known fix for this problem: we can group the observations according to their value of Z , and then permute within groups. For example, if $Z \in \{0, 1\}$ is binary, we could draw a permutation π that permutes the X_i ’s within the set of indices $\{i : Z_i = 0\}$, and separately permutes the X_i ’s within the set $\{i : Z_i = 1\}$. However, this strategy cannot be applied directly in the case where Z is continuously distributed, or where Z is discrete but with few repeated values (note that when Z is high-dimensional, even if it is discrete, each observation i will typically have a unique feature vector Z_i). In these settings, it is common to use a binning strategy, where first Z is discretized to fall into finitely many bins, and then the “permute within groups” strategy is deployed. However, Type I error control is no longer guaranteed, since the null hypothesis $H_0 : X \perp\!\!\!\perp Y|Z$ does not imply that $X \perp\!\!\!\perp Y|(Z \in \text{bin } b)$; the best we can usually hope for is that the latter statement would be approximately true under the null. Furthermore, in a high-dimensional setting, choosing these bins can itself be very challenging.

Apart from independence testing, permutation tests are also popular in other settings in which the null hypothesis is exchangeable [11]. Moreover, Roach and Valdar [18]

develop a theory of generalized permutation tests, primarily in the context of testing simple hypotheses, for non-exchangeable null models where the weights assigned to permutations are non-uniform.

2.2 The conditional randomization test

The conditional randomization test (CRT), proposed by Candès et al. [7], works in a setting where no assumptions are made about the distribution of the response variable Y , but instead, it is assumed that the conditional distribution of X given Z is known. In practice, in semi-supervised learning settings where unlabeled data (X, Z) are easier to obtain than labeled data (X, Y, Z) , it may be possible to obtain a very accurate estimate of the conditional distribution of $X|Z$, but testing for independence with Y remains challenging due to limited sample size of the labeled data. Candès et al. [7, Section 1.3] give examples of applications where unlabeled (X, Z) data is amply available while labeled data (X, Y, Z) is scarce—for example, genome-wide association studies (GWAS), where it is important to determine whether a particular genetic variant, X , affects a response Y such as disease status or some other phenotype, even after controlling for the rest of the genome, encoded in Z . Human genome data, i.e., (X, Z) data, is now plentiful, but labeled data (X, Y, Z) is expensive; if we do not know the disease status Y of the individuals in previously collected samples, we need to obtain the (X, Y, Z) samples ourselves.

Assuming then that the distribution of $X|Z$ is known (or is estimated accurately from a large sample of unlabeled data), the CRT operates by sampling a new copy of the X values in the data set. Letting $Q(\cdot|z)$ denote the distribution of X given $Z = z$, conditional on Z_1, \dots, Z_n , the CRT draws

$$X_i^{(1)} \sim Q(\cdot|Z_i),$$

independently for each $i = 1, \dots, n$, and independently of the observed X_i 's and Y_i 's. (In the special case where X is binary, earlier work by Rosenbaum [19] proposed a related test, referred to as a “conditional permutation test” but which in fact resamples X by estimating $\mathbb{P}\{X = 1 | Z\}$ with a logistic model.)

Under the null hypothesis H_0 that $X \perp\!\!\!\perp Y|Z$, we see that

$$(X|Y = y, Z = z) \stackrel{d}{=} (X|Z = z) \sim Q(\cdot|z),$$

where $\stackrel{d}{=}$ denotes equality in distribution. This means that

$$(\mathbf{X}^{(1)}, \mathbf{Y}, \mathbf{Z}) \stackrel{d}{=} (\mathbf{X}, \mathbf{Y}, \mathbf{Z}) \text{ under } H_0,$$

where $\mathbf{X}^{(1)} = (X_1^{(1)}, \dots, X_n^{(1)})$. Any large differences between these two triples—for instance, if \mathbf{Y} is highly correlated with \mathbf{X} but not with $\mathbf{X}^{(1)}$ —can therefore be interpreted as evidence against the null hypothesis. In order to construct a test of H_0 , then,

the CRT repeats this process M times, sampling

$$(X_i^{(m)} | \mathbf{X}, \mathbf{Y}, \mathbf{Z}) \sim Q(\cdot | Z_i), \text{ independently for } i = 1, \dots, n \text{ and } m = 1, \dots, M$$

to form control vectors $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)}$. Under the null hypothesis, the triples $(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$, $(\mathbf{X}^{(1)}, \mathbf{Y}, \mathbf{Z}), \dots, (\mathbf{X}^{(M)}, \mathbf{Y}, \mathbf{Z})$ are all identically distributed; in fact, they are exchangeable. For any statistic $T = T(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$ that is chosen in advance (or, at least, without looking at \mathbf{X}), the random variables

$$T(\mathbf{X}, \mathbf{Y}, \mathbf{Z}), T(\mathbf{X}^{(1)}, \mathbf{Y}, \mathbf{Z}), \dots, T(\mathbf{X}^{(M)}, \mathbf{Y}, \mathbf{Z}) \quad (1)$$

are therefore exchangeable as well. We can compute a p-value by ranking the value obtained from the true \mathbf{X} vector against the values obtained from the CRT's copies:

$$p = \frac{1 + \sum_{m=1}^M \mathbb{1} \{T(\mathbf{X}^{(m)}, \mathbf{Y}, \mathbf{Z}) \geq T(\mathbf{X}, \mathbf{Y}, \mathbf{Z})\}}{1 + M}.$$

The exchangeability of the random variables in (1) ensures that this is a valid p-value under the null, i.e., it satisfies $\mathbb{P} \{p \leq \alpha\} \leq \alpha$ for all $\alpha \in [0, 1]$ if the null hypothesis H_0 is true.

The “model-X knockoffs” framework of Candès et al. [7] also extends the CRT technique to the high-dimensional variable selection setting, where each of p features is tested in turn for conditional independence with the response Y , with the goal of false discovery rate control. In this framework, only a single copy of each feature is created. The robustness of the model-X knockoffs method, with respect to errors in the conditional distributions used to construct the knockoff copies of each feature (analogous to the $\mathbf{X}^{(m)}$'s above), was studied by Barber et al. [3].

2.3 Other tests of conditional independence

Before introducing our new work, we give a brief overview of some additional conditional independence testing methods proposed in the literature. Many methods assume some parametric model for the response Y , such as a linear model, $Y = \alpha X + \beta^\top Z + (\text{noise})$, in which case the problem reduces to testing whether $\alpha = 0$. This can be tested by, for instance, computing an estimate $\hat{\beta}$ and testing whether the residual $Y - \hat{\beta}^\top Z$ is correlated with X . Belloni et al. [4] propose a variant on this approach, which assumes approximate linear models for both Y and X . Their method regresses both X and Y on Z , then tests for correlation between the two resulting residual vectors; this “double regression” offers superior performance by removing much of the bias coming from errors in estimating the effect of Z . Shah and Peters [22] consider a more general double regression framework, assuming that the conditional means $\mathbb{E}[X | Z = z]$ and $\mathbb{E}[Y | Z = z]$ can be estimated at a sufficiently fast rate.

Away from the regression setting, many proposed methods are based on using kernel representations or low-dimensional projections of the data. Tests based on embedding

the data into reproducing kernel Hilbert spaces are studied in, for example, Fukumizu et al. [12], Zhang et al. [33] and Strobl et al. [25]. Other works use permutations of the data, including Doran et al. [10] and Sen et al. [21], where the methods have the flavor of binning Z and then permuting within groups. Bergsma [5], Song [23] and Veraverbeke et al. [31] study copula methods for testing conditional independence. There is also a large literature on extending measures of marginal independence to the conditional setting, including partial distance covariance [29]; conditional mutual information [20]; characteristic functions [26]; Hellinger distances [27]; and smoothed empirical likelihoods [28].

A related problem is that of testing the null hypothesis that a certain treatment has no effect in a randomized experiment. In the treatment effects literature it is common to calculate p-values by comparing a test statistic to null statistics based on randomly reassigning treatments in the data. However, in some situations, uniformly random reassignment is inappropriate, and does not result in valid p-values, due to the presence of some underlying structure; see Athey et al. [1] for network dependence and Hennessy et al. [14] for covariate imbalance. In such cases it is sometimes possible to develop non-uniform randomization schemes that result in valid p-values, as with the CPT and the CRT.

3 The conditional permutation test (CPT)

Recall that the conditional randomization test (CRT) [7] creates copies $\mathbf{X}^{(m)}$ of the vector \mathbf{X} sampled under the null hypothesis that $X \perp\!\!\!\perp Y|Z$, by drawing

$$\mathbf{X}^{(m)}|\mathbf{X}, \mathbf{Y}, \mathbf{Z} \sim Q^n(\cdot|\mathbf{Z}), \text{ independently for } m = 1, \dots, M, \quad (2)$$

where we define $Q^n(\cdot|\mathbf{Z}) := Q(\cdot|Z_1) \times \dots \times Q(\cdot|Z_n)$. This mechanism creates copies $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)}$ that are exchangeable with the original vector \mathbf{X} under the null hypothesis that $X \perp\!\!\!\perp Y|Z$.

Our proposed method, the conditional permutation test (CPT), is a variant on the CRT, with $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)}$ drawn as in (2) but under the constraint that each $\mathbf{X}^{(m)}$ must be a permutation of the original vector \mathbf{X} . Once we have drawn $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)}$, they will then be used exactly as for the CRT—given some predefined statistic $T = T(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$, our p-value is given by

$$p = \frac{1 + \sum_{m=1}^M \mathbb{1} \{T(\mathbf{X}^{(m)}, \mathbf{Y}, \mathbf{Z}) \geq T(\mathbf{X}, \mathbf{Y}, \mathbf{Z})\}}{1 + M}. \quad (3)$$

All that remains, then, is to specify how these permuted copies $\mathbf{X}^{(m)}$ will be drawn.

In order to draw the $\mathbf{X}^{(m)}$'s, we first need to define some notation. Let \mathcal{S}_n denote the set of permutations on the indices $\{1, \dots, n\}$. Given any vector $\mathbf{x} = (x_1, \dots, x_n)$ and

any permutation $\pi \in \mathcal{S}_n$, define $\mathbf{x}_\pi = (x_{\pi(1)}, \dots, x_{\pi(n)})$, i.e., the vector \mathbf{x} with its entries reordered according to the permutation π .

The CPT copies $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)}$ are then drawn as follows: after observing $\mathbf{X}, \mathbf{Y}, \mathbf{Z}$, we draw M permutations $\pi^{(1)}, \dots, \pi^{(M)}$ according to the conditional distribution of $X|Z$, and then apply these permutations to \mathbf{X} . Specifically, let

$$\mathbf{X}^{(m)} = \mathbf{X}_{\pi^{(m)}} \quad \text{where} \quad \mathbb{P} \{ \pi^{(m)} = \pi \mid \mathbf{X}, \mathbf{Y}, \mathbf{Z} \} = \frac{q^n(\mathbf{X}_\pi | \mathbf{Z})}{\sum_{\pi' \in \mathcal{S}_n} q^n(\mathbf{X}_{\pi'} | \mathbf{Z})}. \quad (4)$$

Here we let $q(\cdot|z)$ be the density of the distribution $Q(\cdot|z)$ (i.e., $q(\cdot|z)$ is the conditional density of X given $Z = z$), with respect to some base measure ν on \mathcal{X} that does not depend on z . We write $q^n(\cdot | \mathbf{Z}) := q(\cdot | Z_1) \cdots q(\cdot | Z_n)$ to denote the product density. Note that we are not assuming a continuous distribution necessarily; the base measure may be discrete, allowing X to be discrete as well.

Why is this the right distribution for drawing the permuted copies $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)}$? To understand this, it is helpful to consider a different formulation of the permutation scheme. Let $\mathbf{X}_{()} = (X_{(1)}, \dots, X_{(n)})$ be the order statistics of the list of values $\mathbf{X} = (X_1, \dots, X_n)$.² Define also $\mathbf{X}_{(\pi)} = (X_{(\pi(1))}, \dots, X_{(\pi(n))})$ for each $\pi \in \mathcal{S}_n$, and let $\Pi \in \mathcal{S}_n$ be the permutation given by the ranks of the true observed vector \mathbf{X} , so that $\mathbf{X} = \mathbf{X}_{(\Pi)}$. In other words, $\mathbf{X}_{()}$ gives the order statistics of \mathbf{X} , and Π reveals the ranks; together these two pieces of information are sufficient to reconstruct \mathbf{X} .³

Under the null hypothesis that $X \perp\!\!\!\perp Y|Z$, we can verify that the distribution of the true ranks Π , conditional on \mathbf{Y}, \mathbf{Z} as well as on the order statistics $\mathbf{X}_{()}$, is given by

$$\mathbb{P} \{ \Pi = \pi \mid \mathbf{X}_{()}, \mathbf{Y}, \mathbf{Z} \} = \frac{q^n(\mathbf{X}_{(\pi)} | \mathbf{Z})}{\sum_{\pi' \in \mathcal{S}_n} q^n(\mathbf{X}_{(\pi')} | \mathbf{Z})}. \quad (5)$$

Furthermore, examining the definition (4) of the CPT copies $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)}$, we can see that the CPT can equivalently be defined by

$$\mathbf{X}^{(m)} = \mathbf{X}_{(\Pi^{(m)})} \quad \text{where} \quad \Pi^{(m)} | \mathbf{X}_{()}, \mathbf{Y}, \mathbf{Z} \text{ is drawn from (5)}. \quad (6)$$

In fact, comparing with (4), we see that $\Pi^{(m)} = \Pi \circ \pi^{(m)}$.

The following theorem formalizes the above intuition, and verifies that this procedure yields a valid test of H_0 .

²In the setting where $\mathcal{X} = \mathbb{R}$, we can of course use the usual ordering on \mathbb{R} . In the general case we can simply take an arbitrary total ordering on \mathcal{X} ; the choice of ordering is irrelevant as its only role is to allow us to observe the set of values of \mathbf{X} without knowing which one corresponds to which data point.

³If the unlabeled values $X_{(i)}$ are not unique, then formally, we define Π by choosing it uniformly at random from the set of all permutations that satisfy this condition.

Theorem 1. Assume that $H_0 : X \perp\!\!\!\perp Y|Z$ is true, and that the conditional distribution of $X|Z$ is given by $Q(\cdot|Z)$. Suppose that $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)}$ are drawn i.i.d. from the CPT sampling scheme given in (4). Then the $M + 1$ triples

$$(\mathbf{X}, \mathbf{Y}, \mathbf{Z}), (\mathbf{X}^{(1)}, \mathbf{Y}, \mathbf{Z}), \dots, (\mathbf{X}^{(M)}, \mathbf{Y}, \mathbf{Z})$$

are exchangeable. In particular, this implies that for any statistic $T : \mathcal{X}^n \times \mathcal{Y}^n \times \mathcal{Z}^n \rightarrow \mathbb{R}$, the p -value defined in (3) is valid, satisfying $\mathbb{P}\{p \leq \alpha\} \leq \alpha$ for any desired Type I error rate $\alpha \in [0, 1]$ when H_0 is true.

Proof of Theorem 1. Our work above verified that, under H_0 , the true data vector \mathbf{X} and the CPT copies $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)}$ are permutations of \mathbf{X}_0 obtained via i.i.d. draws from (5), conditional on $\mathbf{X}_0, \mathbf{Y}, \mathbf{Z}$. Therefore, after marginalizing over $\mathbf{X}_0, \mathbf{Y}, \mathbf{Z}$, the $M + 1$ triples $(\mathbf{X}, \mathbf{Y}, \mathbf{Z}), (\mathbf{X}^{(1)}, \mathbf{Y}, \mathbf{Z}), \dots, (\mathbf{X}^{(M)}, \mathbf{Y}, \mathbf{Z})$ are exchangeable. \square

3.1 Comparing the CPT and CRT

To compare the construction of the copies $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)}$ in each of the two methods, for the CPT, the copies $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)}$ are i.i.d. draws from the null distribution of \mathbf{X} , conditional on $\mathbf{X}_0, \mathbf{Y}, \mathbf{Z}$. In comparison, the CRT copies defined in (2) are i.i.d. draws from the null distribution of \mathbf{X} conditioned on \mathbf{Y}, \mathbf{Z} —but without conditioning on \mathbf{X}_0 .

Each of these two constructions yields a valid test if the distribution $Q(\cdot|z)$, used to draw the (resampled or permuted) copies $\mathbf{X}^{(m)}$, is correct—that is, if we know the true conditional distribution of $X|Z$. This result is proved in Theorem 1 above for the CPT, while the analogous result for the CRT is proved in Candès et al. [7, Lemma 4.1]. However, if the null hypothesis is *not* true, which method might be more sensitive and better able to detect a non-null? Furthermore, what might occur for these two methods if $Q(\cdot|z)$ is not exactly correct? We next explore the difference between the two methods in greater depth to begin to address these questions.

Use of marginal distribution of \mathbf{X} In terms of how the tests are run, the difference between the CPT and CRT can be described as follows: while both tests use the (true or estimated) conditional distribution $Q(\cdot|Z)$, the CPT additionally uses the marginal distribution of the observed data vector \mathbf{X} , by observing its unlabeled values \mathbf{X}_0 . Intuitively, using this additional information can in some cases make the copies $\mathbf{X}^{(m)}$ more similar to the original \mathbf{X} , than for the CRT. Therefore, the CPT may be somewhat less likely to reject H_0 , which could lead to lower Type I error if H_0 is true, or reduced power to detect when H_0 is false. In Section 5, we will develop theory to examine the two tests’ robustness to errors in estimating the conditional distribution $Q(\cdot|Z)$, and we will compare the tests in terms of both Type I error and power in experiments in Section 6.

Invariance to base measure Since the CPT works only over permutations of the same set of X values, it follows that it is invariant to changes in the base measure on \mathcal{X} . To make this concrete, suppose that $q_1(\cdot|z)$ is another conditional density, with the property that there exist functions $h(\cdot), c(\cdot)$ such that $q_1(x|z) = q(x|z)h(x)c(z)$ for all $x \in \mathcal{X}$ and all $z \in \mathcal{Z}$. (Here we can think of $h(x)$ as changing the base measure on \mathcal{X} , while $c(z)$ adjusts the normalizing constants as needed.)

If this is the case, then running the CPT with q_1 in place of q will have no effect on the outcome—this is because we can calculate

$$q_1^n(\mathbf{X}_\pi|\mathbf{Z}) = \prod_{i=1}^n q(X_{\pi(i)}|Z_i)h(X_{\pi(i)})c(Z_i) = q^n(\mathbf{X}_\pi|\mathbf{Z}) \cdot \prod_{i=1}^n h(X_i)c(Z_i).$$

The first term, $q^n(\mathbf{X}_\pi|\mathbf{Z})$, is the same as for the CPT run with conditional density q , while the second term, $\prod_{i=1}^n h(X_i)c(Z_i)$, does not depend on the permutation π and therefore does not affect the resulting distribution of the sampled permutations. In other words, the CPT sampling distribution given in (4) is unchanged if we replace q with q_1 .

This means that the CPT is a valid test, i.e., the result of Theorem 1 holds, even if the conditional density $q(\cdot|z)$ is correct only up to a change in base measure—that is, Theorem 1 holds whenever the conditional distribution $Q(\cdot|Z)$ has a density of the form $q(x|z)h(x)c(z)$, for some functions $h(\cdot), c(\cdot)$. Indeed, in some settings, it may be substantially simpler to estimate the conditional density only up to base measure—for instance, we can consider a semiparametric model with a conditional density of the form $\exp\{x \cdot z^\top \theta - f(x) - g(z)\}$, in which case the CPT would only need to estimate the parametric component θ . In contrast, running the CRT requires being able to sample from the conditional distribution $Q(\cdot|Z)$, so we would need to approximate the full conditional density.

4 Sampling algorithms for the CPT

In order to run the CPT, we need to be able to sample permutations $\Pi^{(1)}, \dots, \Pi^{(M)}$ from the distribution given in (4). We now turn to the problem of generating such samples efficiently.

One simple approach would be to run a Metropolis–Hastings algorithm with a proposal distribution that, from a current state π , draws its proposed permutation π' uniformly at random. For even a moderate n , however, the acceptance odds ratio

$$\frac{q^n(\mathbf{X}_{\pi'}|\mathbf{Z})}{q^n(\mathbf{X}_\pi|\mathbf{Z})} = \frac{\prod_{i=1}^n q(X_{\pi'(i)}|Z_i)}{\prod_{i=1}^n q(X_{\pi(i)}|Z_i)} \quad (7)$$

will be extremely low for nearly all permutations π' (unless, of course, the dependence of X on Z is very weak). In other words, a uniformly drawn permutation π' is not likely to lead to a plausible vector of X values, leading to slow mixing times.

Algorithm 1 Parallelized pairwise sampler for the CPT

Input: Initial permutation $\Pi^{[0]}$, integer $S \geq 1$.

for $s = 1, 2, \dots, S$ **do**

Sample uniformly without replacement from $\{1, \dots, n\}$ to obtain disjoint pairs

$$(i_{s,1}, j_{s,1}), \dots, (i_{s,\lfloor n/2 \rfloor}, j_{s,\lfloor n/2 \rfloor}).$$

Draw independent Bernoulli variables $B_{s,1}, \dots, B_{s,\lfloor n/2 \rfloor}$ with odds ratios

$$\frac{\mathbb{P}\{B_{s,k} = 1\}}{\mathbb{P}\{B_{s,k} = 0\}} = \frac{q(X_{(\Pi^{[s-1]}(j_{s,k}))} | Z_{i_{s,k}}) \cdot q(X_{(\Pi^{[s-1]}(i_{s,k}))} | Z_{j_{s,k}})}{q(X_{(\Pi^{[s-1]}(i_{s,k}))} | Z_{i_{s,k}}) \cdot q(X_{(\Pi^{[s-1]}(j_{s,k}))} | Z_{j_{s,k}})}. \quad (9)$$

Define $\Pi^{[s]}$ by swapping $\Pi^{[s-1]}(i_{s,k})$ and $\Pi^{[s-1]}(j_{s,k})$ for each k with $B_{s,k} = 1$.

end for

As a second attempt, we can consider a different proposal distribution: from the current state π , we propose the permutation $\pi' = \pi \circ \sigma_{ij}$, where σ_{ij} is the permutation that swaps indices i and j , which are drawn at random. The acceptance odds ratio (7) now simplifies to

$$\frac{q(X_{\pi(j)} | Z_i) \cdot q(X_{\pi(i)} | Z_j)}{q(X_{\pi(i)} | Z_i) \cdot q(X_{\pi(j)} | Z_j)}. \quad (8)$$

The probability of accepting a swap will now be reasonably high; however, each step can only alter two of the n indices, again leading to slow mixing times.

4.1 A parallelized pairwise sampler

To address these issues, we propose a parallelized version of this pairwise algorithm. At each step, we first draw $\lfloor n/2 \rfloor$ disjoint pairs of indices from $\{1, \dots, n\}$. Next, independently and in parallel for each pair, we decide whether or not to swap this pair (i, j) , according to the odds ratio (8). This sampler is defined formally in Algorithm 1. For ease of our theoretical analysis, we will work with the order statistics $\mathbf{X}_{()}$, rather than the original ordered vector \mathbf{X} , in our sampler; this difference is only in the notation, i.e., the algorithm can equivalently be implemented with \mathbf{X} in place of $\mathbf{X}_{()}$.

The next theorem verifies that the resulting Markov chain yields the desired stationary distribution. (The proof of this theorem, and all remaining proofs, are given in Appendix A.)

Theorem 2. *For every initial permutation $\Pi^{[0]}$, the distribution (5) of the permutation Π conditional on $\mathbf{X}_{()}, \mathbf{Y}, \mathbf{Z}$ is a stationary distribution of the Markov chain defined in Algorithm 1. If additionally $q(x|z) > 0$ for all $x \in \mathcal{X}$ and all $z \in \mathcal{Z}$, then it is the unique stationary distribution.*

This result justifies the thought that, if Algorithm 1 is run for a sufficient number of steps S , then the resulting copy $\mathbf{X}_{(\Pi^{[S]})}$ acts as an appropriate control for \mathbf{X} in testing conditional independence. In fact, though, we can make a much stronger statement—since the original permutation Π also follows the distribution (5) conditional on $\mathbf{X}_{()}, \mathbf{Y}, \mathbf{Z}$ under the null, this means that by initializing Algorithm 1 at $\Pi^{[0]} = \Pi$ (that is, at the original data vector \mathbf{X}), we are initializing with a draw from the stationary distribution. Therefore $\mathbf{X}^{[S]} = \mathbf{X}_{(\Pi^{[S]})}$ is a draw from the target distribution at any S , and is a valid control for \mathbf{X} even if the number of steps S is small. Of course, if S is too small, then the control copy will be too similar to the original data vector \mathbf{X} , and our power to reject the null will be low; we explore this empirically in Section 6, and will see that the sampler mixes well at even a moderate S (e.g., in our experiments, we used $S = 50$).

In practice, we want to draw M copies, $\mathbf{X}^{(m)}$ for $m = 1, \dots, M$, and we need to ensure that the original data \mathbf{X} and each of the M permutations $\mathbf{X}^{(m)}$ are all exchangeable with each other. If we sample the permuted vectors $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)}$ sequentially, by running Algorithm 1 for $S \cdot M$ steps and extracting one copy $\mathbf{X}^{(m)}$ after each round of S steps, then we would not achieve exchangeability, since there would be some correlation between adjacent copies in this sequence. (Of course, in practice, if the number of steps S is chosen to be large, then the violation of exchangeability would be very mild.)

Instead, we can construct an exchangeable sampling mechanism with the following algorithm:

Algorithm 2 Exchangeable sampler for multiple draws from the CPT

Input: Initial permutation Π_{init} and integer $S \geq 1$.

Define $\Pi_{\#}$ by running Algorithm 1 initialized at $\Pi^{[0]} = \Pi_{\text{init}}$ for S steps.

for $m = 1, \dots, M$ (independently for each m) **do**

 Define $\Pi^{(m)}$ by running Algorithm 1 initialized at $\Pi^{[0]} = \Pi_{\#}$ for S steps.

end for

Algorithm 2 provides an exchangeable sampling mechanism, since the permutation $\Pi_{\#}$ is at the “center”, lying S steps away from each of the permutations $\Pi, \Pi^{(1)}, \dots, \Pi^{(M)}$. The following result verifies exchangeability:

Theorem 3. *Let $\mathbf{X}_{()}$ and Π be the order statistics and ranks of \mathbf{X} , as defined previously, so that $\mathbf{X} = \mathbf{X}_{(\Pi)}$. Let $\Pi^{(1)}, \dots, \Pi^{(M)}$ be the output of Algorithm 2, when initialized at $\Pi_{\text{init}} = \Pi$, and let $\mathbf{X}^{(m)} = \mathbf{X}_{(\Pi^{(m)})}$ for each $m = 1, \dots, M$. Assume that the null hypothesis that $X \perp\!\!\!\perp Y|Z$ holds, and the conditional distribution of $X|Z$ is given by $Q(\cdot|Z)$, so that the distribution of Π conditional on $\mathbf{X}_{()}, \mathbf{Y}, \mathbf{Z}$ is given by (5). Then the triples $(\mathbf{X}, \mathbf{Y}, \mathbf{Z}), (\mathbf{X}^{(1)}, \mathbf{Y}, \mathbf{Z}), \dots, (\mathbf{X}^{(M)}, \mathbf{Y}, \mathbf{Z})$ are exchangeable.*

This result ensures that the results of Theorem 1 hold when the permuted vectors

$\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)}$ are obtained via the exchangeable sampler.

5 Robustness of the CPT and CRT

We next consider whether the CPT and CRT, based on resampling X from a known or estimated conditional distribution given Z , are robust to slight errors in this distribution. Suppose that the conditional distribution $Q(\cdot|Z)$ that we use for sampling when running the CPT or CRT is only an approximation to the true conditional, denoted by $Q_\star(\cdot|Z)$. In this section we provide bounds on the excess Type I error of the CPT and CRT as a function of the difference between the true conditional Q_\star and its approximation Q . Throughout, we will assume that the statistic $T : \mathcal{X}^n \times \mathcal{Y}^n \times \mathcal{Z}^n \rightarrow \mathbb{R}$ used in the test, as well as the approximation Q to the conditional distribution, are chosen independently of \mathbf{X}, \mathbf{Y} . For instance, in many applications, we may have access to unlabeled data, i.e., draws of (X, Z) without Y , which we can use to construct an estimate Q .

Our first result demonstrates that, conditional on \mathbf{Y}, \mathbf{Z} , the excess Type I error of both the CPT and the CRT is bounded by the total variation distance between Q_\star and Q . (For any two distributions Q_1, Q_2 defined on the same probability space, the total variation distance is defined as $d_{\text{TV}}(Q_1, Q_2) = \sup_A |Q_1(A) - Q_2(A)|$, where the supremum is taken over all measurable sets.)

Theorem 4. *Assume that $H_0 : X \perp\!\!\!\perp Y|Z$ is true, and that the conditional distribution of $X|Z$ is given by $Q_\star(\cdot|Z)$. For a fixed integer $M \geq 1$, let $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)}$ be copies of \mathbf{X} generated either from the CRT (2), from the CPT (4), or from the exchangeable sampler for the CPT (Algorithm 2) with any fixed parameter $S \geq 1$, using an estimate Q of the true conditional distribution Q_\star .*

Then, for any desired Type I error rate $\alpha \in [0, 1]$,

$$\mathbb{P} \{p \leq \alpha \mid \mathbf{Y}, \mathbf{Z}\} \leq \alpha + d_{\text{TV}}(Q_\star^n(\cdot|\mathbf{Z}), Q^n(\cdot|\mathbf{Z})),$$

where p is the p -value computed in (3), and the probability is taken with respect to the distribution of $\mathbf{X}, \mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)}$ conditional on \mathbf{Y}, \mathbf{Z} .

Of course, we can also bound the Type I error rate unconditionally, with

$$\mathbb{P} \{p \leq \alpha\} \leq \alpha + \mathbb{E} [d_{\text{TV}}(Q_\star^n(\cdot|\mathbf{Z}), Q^n(\cdot|\mathbf{Z}))],$$

which we obtain from the result above by marginalizing over \mathbf{Y}, \mathbf{Z} .

This result ensures that, if Q is a good approximation to Q_\star , then both the CPT and CRT will have at most a mild increase in their Type I error. Of course, Theorem 4 is a worst-case result, proved with respect to an arbitrary statistic T which may be chosen adversarially so as to be maximally sensitive to errors in estimating the true

conditional distribution Q_\star . In practice, we might expect that the simple statistics T that we would most often use, such as correlation between \mathbf{X} and \mathbf{Y} , could be more robust to errors than the theorem suggests.

While Theorem 4 provides an upper bound on the Type I error for both the CPT and the CRT, we do not yet have a comparison between the two. The following theorem proves that, for the case of the CRT, the upper bound is in fact tight when the number of copies $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)}$ is large:

Theorem 5. *Under the setting and assumptions of Theorem 4, there exists a statistic $T : \mathcal{X}^n \times \mathcal{Y}^n \times \mathcal{Z}^n \rightarrow \mathbb{R}$ such that, for the CRT,⁴*

$$\sup_{\alpha \in [0,1]} \left(\mathbb{P} \{p \leq \alpha \mid \mathbf{Y}, \mathbf{Z}\} - \alpha \right) \geq d_{\text{TV}}(Q_\star^n(\cdot \mid \mathbf{Z}), Q^n(\cdot \mid \mathbf{Z})) - 0.5(1 + o(1)) \sqrt{\frac{\log(M)}{M}}$$

as $M \rightarrow \infty$.

In other words, if we use the statistic T that is best able to detect errors in our conditional distribution, and choose α adversarially, then the excess Type I error of the CRT is exactly equal to $d_{\text{TV}}(Q_\star^n(\cdot \mid \mathbf{Z}), Q^n(\cdot \mid \mathbf{Z}))$ (up to a vanishing factor), and therefore is at least as high as that of the CPT under *any* statistic.

Unlike for the CRT, we have found that there is no simple characterization of the worst-case scenario for the CPT. In particular, for some specially constructed distributions on (X, Y, Z) , we can show that the CPT achieves the same lower bound as given in Theorem 5 for the CRT (again, under a worst-case choice of the statistic T), but for other joint distributions on (X, Y, Z) we can verify that the CPT cannot achieve this error rate. In particular, since the CPT is invariant to the base measure (as discussed in Section 3.1), if $Q(\cdot \mid z)$ is correct up to the base measure, then the excess Type I error of CRT may be as large as $d_{\text{TV}}(Q_\star^n(\cdot \mid \mathbf{Z}), Q^n(\cdot \mid \mathbf{Z}))$ while the CPT is guaranteed to control Type I error at level α .

It is important to note that the lower bound for the CRT in Theorem 5 applies only to a specific worst-case statistic T , and does not guarantee that the excess error of the CRT will bound that of the CPT when both tests use some other statistic T . However, in Section 6 we will see that empirically, the CPT often yields a far lower Type I error than the CRT in simulations. Thus, we interpret Theorem 5 as giving us a partial theoretical understanding of this phenomenon, since it only addresses the worst-case statistic.

5.1 When is the total variation distance small?

In order for Theorem 4 to have practical implications, we need to verify that there are settings where, although the true distribution Q_\star of $X \mid Z$ is unknown, it can be

⁴To be more precise with the constant, we can replace $0.5(1 + o(1))$ with 2.5 for any $M \geq 2$.

estimated to high accuracy, with $d_{\text{TV}}(Q_{\star}^n(\cdot|\mathbf{Z}), Q^n(\cdot|\mathbf{Z})) = o_p(1)$ (so that excess Type I error is guaranteed to be small). As discussed in Section 2.2, in many applications we may have a large unlabeled data set, say $(X_i^{\text{unlab}}, Z_i^{\text{unlab}}), i = 1, \dots, N$, with which we can compute an estimate Q of Q_{\star} . (In fact, as discussed by Barber and Candès [2] in the setting of model- X knockoffs, the unlabeled data set does not need to have the same distribution over (X, Z) as the labeled data, as long as the conditional distribution of $X|Z$ is the same.)

In this section, we briefly sketch two settings where, given a large unlabeled sample size N , our estimate Q is likely to satisfy $d_{\text{TV}}(Q_{\star}^n(\cdot|\mathbf{Z}), Q^n(\cdot|\mathbf{Z})) = o_p(1)$. Our results here are stated informally, with no technical details, since we aim only to give intuition for the settings where Theorem 4 is useful.

Parametric setting We will use Pinsker’s inequality relating total variation distance to the Kullback–Leibler divergence, namely,

$$d_{\text{TV}}^2(Q_{\star}^n(\cdot|\mathbf{Z}), Q^n(\cdot|\mathbf{Z})) \leq \frac{1}{2} d_{\text{KL}}(Q_{\star}^n(\cdot|\mathbf{Z}), Q^n(\cdot|\mathbf{Z})) = \frac{1}{2} \sum_{i=1}^n d_{\text{KL}}(Q_{\star}(\cdot|Z_i), Q(\cdot|Z_i)).$$

It is therefore sufficient to show that $\sum_{i=1}^n d_{\text{KL}}(Q_{\star}(\cdot|Z_i), Q(\cdot|Z_i)) = o_p(1)$.

In fact, if the true conditional distribution $Q_{\star}(\cdot|z)$ belongs to a parametric family, then this will typically hold whenever the unlabeled sample size satisfies $N \gg n \cdot k$, where k is the number of parameters defining the models in the family. Specifically, we can think of a setting where $Q_{\star}(\cdot|z)$ has density $f_{\theta_{\star}}(\cdot|z)$, where $\theta_{\star} \in \mathbb{R}^k$ is the unknown parameter vector while the family of densities $f_{\theta}(\cdot|z)$ is known. For example, suppose that $\mathcal{Z} = \mathbb{R}^{k-1}$, and the conditional distribution of $X|Z$ is given by

$$X|Z = z \sim \mathcal{N}(z^{\top} \beta_{\star}, \sigma_{\star}^2).$$

Then the unknown parameters are $\theta_{\star} = (\beta_{\star}, \sigma_{\star}^2)$ and standard least squares theory allows us to produce independent (maximum likelihood) estimates $\hat{\beta}, \hat{\sigma}^2$ satisfying

$$\hat{\beta} \sim N_{k-1}(\beta_{\star}, \sigma_{\star}^2 (\mathbf{Z}_{\text{unlab}}^{\top} \mathbf{Z}_{\text{unlab}})^{-1}), \quad \hat{\sigma}^2 \sim \frac{\sigma_{\star}^2}{N} \chi_{N-k+1}^2,$$

where $\mathbf{Z}_{\text{unlab}}$ is the $N \times (k-1)$ matrix with i th row Z_i^{unlab} . Thus, for any $z \in \mathcal{Z}$,

$$\begin{aligned} d_{\text{KL}}(Q_{\star}(\cdot|z), Q(\cdot|z)) &= d_{\text{KL}}(\mathcal{N}(z^{\top} \beta_{\star}, \sigma_{\star}^2), \mathcal{N}(z^{\top} \hat{\beta}, \hat{\sigma}^2)) \\ &= \log \frac{\hat{\sigma}}{\sigma_{\star}} + \frac{\sigma_{\star}^2}{2\hat{\sigma}^2} - \frac{1}{2} + \frac{(z^{\top} \hat{\beta} - z^{\top} \beta_{\star})^2}{2\hat{\sigma}^2} = O_p\left(\frac{1 + \|z\|^2}{N}\right) \end{aligned}$$

under mild conditions on the distribution of Z . Putting everything together, if Z has a finite second moment we then have

$$d_{\text{TV}}(Q_{\star}^n(\cdot|\mathbf{Z}), Q^n(\cdot|\mathbf{Z})) = O_p\left(\sqrt{n \cdot \frac{k}{N}}\right),$$

which is vanishing as long as the unlabeled sample size satisfies $N \gg n \cdot k$.

Nonparametric setting with binary data As a second example, suppose that $\mathcal{X} = \{0, 1\}$, so that estimating $Q_\star(\cdot|z)$ is equivalent to estimating the regression function $p_\star(z) := \mathbb{P}\{X = 1 \mid Z = z\}$. Assuming that this probability is bounded away from 0 and 1, and again applying Pinsker’s inequality, we see that, under mild conditions,

$$d_{\text{TV}}^2(Q_\star^n(\cdot|\mathbf{Z}), Q^n(\cdot|\mathbf{Z})) \leq \frac{1}{2} \sum_{i=1}^n d_{\text{KL}}(Q_\star(\cdot|Z_i), Q(\cdot|Z_i)) \asymp \sum_{i=1}^n (\widehat{p}(Z_i) - p_\star(Z_i))^2,$$

where $\widehat{p}(z)$ is our estimate of $p_\star(z) = \mathbb{P}\{X = 1 \mid Z = z\}$ based on the unlabeled sample.

Since we are working in a nonparametric setting, suppose that we estimate $p_\star(z) = \mathbb{P}\{X = 1 \mid Z = z\}$ via a kernel method, working in a low-dimensional space $\mathcal{Z} = \mathbb{R}^k$. Then standard nonparametric theory ensures that, at “most” values z , we can achieve error

$$(\widehat{p}(z) - p_\star(z))^2 \sim N^{-a_k},$$

where the exponent a_k is a small positive value, depending on both the ambient dimension k and the properties of the function $z \mapsto p_\star(z)$ (e.g., smoothness or Lipschitz properties). Therefore, we can expect to have

$$d_{\text{TV}}(Q_\star^n(\cdot|\mathbf{Z}), Q^n(\cdot|\mathbf{Z})) \lesssim \sqrt{n \cdot N^{-a_k}},$$

which is vanishing whenever the unlabeled sample size N is sufficiently large relative to the labeled sample size n .

6 Empirical results

We next examine the empirical performance of the CPT and CRT on simulated data, and on real data from the Capital Bikeshare system. Code for reproducing all experiments is available on the authors’ websites.⁵

6.1 Simulated data: power and error control

The results of Section 5 show that the CPT is more robust than the CRT to errors in the estimated conditional distribution $Q(\cdot|Z)$, when the worst case test statistics $T(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$ are used. Our first aim here is to provide evidence to validate this result, and to show that this extra robustness is not only exhibited by the worst case test statistic but also for practical and simple choices of T . Our second aim is to examine the power of the CPT and CRT to detect deviations from the null hypothesis.

⁵Available at <http://www.stat.uchicago.edu/~rina/cpt.html>.

In all of our simulations we set $\alpha = 0.05$ as the desired Type I error rate, and use marginal absolute correlation $T(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) = |\text{Corr}(\mathbf{X}, \mathbf{Y})|$ as our test statistic. We generate $M = 500$ copies of \mathbf{X} under either CPT or CRT. To run the CPT, we use Algorithm 2 with $S = 50$ steps. All results are shown averaged over 1000 trials.

6.1.1 Simulations under the null

First we test whether the CPT and CRT show large increases in Type I error when the conditional distribution estimate $Q(\cdot|Z)$ is incorrect, in a setting where the null hypothesis $H_0 : X \perp\!\!\!\perp Y|Z$ holds.

We will have $X, Y \in \mathbb{R}$ and $Z \in \mathbb{R}^p$ for $p = 20$. We first draw independent parameter vectors

$$a, b \sim \mathcal{N}_p(0, \mathbf{I}_p).$$

The variables (X, Y, Z) are then generated as

$$Z \sim \mathcal{N}_p(0, \mathbf{I}_p), \quad X|Z \sim Q_\star(\cdot|Z), \quad Y|X, Z \sim \mathcal{N}(p^{-1}a^\top Z, 1),$$

where $Q_\star(\cdot|Z)$ will be specified below. (Note that $Y|X, Z$ depends on Z only, since we are working under the null hypothesis that $X \perp\!\!\!\perp Y|Z$.)

Throughout, the estimated conditional distribution of $X|Z$ will be given by $Q(\cdot|Z) = \mathcal{N}(b^\top Z, 1)$, but this estimate might not be exactly correct. We will consider several different sources of error in this model:

1. Nonlinear mean. One source of error comes from assuming a linear relationship between variables where this is in fact not the case. We choose sample size $n = 50$, and try three different simple examples, taking $Q_\star(\cdot|z) = \mathcal{N}(\mu(z), 1)$, where $\mu(\cdot)$ is given by:

(a) Quadratic: $\mu(z) = b^\top z + \theta(b^\top z)^2$,

(b) Cubic: $\mu(z) = b^\top z - \theta(b^\top z)^3$,

(c) Tanh: $\mu(z) = \tanh(\theta \cdot b^\top z)/\theta$.

In each case, $\theta \geq 0$ is the model misspecification parameter. Note that $\theta = 0$ corresponds to the case that $Q(\cdot|Z) = Q_\star(\cdot|Z)$, i.e., the estimate is indeed correct, while larger values of θ correspond to increasing errors.

2. Coefficients estimated on unlabeled data. Even if the form of the model for $X|Z$ is correct, the coefficients b may not be known perfectly. As described earlier, in many practical settings we may have access to ample unlabeled data (X, Z) , separate from our labeled data set of points (X, Y, Z) used to test the hypothesis of conditional independence. For this setting, we estimate the unknown coefficient vector b with \hat{b} , defined as the least-squares estimate using an

unlabeled sample $(X_i^{\text{unlab}}, Z_i^{\text{unlab}}), i = 1, \dots, N$, generated independently of the data points (X_i, Y_i, Z_i) . This experiment is repeated for unlabeled sample sizes $N = 50, 100, \dots, 500$. The labeled sample size is given by $n = 50$ in each case.

3. Coefficients estimated by reusing the data. Finally, in settings where unlabeled data may not be available, we may be tempted to estimate the model of $X|Z$ simply using our data points $(X_i, Y_i, Z_i), i = 1, \dots, n$. This approach is not covered by our theory (since the conditional distribution $Q(X|Z)$ is data-dependent in this case), but it is certainly of practical interest to see how the method performs in this setting. We test sample sizes $n = 50, 100, \dots, 500$, in each case estimating the unknown true coefficient vector b with \hat{b} , which in this case is now given by the least-squares regression of X on Z trained on the *same* data set, $(X_1, Z_1), \dots, (X_n, Z_n)$.

Results The plots in Figures 1 and 2 show the results of these experiments when we have a nonlinear mean, and when we estimate the coefficients using unlabeled data or reusing data, respectively. As the null hypothesis, $H_0 : X \perp Y|Z$, is true in all of these experiments, we would hope for the probability of rejection to be close to the nominal level of $\alpha = 0.05$, at least when the model misspecification parameter θ is not too large (for the nonlinear mean setting) or when the unlabeled sample size N or labeled sample size n is not too small (when the model coefficients are trained on unlabeled data or reused data).

For the nonlinear mean experiments, in Figure 1 we see that in many cases the CPT is significantly more robust than the CRT. The $\theta = 0$ cases confirm that both tests achieve the nominal Type I error level $\alpha = 0.05$ when the assumed distribution Q is correct. As the misspecification parameter θ increases (so that the model $Q(\cdot|z)$ that we use for running CPT or CRT, grows farther from the true model $Q_\star(\cdot|z)$), we see that both methods suffer an inflation of the Type I error level, but for the CPT the excess Type I error is substantially lower than that of the CRT.

Next, we turn to the setting where the estimated model $Q(\cdot|z)$ is obtained by regressing X on Z using either a separate unlabeled data set, shown in Figure 2a, or by reusing the same data set, shown in Figure 2b. The results are encouraging, showing that, when using unlabeled data, the Type I error is already very close to the nominal level as soon as the unlabeled sample size N is larger than n . When reusing the data, the method in fact appears to be somewhat conservative at smaller sample sizes n —the cause of this phenomenon is an interesting question we hope to study in future work.

6.1.2 Simulations under the alternative

Our final simulation concerns the power of the tests. Here we generate Z as before, and generate $X|Z \sim \mathcal{N}(b^\top Z, 1)$, exactly according the assumed distribution $Q(\cdot|Z)$, so

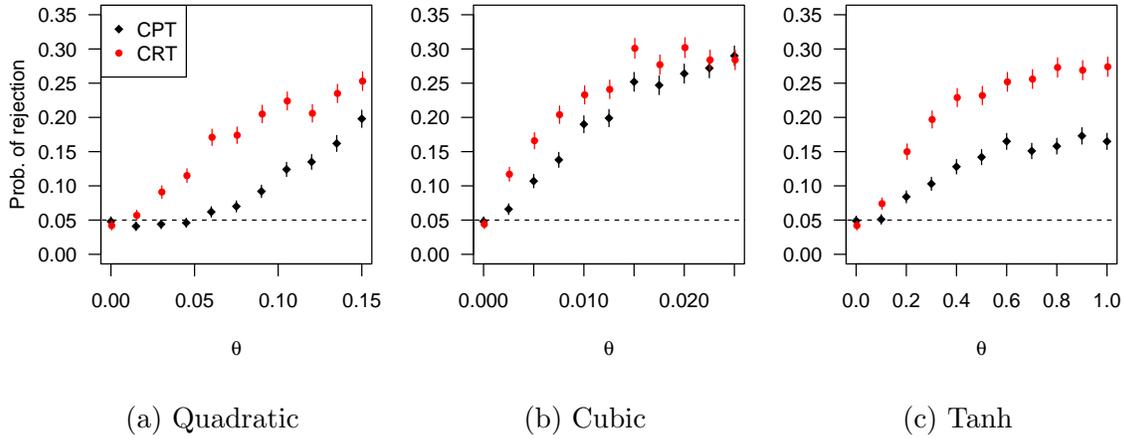


Figure 1: Simulation results for robustness to misspecification of the mean function. The figures show the probability of rejection (i.e., the Type I error rate), plotted against the model misspecification parameter θ . The plots show the average rejection probability with standard error bars computed over 1000 trials for the CPT and CRT. The dashed line indicates the nominal level $\alpha = 0.05$.

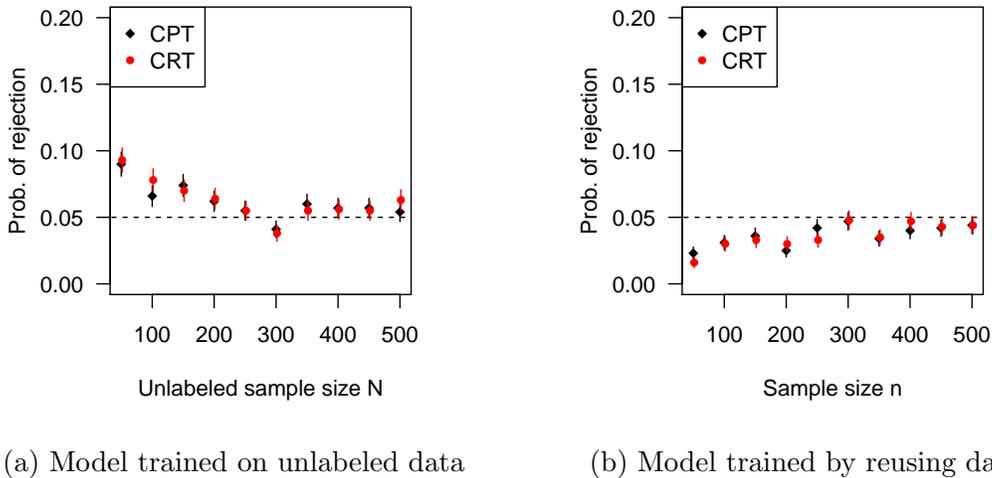


Figure 2: Simulation results for robustness to models trained on unlabeled data or by reusing the data. Details as for Figure 1.

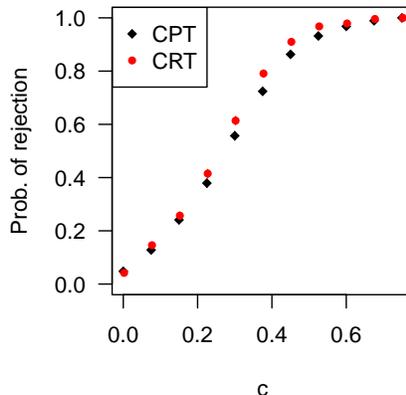


Figure 3: Simulation results testing power under the alternative. The figures show the probability of rejection (i.e., the power), plotted against the signal strength parameter c . The plots show the average rejection probability with standard error bars computed over 1000 trials for the CPT and CRT. The tests are run at level $\alpha = 0.05$.

that both tests have the nominal Type I error level $\alpha = 0.05$. Unlike the null setting, we now generate $Y|X, Z \sim \mathcal{N}(a^\top Z + cX, 1)$. The strength of the signal is controlled by the parameter $c \geq 0$, where $c = 0$ corresponds to the null hypothesis being true while larger values of c move farther away from the null. The results, shown in Figure 3, reveal that the CPT is slightly less powerful than the CRT across a range of values of c , but overall shows fairly similar performance. Thus there is only a small price to pay for the additional robustness of the CPT.

6.2 Simulated data: mixing of the CPT sampler

In practice, we cannot implement the CPT method as defined in (4) (unless, of course, the sample size n is so small that we can simply enumerate all $n!$ possible permutations). Instead, in our experiments, we use the exchangeable MCMC sampler, defined in Algorithm 2. All of our simulations and real data experiments implement this sampler with $S = 50$, meaning that the Markov chain is run for 50 steps for each new permuted copy $\mathbf{X}^{(m)}$ of the data. Is this moderate number of steps sufficient to ensure that the chain has mixed well, or are we producing highly correlated data that will lead to reduced power? To examine this question, we generate one data set, consisting of confounders Z and feature X generated exactly as in Section 6.1.2, and then run the parallel pairwise sampler (Algorithm 1) independently for 20 trials (i.e., each time initializing at the same original data). At each iteration, setting $\mathbf{X}^{[s]} = \mathbf{X}_{(\Pi^{[s]})}$ to be our current CPT copy of the original data vector \mathbf{X} , we track the log-likelihood,

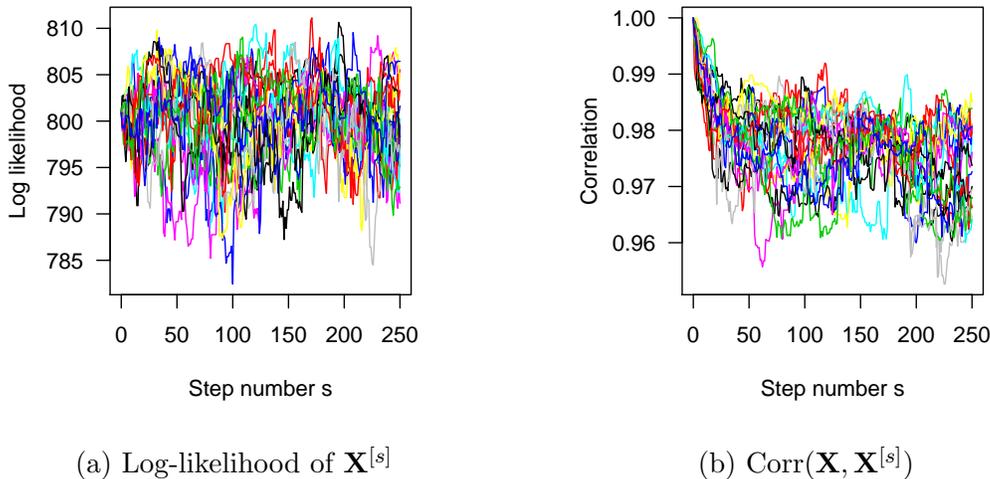


Figure 4: Simulation results showing trace plots for the CPT sampler, examining the CPT copy $\mathbf{X}^{[s]}$ at step s of Algorithm 1.

$\sum_{i=1}^n q(X_i^{[s]}|Z_i)$, and the correlation with the original data vector, $\text{Corr}(\mathbf{X}, \mathbf{X}^{[s]})$. (Note that, since X is strongly dependent with Z , it is to be expected that two draws of the data, i.e., \mathbf{X} and $\mathbf{X}^{[s]}$, will necessarily have a high correlation.) The trace plots of these two quantities, plotted over $s = 0, 1, 2, \dots, 250$ in Figure 4, demonstrate that, in this simulation, the Markov chain appears to mix quickly, within about 50 or 100 iterations. Of course, this will be affected by factors such as the strength of the dependence between X and Z , and the sample size n .

6.3 Capital bikeshare data set

We next implement the CPT and CRT on the Capital Bikeshare data set.⁶ Capital Bikeshare is a bike sharing program in Washington, D.C., where users may check out a bike from one of their locations and return at any other location. The data set contains each ride ever taken, recording the start time and location, end time and location, bike ID number, and a user type which can be “Member” (i.e., purchasing a long-term membership in the system) or “Casual” (i.e., paying for one-time rental or a short-term pass). We use the following data:

- Test data set: all rides taken on weekdays (Monday through Friday) in October 2011. Sample size $n = 7,346$ rides, after an initial screening step (details below).
- Training data set (for fitting the conditional distribution $Q(\cdot|Z)$): all rides taken

⁶Data obtained from <https://www.capitalbikeshare.com/system-data>.

on weekdays in September 2011 and November 2011. Sample size $n_{\text{train}} = 149,912$ rides.

In our experiments, we are interested in determining whether the duration, X , of the ride is dependent on various factors Y , such as user type (“Member” or “Casual”). Of course, the duration of the ride will be heavily dependent on the length of the route, in addition to other factors, and so to control for this we let Z encode both the route, i.e., the start and end locations, as well as the time of day at the start of the ride, since varying traffic might also affect the speed of the ride.

In order to implement the CPT and CRT, we will use a conditional normal distribution, i.e., $(X|Z = z) \sim \mathcal{N}(\mu(z), \sigma^2(z))$ as an estimate $Q(\cdot|z)$ of $Q_\star(\cdot|z)$. Before running the CPT or CRT, as an initial screening step we discard any test points for which we do not have a good estimate of the conditional distribution of X , keeping only those test data points where we have ample training data for rides taken along the same route and at similar times of day. The details for fitting $Q(\cdot|Z)$, and for this initial screening step, are given in Appendix B. For both the CPT and CRT, we sample $M = 1000$ copies of \mathbf{X} to produce the p-value. For the CPT, the Monte Carlo sampler given in Algorithm 2 is run with $S = 50$ as the number of steps for producing each copy.

Results We test the null hypothesis $H_0 : X \perp\!\!\!\perp Y|Z$ for several different choices of the response Y :

- User type (“Member” or “Casual”). We might expect that “Casual” users, who are likely to be tourists or infrequent bike riders, may ride at a slower speed.
- Date, treated as continuous. Since the test data set is taken from the single month October 2011, the date of this month is a continuous variable that acts as a proxy for factors such as weather and the time of sunrise and sunset.
- Day of the week (Monday through Friday), treated as categorical. Bike riders’ behavior may differ on different days of the week, for instance, if rides on Friday are more likely to be leisure rides than the other days of the week.

For user type and date, the statistic $T(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$ that we use is the correlation between the vector \mathbf{Y} , and the vector of ride duration residuals after controlling for the effects of Z —in other words, the vector with entries $R_i = X_i - \mathbb{E}_{X \sim Q(\cdot|Z_i)}[X]$. For day of the week, our statistic $T(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$ is given by

$$\max_{y \in \{\text{Mon}, \dots, \text{Fri}\}} \left| \text{Correlation between } (R_1, \dots, R_n) \text{ and } (\mathbb{1}\{Y_1 = y\}, \dots, \mathbb{1}\{Y_n = y\}) \right|.$$

Table 1 shows the resulting p-values for each choice of the variable Y . We can see that the CPT and CRT produce nearly identical p-values in all three cases. We conclude

Variable Y	CPT p-value (std. err.)	CRT p-value (std. err.)
User type	0.0010 (0.0000)	0.0010 (0.0000)
Date	0.1146 (0.0032)	0.1293 (0.0032)
Day of week	0.1980 (0.0037)	0.2063 (0.0032)

Table 1: p-values obtained from the CPT and CRT for the Capital Bikeshare data. The mean p-value and standard error are calculated from 10 trials of each experiment (the randomness comes from the construction of the copies $\mathbf{X}^{(m)}$ for each test).

that the user type and duration of ride are dependent, even after controlling for our various confounding variables; on the other hand there is insufficient evidence to reach the same conclusion for the corresponding tests for the date and the day of the week.

7 Discussion

In this work, we have developed a conditional permutation test that modifies the standard permutation test of independence between X and Y in order to account for a known dependence of X on potentially relevant confounding variables Z . Our theoretical results prove finite-sample Type I error control, even when the distribution of $X|Z$ is not known exactly.

We have shown that, empirically, resampling from the set of observed X values preserves better Type I error control under mild errors in our model, and does not lose much power, in settings where we use intuitive statistics such as correlation between Y and X . In contrast, our theoretical understanding of Type I error control covers the worst-case scenario over all possible statistics, and it may be the case that the simple statistics used in practical analyses suffer much less inflation of the Type I error. We hope to bridge this gap in future work, and also to provide some theoretical insight into the power of the CPT method, as well as to study the efficiency of the Monte Carlo sampler for the CPT and examine whether proposing swaps non-uniformly may improve the speed at which we can obtain copies $\mathbf{X}^{(m)}$ that are not too correlated with each other.

Furthermore, while in many applications we have access to plenty of unlabeled data, there will certainly be some domains where this is not the case and it may not be possible to estimate the conditional distribution of $X|Z$ independently of the data. If only a small labeled data set (X, Y, Z) is available, with no additional unlabeled data (X, Z) with which to estimate this distribution, we would not want to split the data set to use one half for fitting $Q(X|Z)$ and the remaining half to run the CPT, since this would incur both a substantial loss in the Type I error control (under the theoretical results of Section 5.1), and loss of power when the sample size is limited. It is therefore important to consider how the CPT (and the CRT) can retain their validity when the

data is used for estimating $Q(X|Z)$ and then reused for testing $H_0 : X \perp\!\!\!\perp Y|Z$. It is possible that tools from the selective inference literature may allow us to develop theory towards addressing this question.

Finally, both the CPT and the CRT are based in a setting where it is assumed that modeling $X|Z$ is easy while modeling $Y|X, Z$ is hard—that is, our estimate $Q(\cdot|Z)$ of the conditional distribution $X|Z$ is assumed to be highly accurate, but testing $H_0 : X \perp\!\!\!\perp Y|Z$ is a substantial challenge. In contrast, many of the asymptotic tests described in Section 2.3 treat the X and Y variables symmetrically when testing $X \perp\!\!\!\perp Y|Z$. Are there settings in which we can construct methods offering finite-sample guarantees in the style of the CPT and CRT while taking a more symmetric approach to this testing problem?

A Proofs

A.1 Proving validity of the sampling mechanisms

Proof of Theorem 2. This proof consists of simply checking the detailed balance equations for the Markov chain defined by the algorithm.

Let \mathcal{P} be the set of all partitions of $\{1, \dots, n\}$ into $\lfloor n/2 \rfloor$ disjoint pairs. For any $p \in \mathcal{P}$ and any permutations π, π' , we write $\pi \sim_p \pi'$ if π can be transformed to π' by swapping any subset of the pairs in the partition p . For example, if $(i, j), (k, \ell)$ are two of the disjoint pairs in the partition p , and π and π' are related via $\pi' = \pi \circ \sigma_{ij} \circ \sigma_{k\ell}$, then $\pi \sim_p \pi'$ (recall that σ_{ij} is the permutation that swaps i and j). We note that \sim_p defines an equivalence relation on the set of permutations.

We now compute the transition probability matrix of the Markov chain defined by Algorithm 1. For ease of notation, for the remainder of this proof, we will condition on $\mathbf{X}_(), \mathbf{Y}, \mathbf{Z}$ implicitly. In particular, all probabilities $\mathbb{P}\{\cdot\}$ or $\mathbb{P}\{\cdot|\cdot\}$ should be interpreted as $\mathbb{P}\{\cdot | \mathbf{X}_(), \mathbf{Y}, \mathbf{Z}\}$ or $\mathbb{P}\{\cdot|\cdot, \mathbf{X}_(), \mathbf{Y}, \mathbf{Z}\}$.

For any permutations π, π' , we have

$$\mathbb{P}\{\Pi^{[t]} = \pi' | \Pi^{[t-1]} = \pi\} = \frac{1}{|\mathcal{P}|} \sum_{p \in \mathcal{P}} \mathbb{P}\{\Pi^{[t]} = \pi' | \Pi^{[t-1]} = \pi, t\text{th partition} = p\},$$

since at each time t , Algorithm 1 begins by drawing a partition $p \in \mathcal{P}$ uniformly at random. Next, given p and $\Pi^{[t-1]} = \pi$, $\Pi^{[t]}$ must satisfy $\Pi^{[t]} \sim_p \pi$ by definition of the next step of the algorithm which can only swap pairs of indices in the partition p . By examining the odds ratio defined for each $B_{t,k}$ in (9), we see that for any $\pi', \pi'' \sim_p \pi$,

$$\frac{\mathbb{P}\{\Pi^{[t]} = \pi' | \Pi^{[t-1]} = \pi, t\text{th partition} = p\}}{\mathbb{P}\{\Pi^{[t]} = \pi'' | \Pi^{[t-1]} = \pi, t\text{th partition} = p\}} = \prod_i \frac{q(X_{(\pi'(i))}|Z_i)}{q(X_{(\pi''(i))}|Z_i)} = \frac{\mathbb{P}\{\Pi = \pi'\}}{\mathbb{P}\{\Pi = \pi''\}},$$

where in the last step we refer to the distribution (5) of the permutation Π conditional on $\mathbf{X}_0, \mathbf{Y}, \mathbf{Z}$. Therefore,

$$\mathbb{P}\{\Pi^{[t]} = \pi' \mid \Pi^{[t-1]} = \pi\} = \frac{1}{|\mathcal{P}|} \sum_{p \in \mathcal{P}} \frac{\mathbb{1}\{\pi' \sim_p \pi\} \cdot \mathbb{P}\{\Pi = \pi'\}}{\sum_{\pi''} \mathbb{1}\{\pi'' \sim_p \pi\} \cdot \mathbb{P}\{\Pi = \pi''\}}.$$

Thus, for any π, π' , since \sim_p forms an equivalence relation over permutations, we have

$$\begin{aligned} & \mathbb{P}\{\Pi = \pi\} \cdot \mathbb{P}\{\Pi^{[t]} = \pi' \mid \Pi^{[t-1]} = \pi\} \\ &= \frac{1}{|\mathcal{P}|} \sum_{p \in \mathcal{P}} \mathbb{P}\{\Pi = \pi\} \cdot \frac{\mathbb{1}\{\pi' \sim_p \pi\} \cdot \mathbb{P}\{\Pi = \pi'\}}{\sum_{\pi''} \mathbb{1}\{\pi'' \sim_p \pi\} \cdot \mathbb{P}\{\Pi = \pi''\}} \\ &= \frac{1}{|\mathcal{P}|} \sum_{p \in \mathcal{P}} \mathbb{P}\{\Pi = \pi'\} \cdot \frac{\mathbb{1}\{\pi \sim_p \pi'\} \cdot \mathbb{P}\{\Pi = \pi\}}{\sum_{\pi''} \mathbb{1}\{\pi'' \sim_p \pi'\} \cdot \mathbb{P}\{\Pi = \pi''\}} \\ &= \mathbb{P}\{\Pi = \pi'\} \cdot \mathbb{P}\{\Pi^{[t]} = \pi \mid \Pi^{[t-1]} = \pi'\}. \end{aligned}$$

This verifies the detailed balance equations, and so the Markov chain is reversible and has stationary distribution given by (5). Finally, it is trivial to see that this Markov chain is aperiodic and irreducible when $q(x|z)$ is positive for all $x \in \mathcal{X}$ and $z \in \mathcal{Z}$, and so in this case, the stationary distribution is unique. \square

Proof of Theorem 3. This result follows directly from the fact that the Markov chain defined in Algorithm 1 is reversible, as shown in the proof of Theorem 2. This means that, under H_0 , the permutations $\Pi, \Pi_{\sharp}, \Pi^{(1)}, \dots, \Pi^{(M)}$ can equivalently be drawn as follows: first draw Π_{\sharp} from the distribution (5) conditional on $\mathbf{X}_0, \mathbf{Y}, \mathbf{Z}$, then draw $\Pi, \Pi^{(1)}, \dots, \Pi^{(M)}$ via $M + 1$ independent runs of Algorithm 1 for S steps initialized at $\Pi^{[0]} = \Pi_{\sharp}$. Thus $\Pi, \Pi^{(1)}, \dots, \Pi^{(M)}$ are i.i.d. conditional on $\Pi_{\sharp}, \mathbf{X}_0, \mathbf{Y}, \mathbf{Z}$, and are therefore exchangeable. \square

A.2 Proving robust Type I error control

Proof of Theorem 4. First we prove the result for the CRT. Let $\check{\mathbf{X}}$ be an additional copy drawn also from $Q(\cdot|\mathbf{Z})$, independently of \mathbf{Y} and of $\mathbf{X}, \mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)}$. Then, since conditional on \mathbf{Y}, \mathbf{Z} the copies $\mathbf{X}, \check{\mathbf{X}}, \mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)}$ are independent, we have

$$\begin{aligned} & d_{\text{TV}}\left(\left((\mathbf{X}, \mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)})|\mathbf{Y}, \mathbf{Z}\right), \left((\check{\mathbf{X}}, \mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)})|\mathbf{Y}, \mathbf{Z}\right)\right) \\ &= d_{\text{TV}}\left((\mathbf{X}|\mathbf{Y}, \mathbf{Z}), (\check{\mathbf{X}}|\mathbf{Y}, \mathbf{Z})\right) = d_{\text{TV}}(Q_{\star}^n(\cdot|\mathbf{Z}), Q^n(\cdot|\mathbf{Z})). \end{aligned}$$

Now let $A_{\alpha} \subseteq (\mathcal{X}^n)^{M+1}$ be defined as

$$A_{\alpha} := \left\{ (\mathbf{x}, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)}) : \frac{1 + \sum_{m=1}^M \mathbb{1}\{T(\mathbf{x}^{(m)}, \mathbf{Y}, \mathbf{Z}) \geq T(\mathbf{x}, \mathbf{Y}, \mathbf{Z})\}}{1 + M} \leq \alpha \right\},$$

i.e., the set where we would obtain a p-value $p \leq \alpha$. Then

$$\begin{aligned}
\mathbb{P}\{p \leq \alpha \mid \mathbf{Y}, \mathbf{Z}\} &= \mathbb{P}\{(\mathbf{X}, \mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)}) \in A_\alpha \mid \mathbf{Y}, \mathbf{Z}\} \\
&\leq \mathbb{P}\{(\check{\mathbf{X}}, \mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)}) \in A_\alpha \mid \mathbf{Y}, \mathbf{Z}\} \\
&\quad + d_{\text{TV}}\left(\left((\mathbf{X}, \mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)}) \mid \mathbf{Y}, \mathbf{Z}\right), \left((\check{\mathbf{X}}, \mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)}) \mid \mathbf{Y}, \mathbf{Z}\right)\right) \\
&= \mathbb{P}\{(\check{\mathbf{X}}, \mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)}) \in A_\alpha \mid \mathbf{Y}, \mathbf{Z}\} + d_{\text{TV}}(Q_\star^n(\cdot \mid \mathbf{Z}), Q^n(\cdot \mid \mathbf{Z})).
\end{aligned}$$

Finally, since $\check{\mathbf{X}}, \mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)}$ are clearly i.i.d. after conditioning on \mathbf{Y}, \mathbf{Z} , and are therefore exchangeable, by definition of A_α we must have

$$\mathbb{P}\{(\check{\mathbf{X}}, \mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)}) \in A_\alpha \mid \mathbf{Y}, \mathbf{Z}\} \leq \alpha,$$

proving the desired bound for the CRT.

Next we turn to the CPT, for which the analysis is more complicated since the $\mathbf{X}^{(m)}$'s depend on the observed values in the vector \mathbf{X} . We will use the fact that,

$$\begin{aligned}
\text{For any } (U, V) \text{ and } (U', V'), \text{ if } (V \mid U = u) \stackrel{d}{=} (V' \mid U' = u) \text{ for any } u, \\
\text{then } d_{\text{TV}}\left((U, V), (U', V')\right) = d_{\text{TV}}(U, U').
\end{aligned} \tag{10}$$

Let $\check{\mathbf{X}}$ be drawn from $Q(\cdot \mid \mathbf{Z})$, independently of \mathbf{Y} , and let $\check{\mathbf{X}}^{(1)}, \dots, \check{\mathbf{X}}^{(M)}$ be draws from the CPT when we sample from the values of $\check{\mathbf{X}}$ instead of \mathbf{X} . That is, independently for each $m = 1, \dots, M$, we draw

$$\check{\mathbf{X}}^{(m)} = \check{\mathbf{X}}_{(\check{\Pi}^{(m)})} \quad \text{where } \mathbb{P}\{\check{\Pi}^{(m)} = \pi \mid \check{\mathbf{X}}_0, \mathbf{Y}, \mathbf{Z}\} \propto q^n(\check{\mathbf{X}}_{(\pi)} \mid \mathbf{Z}),$$

where $\check{\mathbf{X}}_0$ and $\check{\mathbf{X}}_{(\pi)}$ are defined analogously to \mathbf{X}_0 and $\mathbf{X}_{(\pi)}$ from Section 3. Next, by comparing to the CPT sampling mechanism (6), we observe that the $\check{\mathbf{X}}^{(m)}$'s, conditional on $\check{\mathbf{X}}$, are generated with the same mechanism as the $\mathbf{X}^{(m)}$'s conditional on \mathbf{X} . In other words, for any $\mathbf{x} \in \mathcal{X}^n$, we have

$$\left((\check{\mathbf{X}}^{(1)}, \dots, \check{\mathbf{X}}^{(M)}) \mid \check{\mathbf{X}} = \mathbf{x}, \mathbf{Y}, \mathbf{Z}\right) \stackrel{d}{=} \left((\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)}) \mid \mathbf{X} = \mathbf{x}, \mathbf{Y}, \mathbf{Z}\right).$$

We can verify that the same equality in distribution holds if we instead use the exchangeable sampler (Algorithm 2) with some choice $S \geq 1$ of the number of steps.

In either case, then, applying (10) we have

$$\begin{aligned}
d_{\text{TV}}\left(\left((\mathbf{X}, \mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)}) \mid \mathbf{Y}, \mathbf{Z}\right), \left((\check{\mathbf{X}}, \check{\mathbf{X}}^{(1)}, \dots, \check{\mathbf{X}}^{(M)}) \mid \mathbf{Y}, \mathbf{Z}\right)\right) \\
= d_{\text{TV}}\left((\mathbf{X} \mid \mathbf{Y}, \mathbf{Z}), (\check{\mathbf{X}} \mid \mathbf{Y}, \mathbf{Z})\right) = d_{\text{TV}}(Q_\star^n(\cdot \mid \mathbf{Z}), Q^n(\cdot \mid \mathbf{Z})).
\end{aligned}$$

From this point on, we proceed as for the CRT—we have

$$\mathbb{P}\{p \leq \alpha \mid \mathbf{Y}, \mathbf{Z}\} \leq \mathbb{P}\{(\tilde{\mathbf{X}}, \tilde{\mathbf{X}}^{(1)}, \dots, \tilde{\mathbf{X}}^{(M)}) \in A_\alpha \mid \mathbf{Y}, \mathbf{Z}\} + d_{\text{TV}}(Q_\star^n(\cdot \mid \mathbf{Z}), Q^n(\cdot \mid \mathbf{Z})),$$

and since $\tilde{\mathbf{X}}, \tilde{\mathbf{X}}^{(1)}, \dots, \tilde{\mathbf{X}}^{(M)}$ are exchangeable after conditioning on \mathbf{Y}, \mathbf{Z} , we see that $\mathbb{P}\{(\tilde{\mathbf{X}}, \tilde{\mathbf{X}}^{(1)}, \dots, \tilde{\mathbf{X}}^{(M)}) \in A_\alpha \mid \mathbf{Y}, \mathbf{Z}\} \leq \alpha$, proving the desired bound for the CPT (with permutations drawn either i.i.d. as in (6), or from the exchangeable sampler given in Algorithm 2). \square

Proof of Theorem 5. For convenience we will write

$$d_{\text{TV}} = d_{\text{TV}}(Q_\star^n(\cdot \mid \mathbf{Z}), Q^n(\cdot \mid \mathbf{Z}))$$

throughout this proof. First, by a standard property of the total variation distance, there exists a subset $A(\mathbf{Z}) \subseteq \mathcal{X}^n$ such that

$$\mathbb{P}_{Q_\star^n(\cdot \mid \mathbf{Z})}\{\mathbf{X} \in A(\mathbf{Z}) \mid \mathbf{Z}\} = \mathbb{P}_{Q^n(\cdot \mid \mathbf{Z})}\{\mathbf{X} \in A(\mathbf{Z}) \mid \mathbf{Z}\} + d_{\text{TV}}.$$

Fix any $M \geq 2$, and define

$$\alpha_0(\mathbf{Z}) := \mathbb{P}_{Q^n(\cdot \mid \mathbf{Z})}\{\mathbf{X} \in A(\mathbf{Z}) \mid \mathbf{Z}\}, \quad \alpha(\mathbf{Z}) := \alpha_0(\mathbf{Z}) + 0.5\sqrt{\frac{\log(M)}{M}}.$$

Now, by definition of the setting and the CRT, we know that conditional on \mathbf{Z} , we have $\mathbf{X} \sim Q_\star^n(\cdot \mid \mathbf{Z})$ and independently, $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)} \sim Q^n(\cdot \mid \mathbf{Z})$. Therefore,

$$\left(\mathbb{1}\{\mathbf{X} \in A(\mathbf{Z})\} \mid \mathbf{Y}, \mathbf{Z}\right) \sim \text{Bernoulli}\left(\alpha_0(\mathbf{Z}) + d_{\text{TV}}\right),$$

and independently,

$$\left(\sum_{m=1}^M \mathbb{1}\{\mathbf{X}^{(m)} \in A(\mathbf{Z})\} \mid \mathbf{Y}, \mathbf{Z}\right) \sim \text{Binomial}(M, \alpha_0(\mathbf{Z})).$$

We will work with the statistic $T(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) = \mathbb{1}\{\mathbf{X} \in A(\mathbf{Z})\}$. We have

$$\begin{aligned} & \mathbb{P}\{p \leq \alpha(\mathbf{Z}) \mid \mathbf{Y}, \mathbf{Z}\} \\ &= \mathbb{P}\left\{\frac{1 + \sum_{m=1}^M \mathbb{1}\{T(\mathbf{X}^{(m)}, \mathbf{Y}, \mathbf{Z}) \geq T(\mathbf{X}, \mathbf{Y}, \mathbf{Z})\}}{1 + M} \leq \alpha(\mathbf{Z}) \mid \mathbf{Y}, \mathbf{Z}\right\} \\ &\geq \mathbb{P}\left\{\mathbf{X} \in A(\mathbf{Z}) \text{ and } \sum_{m=1}^M \mathbb{1}\{\mathbf{X}^{(m)} \in A(\mathbf{Z})\} \leq \alpha(\mathbf{Z}) \cdot (M + 1) - 1 \mid \mathbf{Y}, \mathbf{Z}\right\} \\ &= \left(\alpha_0(\mathbf{Z}) + d_{\text{TV}}\right) \cdot \mathbb{P}\{\text{Binomial}(M, \alpha_0(\mathbf{Z})) \leq \alpha(\mathbf{Z}) \cdot (M + 1) - 1 \mid \mathbf{Z}\} \\ &\geq \alpha(\mathbf{Z}) + d_{\text{TV}} - 0.5\sqrt{\frac{\log(M)}{M}} - \mathbb{P}\{\text{Binomial}(M, \alpha_0(\mathbf{Z})) > \alpha(\mathbf{Z}) \cdot (M + 1) - 1 \mid \mathbf{Z}\}, \end{aligned} \tag{11}$$

where the last step holds by definition of $\alpha(\mathbf{Z})$, $\alpha_0(\mathbf{Z})$, and the fact that $\alpha_0(\mathbf{Z}) + d_{\text{TV}} \leq 1$. Finally, it suffices to bound this binomial probability. By Bennett's inequality, writing $h(u) = (1+u)\log(1+u) - u$, for any $t \in [0, 1]$ we have

$$\begin{aligned}
& \mathbb{P} \left\{ \text{Binomial}(M, t) > \left(t + 0.5\sqrt{\frac{\log(M)}{M}} \right) \cdot (M+1) - 1 \right\} \\
&= \mathbb{P} \left\{ \text{Binomial}(M, t) - Mt > t + 0.5\sqrt{\frac{\log(M)}{M}} \cdot (M+1) - 1 \right\} \\
&\leq \exp \left\{ -Mt(1-t) \cdot h \left(\frac{t + 0.5\sqrt{\frac{\log(M)}{M}} \cdot (M+1) - 1}{Mt(1-t)} \right) \right\} \\
&\leq \exp \left\{ -\frac{M}{4} h \left(\frac{0.5\sqrt{\frac{\log(M)}{M}} \cdot (M+1) - 1}{M/4} \right) \right\}, \tag{12}
\end{aligned}$$

where the last step holds since h is an increasing function, while $c \mapsto c \cdot h(a/c)$ is decreasing in $c > 0$, for any $a > 0$, and $t(1-t) \leq 1/4$.

Finally, as $\epsilon \rightarrow 0$, we have $h(\epsilon) = \epsilon^2/2 + O(\epsilon^3)$, so as $M \rightarrow \infty$ we have

$$\begin{aligned}
\exp \left\{ -\frac{M}{4} h \left(\frac{0.5\sqrt{\frac{\log(M)}{M}} \cdot (M+1) - 1}{M/4} \right) \right\} &= \exp \left\{ -\frac{1}{2} \log(M) + o(1) \right\} \\
&= \frac{1}{\sqrt{M}} = o(1) \cdot 0.5\sqrt{\frac{\log(M)}{M}}.
\end{aligned}$$

Returning to (11), we see that

$$\mathbb{P} \{ p \leq \alpha(\mathbf{Z}) \mid \mathbf{Y}, \mathbf{Z} \} \geq \alpha(\mathbf{Z}) + d_{\text{TV}} - \sqrt{\frac{\log(M)}{M}} \cdot 0.5(1 + o(1)).$$

More concretely, for any $M \geq 2$ we can verify numerically that the quantity in (12) is bounded by $2\sqrt{\frac{\log(M)}{M}}$, which shows that the term $0.5(1 + o(1))$ above can be replaced with 2.5 for any $M \geq 2$. \square

B Details for bikeshare data experiment

We will write $Z = (Z_{\text{route}}, Z_{\text{time}})$, where the route encodes both the start and end locations and is treated as categorical.

To estimate a conditional distribution $Q(\cdot|Z)$, we assume that $X|Z$ is normally distributed, and we fit the conditional mean and variance on the training data by grouping

rides according to their route and taking a Gaussian kernel over their start time: for any $z = (z_{\text{route}}, z_{\text{time}})$,

$$\hat{\mu}(z) = \sum_i \frac{w(z, Z_i^{\text{train}})}{\sum_{i'} w(z, Z_{i'}^{\text{train}})} \cdot X_i^{\text{train}}, \quad \hat{\sigma}^2(z) = \sum_i \frac{w(z, Z_i^{\text{train}})}{\sum_{i'} w(z, Z_{i'}^{\text{train}})} \cdot (X_i^{\text{train}})^2 - (\hat{\mu}(z))^2,$$

where the weights are given by grouping observations by route and applying a Gaussian kernel to the time, i.e.

$$w(z, Z_i^{\text{train}}) = \mathbb{1} \{ (Z_i^{\text{train}})_{\text{route}} = z_{\text{route}} \} \cdot \exp \left\{ - \left((Z_i^{\text{train}})_{\text{time}} - z_{\text{time}} \right)^2 / (2h^2) \right\}$$

for a bandwidth h of 20 minutes. Time of day is on a continuous 24 hour clock, that is, if $z_{\text{time}} = 11:00\text{pm}$ and $(Z_i^{\text{train}})_{\text{time}} = 1:00\text{am}$ then the difference between them is two hours, not 22 hours.

Our conditional distribution estimate $Q(\cdot|Z)$ is then given by

$$(X|Z = z) \sim \mathcal{N}(\hat{\mu}(z), \hat{\sigma}^2(z)).$$

However, since the popularity of various routes and different times of day varies widely, there are some values z where our estimate of the conditional mean and variance of X is unreliable due to scarce data. To check this, for any z we define

$$N(z) = \sum_i w(z, Z_i^{\text{train}}),$$

where a larger $N(z)$ means that there are a larger number of rides in the training data that were taken along the same route z_{route} , and at a time of day similar to z_{time} . For the test data, we then keep only those data points (X_i, Y_i, Z_i) for which $N(Z_i) \geq 20$. Since this screening step uses the value of Z_i but not the value of X_i , the X_i 's are still unobserved even after screening, and their distribution conditional on Z_i is unchanged; therefore the CPT and CRT tests are valid even on this screened data.

Acknowledgements

R.F.B. was partially supported by the NSF via grant DMS-1654076 and by an Alfred P. Sloan fellowship. T.B.B. and R.J.S. were supported by an EPSRC Programme grant. R.J.S. was also supported by an EPSRC Fellowship and a grant from the Leverhulme Trust. The authors would like to thank the Isaac Newton Institute for Mathematical Sciences for its hospitality during the programme Statistical Scalability which was supported by EPSRC Grant Number: LNAG/036, RG91310. The authors thank Samir Khan for help implementing code for our algorithms. We thank the anonymous reviewers for helpful and constructive feedback on an earlier draft.

References

- [1] Susan Athey, Dean Eckles, and W Imbens, Guido. Exact p-values for network inference. *Journal of the American Statistical Association*, 113(521):230–240, 2018.
- [2] Rina Foygel Barber and Emmanuel Candès. On the construction of knockoffs in case–control studies. *Stat*, 8(1):e225, 2019.
- [3] Rina Foygel Barber, Emmanuel J Candès, and Richard J Samworth. Robust inference with knockoffs. *Annals of Statistics*, to appear, 2019.
- [4] Alexandre Belloni, Victor Chernozhukov, and Christian Hansen. Inference on treatment effects after selection among high-dimensional controls. *The Review of Economic Studies*, 81(2):608–650, 2014.
- [5] Wicher Pieter Bergsma. *Testing conditional independence for continuous random variables*. Eurandom, 2004.
- [6] Thomas B Berrett and Richard J Samworth. Nonparametric independence testing via mutual information. *Biometrika*, to appear, 2019.
- [7] Emmanuel Candès, Yingying Fan, Lucas Janson, and Jinchi Lv. Panning for gold: model-X knockoffs for high dimensional controlled variable selection. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 80(3):551–577, 2018.
- [8] Thomas M Cover and Joy A Thomas. *Elements of information theory*. John Wiley & Sons, 2012.
- [9] A Philip Dawid. Conditional independence in statistical theory. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 41(1):1–31, 1979.
- [10] Gary Doran, Krikamol Muandet, Kun Zhang, and Bernhard Schölkopf. A permutation-based kernel conditional independence test. *Uncertainty In Artificial Intelligence*, 30:132–141, 2014.
- [11] Michael D. Ernst. Permutation methods: a basis for exact inference. *Statistical Science*, 19(4):676–685, 2004.
- [12] Kenji Fukumizu, Arthur Gretton, Xiaohai Sun, and Bernhard Schölkopf. Kernel measures of conditional dependence. *Advances in Neural Information Processing Systems*, 20:489–496, 2008.
- [13] Arthur Gretton, Olivier Bousquet, Alex Smola, and Bernhard Schölkopf. Measuring statistical dependence with Hilbert–Schmidt norms. *International Conference on Algorithmic Learning Theory*, 16:63–77, 2005.

- [14] Jonathan Hennessy, Tirthankar Dasgupta, Miratrix Luke, Cassandra Pattanayak, and Pradipta Sarkar. A conditional randomization test to account for covariate imbalance in randomized experiments. *Journal of Causal Inference*, 4(1):61–80, 2016.
- [15] Julie Josse and Susan Holmes. Measures of dependence between random vectors and tests of independence. literature review. *arXiv preprint arXiv:1307.7383*, 2013.
- [16] Ivan Kojadinovic and Mark Holmes. Tests of independence among continuous random vectors based on Cramér–von Mises functionals of the empirical copula process. *Journal of Multivariate Analysis*, 100(6):1137–1154, 2009.
- [17] Niklas Pfister, Peter Bühlmann, Bernhard Schölkopf, and Jonas Peters. Kernel-based tests for joint independence. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 80(1):5–31, 2018.
- [18] Jeffrey Roach and William Valdar. Permutation tests of non-exchangeable null models. *arXiv preprint arXiv:1808.10483*, 2018.
- [19] Paul R Rosenbaum. Conditional permutation tests and the propensity score in observational studies. *Journal of the American Statistical Association*, 79(387):565–574, 1984.
- [20] Jakob Runge. Conditional independence testing based on a nearest-neighbor estimator of conditional mutual information. *Proceedings of the 21st International Conference on Artificial Intelligence and Statistics*, 84:938–947, 2018.
- [21] Rajat Sen, Ananda Theertha Suresh, Karthikeyan Shanmugam, Alexandros G Dimakis, and Sanjay Shakkottai. Model-powered conditional independence test. *Advances in Neural Information Processing Systems*, 31:2955–2965, 2017.
- [22] Rajen D Shah and Jonas Peters. The hardness of conditional independence testing and the generalised covariance measure. *Annals of Statistics*, to appear, 2019.
- [23] Kyungchul Song. Testing conditional independence via Rosenblatt transforms. *The Annals of Statistics*, 37(6B):4011–4045, 2009.
- [24] Stephen M Stigler. Francis Galton’s account of the invention of correlation. *Statistical Science*, 4(2):73–79, 1989.
- [25] Eric V Strobl, Kun Zhang, and Shyam Visweswaran. Approximate kernel-based conditional independence tests for fast non-parametric causal discovery. *Journal of Causal Inference*, to appear, 2019.
- [26] Liangjun Su and Halbert White. A consistent characteristic function-based test for conditional independence. *Journal of Econometrics*, 141(2):807–834, 2007.

- [27] Liangjun Su and Halbert White. A nonparametric Hellinger metric test for conditional independence. *Econometric Theory*, 24(4):829–864, 2008.
- [28] Liangjun Su and Halbert White. Testing conditional independence via empirical likelihood. *Journal of Econometrics*, 182(1):27–44, 2014.
- [29] Gábor J Székely and Maria L Rizzo. Partial distance correlation with methods for dissimilarities. *The Annals of Statistics*, 42(6):2382–2412, 2014.
- [30] Gábor J Székely, Maria L Rizzo, and Nail K Bakirov. Measuring and testing dependence by correlation of distances. *The Annals of Statistics*, 35(6):2769–2794, 2007.
- [31] Noël Veraverbeke, Marek Omelka, and Irène Gijbels. Estimation of a conditional copula and association measures. *Scandinavian Journal of Statistics*, 38(4):766–780, 2011.
- [32] Luca Weihs, Mathias Drton, and Nicolai Meinshausen. Symmetric rank covariances: a generalised framework for nonparametric measures of dependence. *Biometrika*, 105(3):547–562, 2018.
- [33] Kun Zhang, Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. Kernel-based conditional independence test and application in causal discovery. *Uncertainty in Artificial Intelligence*, 27:804–813, 2011.