
A Kernel Test for Three-Variable Interactions with Random Processes (ICML 2016)

Abstract

A wild bootstrap method is applied to the Lancaster three-variable interaction measure in order to detect factorisation of the joint distribution on three variables forming a stationary random process, for which existing permutation bootstrap methods fail. As in the *iid* case, the Lancaster test is found to outperform existing tests in cases for which two independent variables individually have a weak influence on a third, but that when considered jointly the influence is strong. The main contributions of this paper are twofold: first, we show that the Lancaster statistic satisfies the conditions required to use the wild bootstrap; second, the way in which this is proved is novel and is simpler than existing methods, and further may be applied to other statistics.

(An additional minor contribution is that it is also shown that the multiple testing correction proposed in [Lancaster] is too conservative, and a new correction is proposed that increases test power)

1. Introduction

Nonparametric testing of independence or interaction between random variables is a core staple of machine learning and statistics. Many existing methods rely on the assumption that the observed data are drawn *iid*, which for many applications is unrealistic and restrictive. Recent work has begun to extend statistical tests exploiting the theory of Reproducing Kernel Hilbert Spaces (RKHSs) from the *iid* case to the time series case. Of course, these tests also rely on rather restrictive assumptions on the mixing properties of the processes from which the observations are drawn; nonetheless they are a significant relaxation on the *iid* assumption and are a step towards methods capable of handling more general forms of time-dependent data.

The Lancaster interaction is a test statistic capable of de-

tecting dependence between three random variables. The null hypothesis for this test is that the joint distribution of the three variables factorises in some way. In the *iid* case, thresholds for the test statistic are found by using a permutation bootstrap technique: this amounts to shuffling the indices of one or more of the variables and recalculating the test statistic on this bootstrapped data set, providing a valid sample of test statistic under the null hypothesis. When our samples exhibit temporal dependence, shuffling the time indices destroys this dependence and so does not provide a valid sample. The wild bootstrap is a method that has been proposed to correctly resample from the null distribution, subject to certain conditions on the test statistic and the processes from which the observations have been drawn.

In this paper we show that the Lancaster interaction test statistic satisfies the conditions required to apply the wild bootstrap procedure; moreover, we provide a proof that is simpler than existing techniques, and may be adapted to prove similar properties of other kernel statistics. In particular, the proof of a similar property for the HSIC test statistic given in [Kacper] can be substantially simplified, using a version of the Central Limit Theorem for Hilbert space valued random variables instead of the Hoeffding decomposition and theory of U- and V-statistics.

2. Background

In this section we briefly introduce the theory and definitions required to understand the statement and proof of our main result.

2.1. Kernels and RKHS notation

Throughout this paper we will stick to the convention that X, Y and Z are random variables taking value in \mathcal{X}, \mathcal{Y} and \mathcal{Z} , on which we define k, l and m respectively to be kernels. We will assume that our kernels are characteristic and bounded. We describe some notation relevant to the kernel k ; similar notation holds for l and m .

Associated with the kernel k is a Hilbert space \mathcal{H}_k of functions on \mathcal{X} and a feature map $\phi_X : \mathcal{X} \rightarrow \mathcal{H}_k$ such that $k(x, x') = \langle \phi_X(x), \phi_X(x') \rangle$. Given observations $\{X_i\}_{i=1}^n$, we write K to be the *Gram matrix* with entries

$$K_{ij} = k(X_i, X_j).$$

We write $\mu_X := \mathbb{E}_X k(X, \cdot) \in \mathcal{H}_k$ which we call the *mean embedding* of the random variable X . When k is *characteristic*, the mapping from the set of probability distributions to \mathcal{H}_k given by $\mathbb{P}_X \mapsto \mu_X$ is injective. When k is bounded we can think of μ_X as the expectation of the \mathcal{H}_k -valued random variable $\phi_X(X)$. We write $\tilde{\mu}_X = \frac{1}{n} \sum_{i=1}^n k(X_i, \cdot)$ and remark that $\tilde{k}(x, x') = \langle \phi_X(x) - \tilde{\mu}_X, \phi_X(x') - \tilde{\mu}_X \rangle$ is a kernel with feature map $\phi_X(X) = \phi_X(X) - \tilde{\mu}_X$. We denote by \tilde{K} the Gram matrix with respect to \tilde{k} and call this the *empirically centred Gram matrix*. We note also that $\bar{k}(x, x') = \langle \phi_X(x) - \mu_X, \phi_X(x') - \mu_X \rangle$ is a kernel with feature map $\phi_X(X) = \phi_X(X) - \mu_X$. We write $\bar{\mu}_X = \tilde{\mu}_X - \mu_X$, the empirical mean embedding with respect to \tilde{k} . We denote by \bar{K} the Gram matrix with respect to \bar{k} and call this the *population centred Gram matrix*.

If k and l are kernels on \mathcal{X} and \mathcal{Y} , then $k \otimes l$ is a kernel on $\mathcal{X} \times \mathcal{Y}$. We write $C_{XY} = \mathbb{E}_{XY} \bar{\phi}_X(X) \otimes \bar{\phi}_Y(Y)$ called the *population centred covariance operator* and define $\bar{C}_{XY} = \frac{1}{n} \sum_{i=1}^n \bar{\phi}_X(X_i) \otimes \bar{\phi}_Y(Y_i)$ to be its empirical counterpart. Note that we can consider C_{XY} to be an operator $\mathcal{H}_l \rightarrow \mathcal{H}_k$, or as an element of the Hilbert space $\mathcal{H}_{k \otimes l}$. In the latter case we consider it to be the difference of the two mean embeddings $C_{XY} = \mu_{XY} - \mu_X \otimes \mu_Y$.

2.2. Lancaster

The idea of injectively embedding measures into a Hilbert space can be exploited to design statistical tests to test certain properties of the distributions from which observations are drawn. For example, HSIC is an independence test for two random variables with test statistic $\|\tilde{\mu}_{\mathbb{P}_{XY}} - \tilde{\mu}_{\mathbb{P}_X \mathbb{P}_Y}\|^2$, using the fact that $\mu_{\mathbb{P}_{XY}} = \mu_{\mathbb{P}_X \mathbb{P}_Y}$ iff $\mathbb{P}_{XY} = \mathbb{P}_X \mathbb{P}_Y$ to understand the asymptotic properties of the statistic under the null and alternative hypotheses.

The Lancaster test is an extension from the two variable case to consider properties of three variables. The Lancaster statistic on the triple of variables (X, Y, Z) is defined as the signed measure $\Delta_L P = \mathbb{P}_{XYZ} - \mathbb{P}_{XY} \mathbb{P}_Z - \mathbb{P}_{XZ} \mathbb{P}_Y - \mathbb{P}_{YZ} \mathbb{P}_X + 2\mathbb{P}_X \mathbb{P}_Y \mathbb{P}_Z$. It is straightforward to show that if any variable is independent of the other two (equivalently, if the joint distribution \mathbb{P}_{XYZ} factorises into a product of marginals in any way), then $\Delta_L P = 0$. That is, writing $\mathcal{H}_X = \{X \perp (Y, Z)\}$ and similar for \mathcal{H}_Y and \mathcal{H}_Z , we have that

$$\mathcal{H}_X \vee \mathcal{H}_Y \vee \mathcal{H}_Z \Rightarrow \Delta_L P = 0$$

Given a finite sample $(X_i, Y_i, Z_i)_{i=1}^n$, the mean embedding of the Lancaster interaction can be empirically estimated as $\Delta_L \hat{P} = \hat{\mu}_{\mathbb{P}_{XYZ}} - \hat{\mu}_{\mathbb{P}_{XY} \mathbb{P}_Z} - \hat{\mu}_{\mathbb{P}_{XZ} \mathbb{P}_Y} - \hat{\mu}_{\mathbb{P}_X \mathbb{P}_{YZ}} + 2\hat{\mu}_{\mathbb{P}_X \mathbb{P}_Y \mathbb{P}_Z}$. We use the squared RKHS norm of this quan-

tity as a test statistic to test the following hypothesis:

$$\mathcal{H}_0 : \mathcal{H}_X \vee \mathcal{H}_Y \vee \mathcal{H}_Z$$

$$\mathcal{H}_1 : \mathbb{P}_{XYZ} \text{ does not factorise in any way}$$

By [Lancaster], we can write

$$\|\Delta_L \hat{P}\|_{k \otimes l \otimes m}^2 = \frac{1}{n^2} \left(\tilde{K} \circ \tilde{L} \circ \tilde{M} \right)_{++}$$

where \circ is the Hadamard (element-wise) product and $A_{++} = \sum_{ij} A_{ij}$.

The next part of the statistical test is to find threshold values of the statistic beyond which we would reject the null hypothesis. Since our null hypothesis is a composite of three ‘sub-hypotheses’, we must test each of them separately and we reject the composite null hypothesis if and only if we reject all three of the components. In the case that the observations are drawn *iid*, the data can be resampled using permutation bootstrap method. For example, under \mathcal{H}_X , the bootstrapped dataset $\{X_{\pi(i)}, Y_i, Z_i\}_{i=1}^n$ has the same likelihood as the original dataset, and so this can be used to generate a valid resample of the test statistic under \mathcal{H}_X . For more information on the details of the bootstrapping method, see [Lancaster].

2.3. Time series

In this paper we are extending the existing Lancaster test from the *iid* case to a case in which our observations are drawn from a random process. There are various formalisations of memory or ‘mixing’ of a random process; of relevance to this paper are the following two:

Definition 1. A process $(X_t)_t$ is τ -mixing if $\tau(r) \rightarrow 0$ as $r \rightarrow \infty$, where

$$\tau(r) = \sup_{l \in \mathbb{N}} \frac{1}{l} \sup_{r \leq i_1 \leq \dots \leq i_l} \tau(\mathcal{F}_0, (X_{i_1}, \dots, X_{i_l})) \rightarrow 0$$

where

$$\tau(\mathcal{M}, X) = \mathbb{E}(\sup_{g \in \Lambda} \left| \int g(t) \mathbb{P}_{X|\mathcal{M}}(dt) - \int g(t) \mathbb{P}_X(dt) \right|)$$

Definition 2. A process $(X_t)_t$ is β -mixing (also known as absolutely regular) if $\beta(m) \rightarrow 0$ as $m \rightarrow \infty$, where

$$\beta(m) = \frac{1}{2} \sup_n \sup_{i=1}^I \sum_{j=1}^J |\mathbb{P}(A_i \cap B_j) - \mathbb{P}(A_i) \mathbb{P}(B_j)|$$

where the second supremum is taken over all finite partitions $\{A_1, \dots, A_I\}$ and $\{B_1, \dots, B_J\}$ of the sample space such that $A_i \in \mathcal{H}_1^n$ and $B_j \in \mathcal{H}_{n+m}^\infty$ and $\mathcal{H}_b^c = \sigma(X_b, X_{b+1}, \dots, X_c)$

The concept of β -mixing will be invoked when applying a central limit theorem in the next section. We will also need the following lemma:

Lemma 1. Suppose that the process $(X_t, Y_t, Z_t)_t$ is β -mixing. Then any ‘sub-process’ is also β -mixing (for example $(X_t, Y_t)_t$ or $(X_t)_t$)

2.4. V-statistics

A V-statistic of a k -argument, symmetric function f given iid observations $\mathcal{S}_n = \{S_1, \dots, S_n\}$ where each $S_i \sim \mathbb{P}$ is written

$$V(f, \mathcal{S}) = \frac{1}{n^k} \sum_{1 \leq i_1, \dots, i_k \leq n} f(S_{i_1}, \dots, S_{i_k})$$

In this case, $V(f, \mathcal{S})$ is a biased (but asymptotically unbiased) estimator of $\mathbb{E}_{S_{i_1}, \dots, S_{i_k} \sim \mathbb{P}}[f(S_{i_1}, \dots, S_{i_k})]$

In this paper we are only concerned with V-statistics for which $k = 2$. We call $nV(f, \mathcal{S})$ *normalised*. We call f the *core* of V and we say that f is *degenerate* if, for any $s_1, s_2 \sim \mathbb{P}$, $\mathbb{E}[f(s_1, s_2)] = 0$, in which case we say that V is a *degenerate V-statistic*.

Many kernel test statistics can be viewed as normalised V-statistics which, under the null hypothesis, are degenerate. If moreover the test statistics diverge under the alternative hypothesis, the test would be consistent. Our main result is to prove that, under the null hypothesis, the Lancaster statistic is asymptotically a degenerate V-statistic.

2.5. Wild Bootstrap

In many frequentist statistical tests, estimates of the test statistic threshold required to achieve a given Type I error are obtained through a bootstrap resampling method. In the case of the Lancaster and HSIC tests with iid observations, this is done by permuting the time indices of one of the variables to simulate samples from the distribution in which the permuted variable is independent of the other(s). However, this procedure relies on the iid assumption of the data generating process - if, in fact, subsequent samples are *not* independent of previous samples, then permuting the order of the time indices destroys any backward dependence.

If our test statistic has the form of a normalised V-statistic, then provided certain extra conditions are met, the wild bootstrap is a method to directly resample the test statistic under the null hypothesis (in contrast to other methods that first generate a new simulated dataset and then compute the test statistic on this dataset). These conditions can be categorised as concerning: (1) Appropriate τ -mixing of the process from which our observations are drawn; (2) The core of the V-statistic. If these conditions are met by the statistic $nV(f, \mathcal{S}_n)$, then [Wild Bootstrap] tell us that a random matrix W can be drawn such that the bootstrapped statistic $nV_b(f, \mathcal{S}_n) = \frac{1}{n} \sum_{i,j,p,q} W_{ij} f(S_j, S_p) W_{pq}$ is distributed according to the null distribution of nV . The con-

dition on $V(f, \mathcal{S})$ that is of crucial importance to this paper is that f must be a degenerate core.

2.6. Hilbert spaced random variable CLT

In this paper we will exploit a Central Limit Theorem for Hilbert space valued random variables that are functions of random processes. One of the conditions required to apply this theorem concerns appropriate β -mixing of the underlying processes. This theorem is used as a black-box, and it is hoped by the authors that as further theorems concerning CLT-properties of Hilbert space random variables, the conditions required of the processes may be weakened.

3. Lancaster Interaction for Random Processes

(Following Kacper’s paper...)

In this section we construct the Lancaster Interaction test for random processes. The major difficulty in doing so is showing that the test statistic asymptotically satisfies the conditions of the Wild Bootstrap under the null hypothesis of the test, and therefore the Wild Bootstrap can be used to resample the test statistic and provide consistent thresholds for desired p-values.

The approach taken in this paper can also be applied to the HSIC test statistic to give a simpler proof that the Wild Bootstrap can be used for HSIC+timeseries than that given in [Kacper].

Lemma 2. Suppose that (X_i) is β -mixing with coefficients $\beta(m)$ satisfying $\sum_{m=1}^{\infty} \beta(m)^{\frac{\delta}{2+\delta}} < \infty$ for some $\delta > 0$ and that k is a bounded kernel on \mathcal{X} . Then $\|\hat{\mu}_X - \mu_X\|_k = O(n^{-\frac{1}{2}})$

Theorem 1. Suppose that $\mathbb{P}_{XYZ} = \mathbb{P}_{XY}\mathbb{P}_Z$ and that $(X_i, Y_i, Z_i)_{i=1}^n$ are drawn from a process that is both:

- β -mixing with coefficients $\beta(m)$ satisfying $\sum_{m=1}^{\infty} \beta(m)^{\frac{\delta}{2+\delta}} < \infty$
- τ -mixing with coefficients $\tau(m)$ satisfying $\sum_{m=1}^{\infty} m^2 \sqrt{\tau(m)} < \infty$

Then, as $n \rightarrow \infty$,

$$n\|\Delta_L \hat{P}\|^2 \xrightarrow{O(n^{-\frac{1}{2}})} \frac{1}{n} \left((\bar{K} \circ \bar{L}) \circ \bar{M} \right)_{++}$$

and this is a normalised degenerate V-statistic.

Corollary 1. Suppose in addition to the above that W is drawn from a process satisfying the conditions of [Wild Bootstrap]. Then asymptotically,

$$\frac{1}{n} \left(W^\top \left(\overline{(\tilde{K} \circ \tilde{L})} \circ \tilde{M} \right) W \right)_{++}$$

has the same distribution as $n\|\Delta_L \hat{P}\|^2$.

We can therefore use this to generate samples of the test statistic $n\|\Delta_L \hat{P}\|^2$ under the null hypothesis \mathcal{H}_Z . Using these samples we can select a threshold value of the test statistic such that the Type I error is bounded by whatever α we choose. By symmetry, we can use a similar procedure to test \mathcal{H}_X and \mathcal{H}_Y .

Multiple testing correction

In the Lancaster test, we use a composite null hypothesis which requires us to test each of the three hypotheses \mathcal{H}_X , \mathcal{H}_Y and \mathcal{H}_Z separately. We reject the null hypothesis \mathcal{H}_0 if and only if we reject all three of the components. In [Lancaster], it is suggested that the Holm-Bonferroni correction be used to account for multiple testing. We show here that more relaxed conditions on the p-values can be used while still bounding the Type I error, thus increasing test power.

Denote by \mathcal{A}_* the event that \mathcal{H}_* is rejected. Then

$$\begin{aligned} \mathbb{P}(\mathcal{A}_0) &= \mathbb{P}(\mathcal{A}_X \wedge \mathcal{A}_Y \wedge \mathcal{A}_Z) \\ &\leq \min\{\mathbb{P}(\mathcal{A}_X), \mathbb{P}(\mathcal{A}_Y), \mathbb{P}(\mathcal{A}_Z)\} \end{aligned}$$

If \mathcal{H}_0 is true, then so must one of the components. WLOG assume that \mathcal{H}_X is true. If we use significance levels of α in each test individually then $\mathbb{P}(\mathcal{A}_X) \leq \alpha$ and thus $\mathbb{P}(\mathcal{A}_0) \leq \alpha$.

Therefore rejecting \mathcal{H}_0 in the event that each test has p-value less than α individually guarantees a Type I error overall of at most α . In contrast, the Holm-Bonferroni method requires that the sorted p-values be lower than $[\frac{\alpha}{3}, \frac{\alpha}{2}, \alpha]$ in order to reject the null hypothesis overall, is therefore more conservative than necessary and thus loses on test power compared to the ‘simple correction’ proposed here.

4. Experiments

The Lancaster test described above amounts to a method to test each of the sub-hypotheses $\mathcal{H}_X, \mathcal{H}_Y, \mathcal{H}_Z$. Rather than using the Lancaster test statistic with wild bootstrap to test each of these, we could instead use HSIC (it has been previously proved in [Kacper wb], but see supplementary material for a simpler proof that the wild bootstrap can be applied to HSIC). For example, by considering the pair of variables (X, Y) and Z with kernels $k \otimes l$ and m respectively, HSIC can be used to test \mathcal{H}_Z . Similar grouping of

Algorithm 1 Test \mathcal{H}_Z with Wild Bootstrap

Input: $\tilde{K}, \tilde{L}, \tilde{M}$, each size $n \times n$, N = number of bootstraps, α = p-value threshold

$$\|\Delta_L \hat{P}\|^2 = \frac{1}{n} \left(\overline{(\tilde{K} \circ \tilde{L})} \circ \tilde{M} \right)_{++}$$

samples = zeros(1,N)

for $i = 1$ **to** N **do**

 Draw random matrix W according to Wild Bootstrap

$$\text{samples}[i] = \frac{1}{n} \left(W^\top \left(\overline{(\tilde{K} \circ \tilde{L})} \circ \tilde{M} \right) W \right)_{++}$$

end for

if $\text{sum}(\|\Delta_L \hat{P}\|^2 > \text{samples}) > \frac{\alpha}{N}$ **then**

 Reject \mathcal{H}_Z

else

 Do not reject \mathcal{H}_Z

end if

the variables can be used to test \mathcal{H}_X and \mathcal{H}_Y . Applying the same multiple testing correction as in the Lancaster test, we derive an alternative test of dependence between three variables. We refer to this HSIC based procedure as *3-way HSIC*.

In the case of *iid* observations, it was shown in [Lancaster] that Lancaster statistical test is more sensitive to dependence between three random variables than the above HSIC-based test when pairwise interaction is weak but joint interaction is strong. In this section, we demonstrate that the same is true in the time series case on synthetic data.

4.1. Weak pairwise interaction, strong joint interaction

Example 2 in thesis. 3-way HSIC in principle should be able to detect the interaction, but Lancaster is much more powerful. See Figure 1.

Synthetic data were generated from autoregressive processes X, Y and Z according to:

$$\begin{aligned} X_t &= \frac{1}{2}X_{t-1} + \epsilon_t \\ Y_t &= \frac{1}{2}Y_{t-1} + \eta_t \\ Z_t &= \frac{1}{2}Z_{t-1} + d|\theta_t|\text{sign}(X_t Y_t) + \zeta_t \end{aligned}$$

where $X_0, Y_0, Z_0, \epsilon_t, \eta_t, \theta_t$ and ζ_t are *iid* $\mathcal{N}(0, 1)$ random variables and $d \in \mathbb{R}$, called the *dependence* coefficient, determines the extent to which the process $(Z_t)_t$ is dependent on $(X_t, Y_t)_t$.

Data were generated according to this definition with varying values for the dependence coefficient. For each value of the dependence coefficient, 500 datasets were generated, each consisting of 2000 consecutive observations of the variables. Gaussian kernels with bandwidth parameter 1 were used on each variable, and 250 bootstrapping procedures were used for each test on each dataset.

Observe that the random variables are pairwise independent but jointly dependent when considered together. Indeed, X and Y are independent, and Z is independent of X and Y when considered separately since the marginal distributions of X and Y are both normal distributions with mean 0 and therefore $\text{sign}(X_t Y_t)$ is either -1 or 1 with equal probability when conditioned upon neither or exactly one of X_t or Y_t . When considered jointly, the three variables are dependent as knowledge of both X and Y gives information about the value of Z .

In this experiment, all of the tests should be able to detect the dependence and therefore reject the null hypothesis in the limit of infinite data. In the finite data regime, the dependence coefficient affects drastically how hard it is to detect the dependence. The results of this experiment are presented in Figure 1. Observe that the Lancaster test is much more sensitive to the dependence than the 3-way HSIC test is and achieves very high test power with weak dependence coefficients while 3-way HSIC has very low test power. Observe also that when using the simple multiple testing correction a higher test power is achieved than with the Holm-Bonferroni correction.

As before, Lancaster with the ‘naive’ correction outperforms Lancaster with Holm-Bonferroni.

4.2. False positive rates

Example 4 in thesis. Comparison of wild bootstrap to permutation bootstrap

The purpose of this example is to demonstrate that in the time series case, existing permutation bootstrap methods fail to control the Type I error, while the wild bootstrap does correctly find test statistic thresholds.

Synthetic data were generated from autoregressive processes X , Y and Z according to:

$$\begin{aligned} X_t &= aX_{t-1} + \epsilon_t \\ Y_t &= aY_{t-1} + \eta_t \\ Z_t &= aZ_{t-1} + \zeta_t \end{aligned}$$

where $X_0, Y_0, Z_0, \epsilon_t, \eta_t$ and ζ_t are iid $\mathcal{N}(0, 1)$ random variables and a , called the *dependence coefficient*, determines

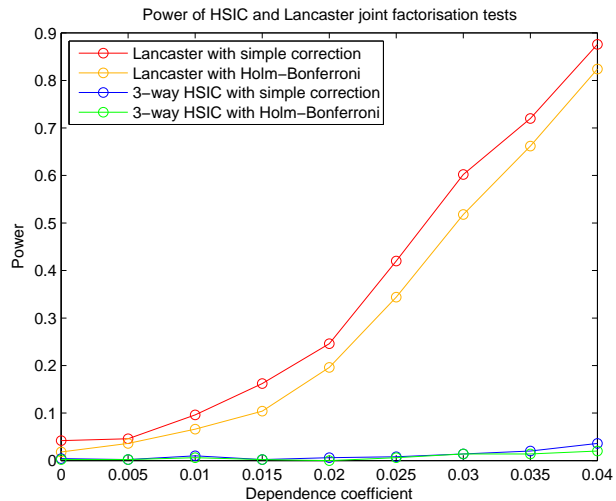


Figure 1. Results of experiment in section 4.1. Observe that the Lancaster test is much more sensitive to dependence than 3-way HSIC. Note also that the test power is higher when using the simple correction when compared to the Holm-Bonferroni multiple testing correction.

how temporally dependent the processes are. Observe that each process is independent of the others and so the null hypothesis is true.

The Lancaster test was performed using both the Wild Bootstrap and the simple permutation bootstrap (used in the *iid* case) in order to sample from the null distributions of the test statistic. We used a fixed desired false positive rate $\alpha = 0.05$ with sample of size 1000, with 500 experiments run for each value of a . Figure 2 shows the false positive rates for these two methods for varying a . It shows that as the processes become more dependent, the false positive rate for the permutation method becomes very large, and is not bounded by the fixed α , whereas the false positive rate for the Wild Bootstrap method is bounded by α .

4.3. Real data

Maybe check this out for some data? https://stat.duke.edu/~mw/ts_data_sets.html

5. Proofs

Proof: (Lemma 2)

We exploit Theorem 1.1 from (?). Using the language of this paper, $\bar{\phi}(X_i)$ is a 1-approximating functional of $(X_i)_i$, following straightforwardly from the definition of 1-approximating functionals given.

Since our kernels are bounded, $\exists C : \|\bar{\phi}(X_i)\| < C$ and

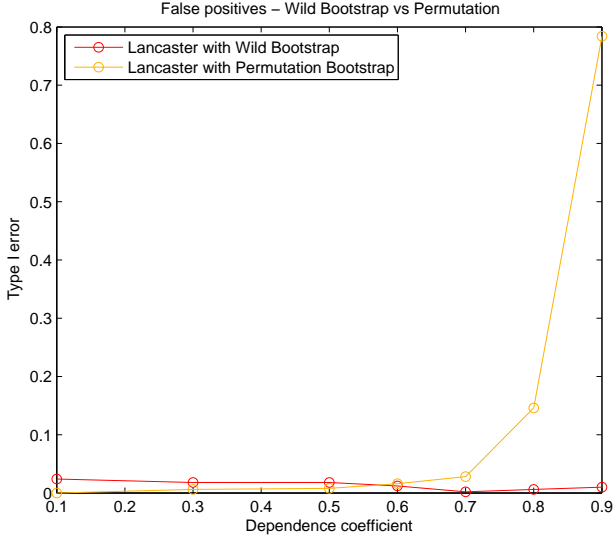


Figure 2. Results of experiment in section 4.2. Whereas the wild bootstrap succeeds in controlling the Type I error across all values of the dependence coefficient, the permutation bootstrap fails to control the Type I error as it does not sample from the correct null distribution as temporal dependence between samples increases.

so

$$\mathbb{E}\|\bar{\phi}(X_1)\|^{2+\delta} < C^{2+\delta} < \infty \quad \forall \delta > 0$$

Thus condition (1) is satisfied.

We can take $f_m = \bar{\phi}(X_0) \quad \forall m$ and so achieve $a_m = 0 \quad \forall m$, thus condition (2) is satisfied.

By assumption on the time series, condition (3) is satisfied.

Thus, by Theorem 1.1 in (?)

$$\sqrt{n}(\bar{\mu}_X - \mu_X) \overset{n \rightarrow \infty}{\rightsquigarrow} N$$

where N is a Hilbert space valued Gaussian random variable. Thus

$$\|\bar{\mu}_X - \mu_X\| = O\left(\frac{1}{\sqrt{n}}\right)$$

■

Proof: (Theorem 1)

By writing

$$\begin{aligned} & \tilde{K}_{ij} \\ &= \langle \phi_X(X_i) - \frac{1}{n} \sum_k \phi_X(X_k), \phi_X(X_j) - \frac{1}{n} \sum_k \phi_X(X_k) \rangle \\ &= \langle \bar{\phi}_X(X_i) - \frac{1}{n} \sum_k \bar{\phi}_X(X_k), \bar{\phi}_X(X_j) - \frac{1}{n} \sum_k \bar{\phi}_X(X_k) \rangle \\ &= \bar{K}_{ij} - \frac{1}{n} \sum_k \bar{K}_{ik} - \frac{1}{n} \sum_k \bar{K}_{jk} + \frac{1}{n^2} \sum_k \bar{K}_{kl} \end{aligned}$$

and expanding \tilde{L} and \tilde{M} in a similar way, we can rewrite the Lancaster test statistic as

$$\begin{aligned} n\|\Delta_L \hat{P}\|^2 &= \frac{1}{n}(\bar{K} \circ \bar{L} \circ \bar{M})_{++} - \frac{2}{n^2}((\bar{K} \circ \bar{L})\bar{M})_{++} \\ &\quad - \frac{2}{n^2}((\bar{K} \circ \bar{M})\bar{L})_{++} - \frac{2}{n^2}((\bar{M} \circ \bar{L})\bar{K})_{++} \\ &\quad + \frac{1}{n^3}(\bar{K} \circ \bar{L})_{++}\bar{M}_{++} + \frac{1}{n^3}(\bar{K} \circ \bar{M})_{++}\bar{L}_{++} \\ &\quad + \frac{1}{n^3}(\bar{L} \circ \bar{M})_{++}\bar{K}_{++} + \frac{2}{n^3}(\bar{M}\bar{K}\bar{L})_{++} \\ &\quad + \frac{2}{n^3}(\bar{K}\bar{L}\bar{M})_{++} + \frac{2}{n^3}(\bar{K}\bar{M}\bar{L})_{++} \\ &\quad + \frac{4}{n^3}\text{tr}(\bar{K}_+ \circ \bar{L}_+ \circ \bar{M}_+) - \frac{4}{n^4}(\bar{K}\bar{L})_{++}\bar{M}_{++} \\ &\quad - \frac{4}{n^4}(\bar{K}\bar{M})_{++}\bar{L}_{++} - \frac{4}{n^4}(\bar{L}\bar{M})_{++}\bar{K}_{++} \\ &\quad + \frac{4}{n^5}\bar{K}_{++}\bar{L}_{++}\bar{M}_{++} \end{aligned}$$

Each of these terms can be expressed as inner products between empirical estimates of population centred covariance operators and tensor products of mean embeddings. Rewriting them as such yields:

$$\begin{aligned} n\|\Delta_L \hat{P}\|^2 &= n\langle \bar{C}_{XYZ}, \bar{C}_{XYZ} \rangle \\ &\quad - 2n\langle \bar{C}_{XYZ}, \bar{C}_{XY} \otimes \bar{\mu}_Z \rangle \\ &\quad - 2n\langle \bar{C}_{XYZ}, \bar{C}_{XZ} \otimes \bar{\mu}_Y \rangle \\ &\quad - 2n\langle \bar{C}_{XYZ}, \bar{C}_{YZ} \otimes \bar{\mu}_X \rangle \\ &\quad + n\langle \bar{C}_{XY} \otimes \bar{\mu}_Z, \bar{C}_{XY} \otimes \bar{\mu}_Z \rangle \\ &\quad + n\langle \bar{C}_{XZ} \otimes \bar{\mu}_Y, \bar{C}_{XZ} \otimes \bar{\mu}_Y \rangle \\ &\quad + n\langle \bar{C}_{YZ} \otimes \bar{\mu}_X, \bar{C}_{YZ} \otimes \bar{\mu}_X \rangle \\ &\quad + 2n\langle \bar{\mu}_Z \otimes \bar{C}_{XY}, \bar{C}_{XZ} \otimes \bar{\mu}_Y \rangle \\ &\quad + 2n\langle \bar{\mu}_X \otimes \bar{C}_{YZ}, \bar{C}_{XY} \otimes \bar{\mu}_Z \rangle \\ &\quad + 2n\langle \bar{\mu}_X \otimes \bar{C}_{ZY}, \bar{C}_{XZ} \otimes \bar{\mu}_Y \rangle \\ &\quad + 4n\langle \bar{C}_{XYZ}, \bar{\mu}_X \otimes \bar{\mu}_Y \otimes \bar{\mu}_Z \rangle \\ &\quad - 4n\langle \bar{C}_{XY} \otimes \bar{\mu}_Z, \bar{\mu}_X \otimes \bar{\mu}_Y \otimes \bar{\mu}_Z \rangle \\ &\quad - 4n\langle \bar{C}_{XZ} \otimes \bar{\mu}_Y, \bar{\mu}_X \otimes \bar{\mu}_Z \otimes \bar{\mu}_Y \rangle \\ &\quad - 4n\langle \bar{C}_{YZ} \otimes \bar{\mu}_X, \bar{\mu}_Y \otimes \bar{\mu}_Z \otimes \bar{\mu}_X \rangle \\ &\quad + 4n\langle \bar{\mu}_X \otimes \bar{\mu}_Y \otimes \bar{\mu}_Z, \bar{\mu}_X \otimes \bar{\mu}_Y \otimes \bar{\mu}_Z \rangle \end{aligned}$$

By assumption, $\mathbb{P}_{XYZ} = \mathbb{P}_{XY}\mathbb{P}_Z$ and thus the expectation operator also factorises similarly. As a consequence,

$$\begin{aligned} C_{XYZ} &= \mathbb{E}_{XYZ}[\bar{\phi}_X(X) \otimes \bar{\phi}_Y(Y) \otimes \bar{\phi}_Z(Z)] \\ &= \mathbb{E}_{XY}[\bar{\phi}_X(X) \otimes \bar{\phi}_Y(Y)] \otimes \mathbb{E}_Z \bar{\phi}_Z(Z) = 0 \end{aligned}$$

Similarly, C_{XZY} , C_{YZX} , C_{XZ} , C_{YZ} are all 0 in their respective Hilbert spaces. Lemma 1 tells us that each subprocess of (X_i, Y_i, Z_i) satisfies the same β -mixing conditions as (X_i, Y_i, Z_i) , thus by applying Lemma 2 to each of the covariance operators at the top of this paragraph we see that each of $\|\bar{C}_{XZY}\|$, $\|\bar{C}_{YZX}\|$, $\|\bar{C}_{XZ}\|$, $\|\bar{C}_{YZ}\|$, $\|\bar{\mu}_X\|$, $\|\bar{\mu}_Y\|$, $\|\bar{\mu}_Z\| = O\left(\frac{1}{\sqrt{n}}\right)$

This can be used to show that

$$\begin{aligned} n\|\Delta_L \hat{P}\|^2 &\rightarrow n\langle \bar{C}_{XYZ}, \bar{C}_{XYZ} \rangle \\ &\quad - 2n\langle \bar{C}_{XYZ}, \bar{C}_{XY} \otimes \bar{\mu}_Z \rangle - 2n\langle \bar{C}_{XZY}, \bar{C}_{XZ} \otimes \bar{\mu}_Y \rangle \\ &= \frac{1}{n}((\bar{K} \circ \bar{L}) \circ \bar{M})_{++} \\ &\quad - \frac{2}{n^2}((\bar{K} \circ \bar{L})\bar{M})_{++} + \frac{1}{n^3}(\bar{K} \circ \bar{L})_{++}\bar{M}_{++} \end{aligned}$$

since all the other terms decay at least as quickly as $O\left(\frac{1}{\sqrt{n}}\right)$ - we show this here for $n\langle \bar{\mu}_X \otimes \bar{C}_{YZ}, \bar{C}_{XY} \otimes \bar{\mu}_Z \rangle$; the proofs for the other terms are similar.

$$\begin{aligned} n\langle \bar{\mu}_X \otimes \bar{C}_{YZ}, \bar{C}_{XY} \otimes \bar{\mu}_Z \rangle &\leq n\|\bar{\mu}_X \otimes \bar{C}_{YZ}\| \|\bar{C}_{XY} \otimes \bar{\mu}_Z\| \\ &= n\sqrt{\langle \bar{\mu}_X \otimes \bar{C}_{YZ}, \bar{\mu}_X \otimes \bar{C}_{YZ} \rangle} \sqrt{\langle \bar{C}_{XY} \otimes \bar{\mu}_Z, \bar{C}_{XY} \otimes \bar{\mu}_Z \rangle} \\ &= n\sqrt{\langle \bar{\mu}_X, \bar{\mu}_X \rangle \langle \bar{C}_{YZ}, \bar{C}_{YZ} \rangle} \sqrt{\langle \bar{C}_{XY}, \bar{C}_{XY} \rangle \langle \bar{\mu}_Z, \bar{\mu}_Z \rangle} \\ &= n\|\bar{\mu}_X\| \|\bar{C}_{YZ}\| \|\bar{C}_{XY}\| \|\bar{\mu}_Z\| \\ &= nO\left(\frac{1}{\sqrt{n}}\right)O\left(\frac{1}{\sqrt{n}}\right)O(1)O\left(\frac{1}{\sqrt{n}}\right) = O\left(\frac{1}{\sqrt{n}}\right) \end{aligned}$$

By treating $\bar{k} \otimes \bar{l}$ as a kernel on the single variable $T := (X, Y)$, we can perform the same recentering trick as before to show that

$$\begin{aligned} n\|\Delta_L \hat{P}\|^2 &\rightarrow \frac{1}{n}((\bar{K} \circ \bar{L}) \circ \bar{M})_{++} \\ &\quad - \frac{2}{n^2}((\bar{K} \circ \bar{L})\bar{M})_{++} + \frac{1}{n^3}(\bar{K} \circ \bar{L})_{++}\bar{M}_{++} \end{aligned}$$

By rewriting the above expression in terms of the operator \bar{C}_{TZ} and mean embeddings μ_T and μ_Z , it can be shown by a similar argument to before that the latter two terms of the above expression tend to 0 at least as $O\left(\frac{1}{n}\right)$, and thus

$$n\|\Delta_L \hat{P}\|^2 \xrightarrow{O\left(\frac{1}{\sqrt{n}}\right)} \frac{1}{n}((\bar{K} \circ \bar{L}) \circ \bar{M})_{++} \text{ as required.}$$

To show that this is a normalised degenerate V-statistic observe that, writing $S_i = (X_i, Y_i, Z_i)$ and $h(S_i, S_j) = \langle \bar{\phi}(X_i) \otimes \bar{\phi}(Y_i) - C_{XY}, \bar{\phi}(X_j) \otimes \bar{\phi}(Y_j) - C_{XY} \rangle \langle \bar{\phi}(Z_i), \bar{\phi}(Z_j) \rangle$, we can write:

$$\frac{1}{n}((\bar{K} \circ \bar{L}) \circ \bar{M})_{++} = \frac{1}{n} \sum_{ij} h(S_i, S_j)$$

And thus it is a normalised V-statistic. To show that it is degenerate, fix any s_i and observe that $\mathbb{E}_{S_j} h(s_i, S_j) = 0$ since $\mathbb{E}_{XYZ} = \mathbb{E}_{XY} \mathbb{E}_Z$. ■

Proof of Lemma 1:

Proof: Let us consider $(X_t, Y_t)_t$. Let us call $\beta_{XYZ}(m)$ the coefficients for the process $(X_t, Y_t, Z_t)_t$, and $\beta_{XY}(m)$ the coefficients for the process $(X_t, Y_t)_t$.

Observe that for $A \in \sigma((X_b, Y_b), \dots, (X_c, Y_c))$, it is the case that $A \times Z \in \sigma((X_b, Y_b, Z_b), \dots, (X_c, Y_c, Z_c))$ and $\mathbb{P}_{XY}(A) = \mathbb{P}_{XYZ}(A \times Z)$.

Thus

$$\begin{aligned} \beta_{XY}(m) &= \frac{1}{2} \sup_n \sup_{\{A_i^{XY}\}, \{B_j^{XY}\}} \sum_{i=1}^I \sum_{j=1}^J |\mathbb{P}_{XY}(A_i^{XY} \cap B_j^{XY}) - \mathbb{P}_{XYZ}(A_i^{XY} \times Z) \mathbb{P}_{XYZ}(B_j^{XY} \times Z)| \\ &= \frac{1}{2} \sup_n \sup_{\{A_i^{XY}\}, \{B_j^{XY}\}} \sum_{i=1}^I \sum_{j=1}^J |\mathbb{P}_{XYZ}((A_i^{XY} \times Z) \cap (B_j^{XY} \times Z)) - \mathbb{P}_{XYZ}(A_i^{XY} \times Z) \mathbb{P}_{XYZ}(B_j^{XY} \times Z)| \\ &\leq \frac{1}{2} \sup_n \sup_{\{A_i^{XYZ}\}, \{B_j^{XYZ}\}} \sum_{i=1}^I \sum_{j=1}^J |\mathbb{P}_{XYZ}(A_i^{XYZ} \cap B_j^{XYZ}) - \mathbb{P}_{XYZ}(A_i^{XYZ}) \mathbb{P}_{XYZ}(B_j^{XYZ})| \\ &= \beta_{XYZ}(m) \end{aligned}$$

Thus we have shown that $\beta_{XYZ}(m) \rightarrow 0 \implies \beta_{XY}(m) \rightarrow 0$. That is, if $(X_t, Y_t, Z_t)_t$ is β -mixing then so is $(X_t, Y_t)_t$

A similar argument holds for any other sub-process. ■

Acknowledgments

cheers!