

---

# A Wild Bootstrap for Degenerate Kernel Tests

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

A wild bootstrap method for a variety of nonparametric hypothesis tests based on kernel distribution embeddings is proposed. This bootstrap results in provably consistent tests that apply to random processes, for which the naive permutation-based bootstrap fails. The wild bootstrap applies to a large group of kernel tests based on V-statistics, which are degenerate under the null hypothesis, and non-degenerate elsewhere. To illustrate this approach, we construct a two-sample test, an instantaneous independence test and a long-range dependence test. In experiments, the wild bootstrap gives strong performance on synthetic examples, on audio data, and in performance benchmarking for the Gibbs sampler.

## 1 Introduction

Statistical tests based on distribution embeddings into reproducing kernel Hilbert spaces have been applied in many contexts, including two sample testing [15, 12, 29], tests of independence [14, 30, 4], tests of conditional independence [11, 30], and tests for higher order (Lancaster) interactions [21]. For these tests, consistency is guaranteed if and only if the observations are independent and identically distributed. Much real-world data fails to satisfy the i.i.d. assumption: audio signals, EEG recordings, text documents, financial time series, and samples obtained when running Markov Chain Monte Carlo, all show significant temporal dependence patterns.

The asymptotic behaviour of kernel test statistics becomes quite different when temporal dependencies exist within the samples. In recent work by on independence testing using the Hilbert-Schmidt Independence Criterion (HSIC) [7], the asymptotic distribution of the statistic under the null hypothesis is obtained for a pair of independent time series, which individually satisfy an absolute regularity or a  $\phi$ -mixing assumption. In this case, the null distribution is shown to be an infinite weighted sum of *dependent*  $\chi^2$ -variables, as opposed to the sum of independent  $\chi^2$ -variables obtained in the i.i.d. setting [14]. The difference in the asymptotic null distributions has important implications in practice: under the i.i.d. assumption, an empirical estimate of the null distribution can be obtained by repeatedly permuting the time indices of one of the signals. This breaks the temporal dependence within the permuted signal, which causes the test to return an elevated number of false positives, when used for testing time series. To address this problem, an alternative estimate of the null distribution is proposed in [7], where the null distribution is simulated by repeatedly *shifting* one signal relative to the other. This preserves the temporal structure within each signal, while breaking the cross-signal dependence.

A serious limitation of the shift procedure in [7] is that it is specific to the problem of independence testing: there is no obvious way to generalise it to other testing contexts. For instance, we might have two time series, with the goal of comparing their marginal distributions - this is a generalization of the two-sample setting to which the shift approach does not apply. We note, however, that many kernel tests have a test statistic with a particular structure: the Maximum Mean Discrepancy (MMD), HSIC, and the Lancaster interaction statistic, each have empirical estimates which can be cast as normalized V-statistics,  $\frac{1}{n^m-1} \sum_{1 \leq i_1, \dots, i_m \leq n} h(Z_{i_1}, \dots, Z_{i_m})$ , where  $Z_{i_1}, \dots, Z_{i_m}$  are draws from a random process at the time points  $\{i_1, \dots, i_m\}$ . Consequently, a method of external random-

ization known as the *wild bootstrap* may be applied [18, 25] to simulate from the null. In brief, the arguments of the above sum are repeatedly multiplied by random, user-defined time series. For a test of level  $\alpha$ , a  $1 - \alpha$  quantile of the empirical distribution of these perturbed statistics serves as the test threshold. This approach has the important advantage over [7] that it may be applied to *all* kernel-based tests for which V-statistics are employed, and not just in the independence setting.

The main result of this paper is to show that the wild bootstrap procedure yields consistent tests for time series in two illustrative problem settings: two-sample testing using MMD, and independence testing using HSIC. Thus, tests based on the wild bootstrap have a Type I error rate (of wrongly rejecting the null hypothesis) approaching the design parameter  $\alpha$ , and a Type II error (of wrongly accepting the null) approaching zero, as the number of samples increases. The procedure will also be used in testing for dependence across a large range of time lags, for which the earlier shift procedure of [7] cannot be applied. This is related to the setting of [4], who test for dependencies which might not occur at a predetermined, single time lag.

We begin our presentation in Section 2, with a review of the  $\tau$ -mixing assumption required of the time series, as well as of generalized V-statistics (of which MMD and HSIC are instances). We also introduce the form taken by the wild bootstrap. In Section 3, we establish a general consistency result for the wild bootstrap procedure on V-statistics, which we apply to MMD and to HSIC in Section 4. Finally, in Section 5, we present a number of empirical comparisons: for the two sample case, we test for differences in audio signals with the same underlying pitch, and present a performance diagnostic for the output of a Gibbs sampler; for the independence case, we test for the dependence of two time series sharing a common variance (a characteristic of FOREX models), and compare against the test of [4] in the case where dependence may occur at multiple, potentially unknown lags. Our tests outperform both the naive approach which neglects the dependence structure within the samples, and the approach of [4] for testing across multiple lags.

## 2 Background

The main results of the paper are based around two concepts:  $\tau$ -mixing [8], which describes the dependence within the time series, and V-statistics [24], which constitute our test statistics. In this section, we review these topics, and introduce the concept of wild bootstrapped V-statistics, which will be the key ingredient in our test construction.

**$\tau$ -mixing.** The notion of  $\tau$ -mixing is used to characterise weak dependence. It is a less restrictive alternative to classical mixing coefficients, and is covered in depth in [8]. Let  $\{Z_t, \mathcal{F}_t\}_{t \in \mathbb{N}}$  be a stationary sequence of integrable random variables, defined on a probability space  $\Omega$  with a probability measure  $P$  and a natural filtration  $\mathcal{F}_t$ . The process is called  $\tau$ -dependent if

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

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

and  $\Lambda$  is the set of all one-Lipschitz real-valued continuous functions on the domain of  $X$ .  $\tau(\mathcal{M}, X)$  can be interpreted as the minimal  $L_1$  distance between  $X$  and  $X^*$  such that  $X \stackrel{d}{=} X^*$  and  $X^*$  is independent of  $\mathcal{M} \subset \mathcal{F}$ . Furthermore, if  $\mathcal{F}$  is rich enough, this  $X^*$  can be constructed (see Proposition 4 in the Appendix).

**V-statistics.** The test statistics considered in this paper are always V-statistics. Given the observations  $Z = \{Z_t\}_{t=1}^n$ , a V-statistic of a symmetric function  $h$  taking  $m$  arguments is given by

$$V(h, Z) = \frac{1}{n^m} \sum_{(i_1, \dots, i_m) \in N^m} h(Z_{i_1}, \dots, Z_{i_m}), \quad (1)$$

where  $N^m$  is a Cartesian power of a set  $N = \{1, \dots, n\}$ . For simplicity, we will often drop the second argument and write simply  $V(h)$ . We will denote the tuple  $(i_1, \dots, i_m)$  by  $i$ .

We will refer to the function  $h$  as to the *core* of the V-statistic  $V(h)$ . While such functions are usually called kernels in the literature, in this paper we reserve the term kernel for positive-definite

functions taking two arguments. A core  $h$  is said to be  $j$ -degenerate if for each  $z_1, \dots, z_j$

$$\mathcal{E}h(z_1, \dots, z_j, Z_{j+1}^*, \dots, Z_m^*) = 0, \quad (2)$$

where  $Z_m^*$  are independent copies of  $Z_0$ . If  $h$  is  $j$ -degenerate for all  $j \leq m-1$ , we will say that it is *canonical*. For a one-degenerate core  $h$ , we define an auxiliary function  $h_2$ , called the second component of the core, and given by

$$h_2(z_1, z_2) = \mathcal{E}h(z_1, z_2, Z_3^*, \dots, Z_{j+1}^*, \dots, Z_m^*). \quad (3)$$

Finally we say that  $nV(h)$  is a normalized  $V$ -statistic, and that a degenerate  $V$ -statistic is a  $V$ -statistic with a one-degenerate core. This degeneracy is common to many kernel statistics when the null hypothesis holds [12, 14, 21].

Our main results will rely on the fact that  $h_2$  governs the asymptotic behaviour of normalized degenerate  $V$ -statistics. Unfortunately, the limiting distribution of such  $V$ -statistics is quite complicated - it is an infinite sum of *dependent*  $\chi^2$ -distributed random variables, with a dependence determined by the temporal dependence structure within the process  $Z$  and by the eigenfunctions of a certain integral operator associated with  $h_2$  [5, 7]. Therefore, we propose a bootstrapped version of the  $V$ -statistics which will allow a consistent approximation of this difficult limiting distribution.

**Bootstrapped V-statistic.** We will study two versions of the bootstrapped  $V$ -statistics

$$V_{b1}(h, Z) = \frac{1}{n^m} \sum_{i \in N^m} W_{i1,n} W_{i2,n} h(Z_{i1}, \dots, Z_{im}), \quad (4)$$

$$V_{b2}(h, Z) = \frac{1}{n^m} \sum_{i \in N^m} \tilde{W}_{i1,n} \tilde{W}_{i2,n} h(Z_{i1}, \dots, Z_{im}), \quad (5)$$

where  $\{W_{t,n}\}_{1 \leq t \leq n}$  is an auxiliary wild bootstrap process and  $\tilde{W}_{t,n} = W_{t,n} - \frac{1}{n} \sum_{j=1}^n W_{j,n}$ . This auxiliary process, proposed by [25, 18], satisfies the following assumption.

*Bootstrap assumption:*  $\{W_{t,n}\}_{1 \leq t \leq n}$  is a row-wise strictly stationary triangular array independent of all  $Z_t$  such that  $\mathcal{E}W_{t,n} = 0$  and  $\sup_n \mathcal{E}|W_{t,n}^{2+\sigma}| < \infty$  for some  $\sigma > 0$ . The autocovariance of the process is given by  $\mathcal{E}W_{s,n}W_{t,n} = \rho(|s-t|/l_n)$  for some function  $\rho$ , such that  $\lim_{u \rightarrow 0} \rho(u) = 1$  and  $\sum_{r=1}^{n-1} \rho(r/l_n) = O(l_n)$ . The sequence  $\{l_n\}$  is taken such that  $l_n = o(n)$  but  $\lim_{n \rightarrow \infty} l_n = \infty$ . The variables  $W_{t,n}$  are  $\tau$ -weakly dependent with coefficients  $\tau(r) \leq C\zeta^{\frac{r}{l_n}}$  for  $r = 1, \dots, n$ ,  $\zeta \in (0, 1)$  and  $C < \infty$ .

As noted in [18, Remark 2], a simple realization of a process that satisfies this assumption is

$$W_{t,n} = e^{-1/l_n} W_{t-1,n} + \sqrt{1 - e^{-2/l_n}} \epsilon_t \quad (6)$$

where  $W_{0,n}$  and  $\epsilon_1, \dots, \epsilon_n$  are independent standard normal random variables. For simplicity, we will drop the index  $n$  and write  $W_t$  instead of  $W_{t,n}$ .

The versions of the bootstrapped  $V$ -statistics in (4) and (5) were previously studied in [18] for the case of canonical cores of degree  $m = 2$ . We extend their results to higher degree cores (common within the kernel testing framework), which are not necessarily one-degenerate. When stating a fact that applies to both  $V_{b1}$  and  $V_{b2}$  we will simply write  $V_b$ , and will drop the argument  $Z$  when there is no ambiguity.

### 3 Asymptotics of wild bootstrapped V-statistics

In this section, we present main Theorems that describe asymptotic behaviour of  $V$ -statistics. In the next section, these results will be used to construct kernel-based statistical tests applicable to dependent observations. Tests are constructed so that the  $V$ -statistic is degenerate under the null hypothesis and non-degenerate under the alternative. Theorem 1 guarantees that the bootstrapped  $V$ -statistic will converge to the same limiting null distribution as the simple  $V$ -statistic. Following [18], since distributions of the bootstrapped statistics are random, we will consider the distributional convergence with the additional qualification “in probability”. This notion can be expressed in terms of convergence in Prokhorov metric  $\varphi$  [10, Section 11.3]. Indeed by [10, Theorem 11.3.3], since values of  $V$ -statistics are real numbers, convergence in distribution is equivalent to the convergence in Prokhorov metric.

**Theorem 1.** Assume that the stationary process  $\{Z_t\}$  is  $\tau$ -dependent with a coefficient  $\tau(i) = O(i^{-6-\epsilon})$  for some  $\epsilon > 0$ . If the core  $h$  is a Lipschitz continuous, one-degenerate, and bounded function of  $m$  arguments and its  $h_2$ -component is a positive definite kernel, then  $\varphi(n^{(m)}_2 V_b(h, Z), nV(h, Z)) \rightarrow 0$  in probability as  $n \rightarrow \infty$ , where  $\varphi$  is Prokhorov metric.

*Proof.* By Lemma 21 and Lemma 20 respectively,  $\varphi(nV_b(h), nV_b(h_2))$  and  $\varphi(nV(h), n^{(m)}_2 V(h_2))$  converge to zero. By [18, Theorem 3.1],  $nV_b(h_2)$  and  $nV(h_2, Z)$  have the same limiting distribution, i.e.,  $\varphi(nV_b(h_2), nV(h_2, Z)) \rightarrow 0$  in probability under certain assumptions. Thus, it suffices to check these assumptions hold: *Assumption A2.* (i)  $h_2$  is one-degenerate and symmetric - this follows from the Lemmas 4 and 3; (ii)  $h_2$  is a kernel - is one of the assumptions of this Theorem; (iii)  $\mathcal{E}h_2(Z_0, Z_0) \leq \infty$  - by Lemma 5,  $h_2$  is bounded and therefore has a finite expected value; (iv)  $h_2$  is Lipschitz continuous - follows from Lemma 5. *Assumption B1.*  $\sum_{i=1}^n i^2 \sqrt{\tau(i)} < \infty$ . Since  $\tau(i) = i^{-6-\epsilon}$  then  $\sum_{i=1}^n i^2 \sqrt{\tau(i)} = \sum_{i=1}^n i^{-1-\epsilon/2} \leq \infty$ . *Assumption B2.* This assumption about the auxiliary process  $\{\tilde{W}_t\}$  is the same as our *Bootstrap assumption*.  $\square$

On the other hand, if the  $V$ -statistic is not degenerate, which is usually true under the alternative, it converges to some non-zero constant. In this setting, Theorem 2 guarantees that the bootstrapped  $V$ -statistic will converge to zero in probability. This property is necessary in testing, as it implies that the test thresholds computed using the bootstrapped  $V$ -statistics will also converge to zero, and so will the corresponding Type II error. The following theorem is due to Lemmas 22 and 23.

**Theorem 2.** Assume that the process  $\{Z_t\}$  is  $\tau$ -dependent with a coefficient  $\tau(i) = O(i^{-6-\epsilon})$ . If the core  $h$  is a Lipschitz continuous, symmetric and bounded function of  $m$  arguments, then  $nV_{b2}(h)$  converges in distribution to some non-zero random variable with finite variance, and  $V_{b1}(h)$  converges to zero in probability.

Although both  $V_{b2}$  and  $V_{b1}$  converge to zero, the rate and the type of convergence are not the same:  $nV_{b2}$  converges in law to some random variable while the behaviour of  $nV_{b1}$  is unspecified. As a consequence, tests that utilize  $V_{b2}$  usually give lower Type II error than the ones that use  $V_{b1}$ . On the other hand,  $V_{b1}$  seems to better approximate  $V$ -statistic distribution under the null hypothesis. This agrees with our experiments in Section 5 as well as with those in [18, Section 5]).

## 4 Applications to Kernel Tests

In this section, we describe how the wild bootstrap for  $V$ -statistics can be used to construct kernel tests for the two-sample problem and for independence, which are applicable to weakly dependent observations. We start by reviewing the main concepts underpinning the kernel testing framework.

For every symmetric, positive definite function, i.e., *kernel*  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ , there is an associated reproducing kernel Hilbert space  $\mathcal{H}_k$  [3, p. 19]. The kernel embedding of a probability measure  $P$  on  $\mathcal{X}$  is an element  $\mu_k(P) \in \mathcal{H}_k$ , given by  $\mu_k(P) = \int k(\cdot, x) dP(x)$  [3, 26]. If a measurable kernel  $k$  is bounded, the mean embedding  $\mu_k(P)$  exists for all probability measures on  $\mathcal{X}$ , and for many interesting bounded kernels  $k$ , including the Gaussian, Laplacian and inverse multi-quadratics, the kernel embedding  $P \mapsto \mu_k(P)$  is injective. Such kernels are said to be *characteristic* [28]. The RKHS-distance  $\|\mu_k(P_x) - \mu_k(P_y)\|_{\mathcal{H}_k}^2$  between embeddings of two probability measures  $P_x$  and  $P_y$  is termed the Maximum Mean Discrepancy (MMD), and its empirical version serves as a popular statistic for non-parametric two-sample testing [12]. Similarly, given a sample of paired observations  $\{(X_i, Y_i)\}_{i=1}^n \sim P_{xy}$ , and kernels  $k$  and  $l$  respectively on  $X$  and  $Y$  domains, the RKHS-distance  $\|\mu_\kappa(P_{xy}) - \mu_\kappa(P_x P_y)\|_{\mathcal{H}_\kappa}^2$  between embeddings of the joint distribution and of the product of the marginals measures dependence between  $X$  and  $Y$ . Here,  $\kappa((x, y), (x', y')) = k(x, x')l(y, y')$  is the kernel on the product space of  $X$  and  $Y$  domains. This quantity is called Hilbert-Schmidt Independence Criterion (HSIC) [13, 14]. When characteristic RKHSs are used, the HSIC is zero iff the variables are independent: this follows from [19, Lemma 3.8] and [27, Proposition 2]. The empirical statistic is written  $\widehat{\text{HSIC}}_k = \frac{1}{n^2} \text{Tr}(KHLH)$  for kernel matrices  $K$  and  $L$  and the centering matrix  $H = I - \frac{1}{n} \mathbf{1}\mathbf{1}^\top$ .

### 4.1 Wild Bootstrap For MMD

Denote the observations by  $\{X_i\}_{i=1}^{n_x} \sim P_x$ , and  $\{Y_j\}_{j=1}^{n_y} \sim P_y$ . Our goal is to test the null hypothesis  $\mathbf{H}_0 : P_x = P_y$  vs. the alternative  $\mathbf{H}_1 : P_x \neq P_y$ . In the case where samples have

equal sizes, i.e.,  $n_x = n_y$ , application of the wild bootstrap from [18] to MMD-based tests on dependent samples is straightforward: the empirical MMD can be written as a V-statistic with the core of degree two on pairs  $z_i = (x_i, y_i)$  given by  $h(z_1, z_2) = k(x_1, x_2) - k(x_1, y_2) - k(x_2, y_1) + k(y_1, y_2)$ . It is clear that whenever  $k$  is Lipschitz continuous and bounded, so is  $h$ . Moreover,  $h$  is a valid positive definite kernel, since it can be represented as an RKHS inner product  $\langle k(\cdot, x_1) - k(\cdot, y_1), k(\cdot, x_2) - k(\cdot, y_2) \rangle_{\mathcal{H}_k}$ . Under the null hypothesis,  $h$  is also one-degenerate, i.e.,  $\mathcal{E}h((x_1, y_1), (X_2, Y_2)) = 0$ . Therefore, we can use the bootstrapped statistics in (4) and (5) to approximate the null distribution and attain a desired test level.

When  $n_x \neq n_y$ , however, it is no longer possible to write the empirical MMD as a one-sample V-statistic. We will therefore require the following bootstrapped version of MMD

$$\begin{aligned} \widehat{\text{MMD}}_{k,b} = & \frac{1}{n_x^2} \sum_{i=1}^{n_x} \sum_{j=1}^{n_x} \tilde{W}_i^{(x)} \tilde{W}_j^{(x)} k(x_i, x_j) - \frac{1}{n_x^2} \sum_{i=1}^{n_y} \sum_{j=1}^{n_y} \tilde{W}_i^{(y)} \tilde{W}_j^{(y)} k(y_i, y_j) \\ & - \frac{2}{n_x n_y} \sum_{i=1}^{n_x} \sum_{j=1}^{n_y} \tilde{W}_i^{(x)} \tilde{W}_j^{(y)} k(x_i, y_j), \end{aligned} \quad (7)$$

where  $\tilde{W}_t^{(x)} = W_t^{(x)} - \frac{1}{n_x} \sum_{i=1}^{n_x} W_i^{(x)}$ ,  $\tilde{W}_t^{(y)} = W_t^{(y)} - \frac{1}{n_y} \sum_{j=1}^{n_y} W_j^{(y)}$ ;  $\{W_t^{(x)}\}$  and  $\{W_t^{(y)}\}$  are two auxiliary wild bootstrap processes that are independent of  $\{X_t\}$  and  $\{Y_t\}$  and also independent of each other, both satisfying the bootstrap assumption in Section 2. The following Proposition shows that the bootstrapped statistic has the same asymptotic null distribution as the empirical MMD. The proof follows that of [18, Theorem 3.1], and is given in the Appendix.

**Proposition 1.** *Let  $k$  be bounded and Lipschitz continuous, and let  $\{X_t\}$  and  $\{Y_t\}$  both be  $\tau$ -dependent with coefficients  $\tau(i) = O(i^{-6-\epsilon})$ , but independent of each other. Further, let  $n_x = \rho_x n$  and  $n_y = \rho_y n$  where  $n = n_x + n_y$ . Then, under the null hypothesis  $P_x = P_y$ ,  $\varphi\left(\rho_x \rho_y n \widehat{\text{MMD}}_{k,b}, \rho_x \rho_y n \widehat{\text{MMD}}_{k,b}\right) \rightarrow 0$  in probability as  $n \rightarrow \infty$ , where  $\varphi$  is the Prokhorov metric.*

## 4.2 Wild Bootstrap For HSIC

Under the null hypothesis, the empirical HSIC is a degenerate V-statistic, and its distribution is an infinite sum of independent  $\chi^2$ -variables [14]. Should  $Z_t = (X_t, Y_t)$  be a random process rather than a sequence of i.i.d. random variables, the limiting distribution of the empirical HSIC becomes more convoluted, as it becomes a function of the temporal dependence within the process  $\{Z_t\}$ . Convergence in probability of empirical HSIC to its population value was shown in [31] for a 1-approximating functional of an absolutely regular process [6]. No asymptotic distributions were obtained, however, nor was a statistical test constructed. The asymptotics of a normalized V-statistic were obtained in [7] for absolutely regular and  $\phi$ -mixing<sup>1</sup> processes [9]. Due to the intractability of the null distribution for the test statistic, the authors propose a procedure to approximate its null distribution using circular shifts of the observations leading to tests of instantaneous independence, i.e., of  $X_t \perp\!\!\!\perp Y_t, \forall t$ . This was shown to be consistent under the null (i.e., leading to the correct Type I error), however consistency of the shift procedure under the alternative is a challenging open question (see [?, Section A.2]chwi2014kernel for further discussion). In contrast, as shown below in Propositions 2 and 3 (which are direct consequences of the Theorems 1 and 2), the wild bootstrap guarantees test consistency under both hypotheses: null and alternative, and is a major advantage of the present approach. In addition, the wild bootstrap can be used in constructing a test for the harder problem of determining independence across multiple lags simultaneously, similar to the one in [4].

Following symmetrisation, it can be shown that the empirical HSIC can be written as a degree four V-statistic with core given by

$$h(z_1, z_2, z_3, z_4) = \frac{1}{4!} \sum_{\pi \in S_4} k(x_{\pi(1)}, x_{\pi(2)}) [l(y_{\pi(1)}, y_{\pi(2)}) + l(y_{\pi(3)}, y_{\pi(4)}) - 2l(y_{\pi(2)}, y_{\pi(3)})],$$

where we denote by  $S_n$  the group of permutations over  $n$  elements. Thus, we can directly apply the theory developed for higher-order V-statistics in Section 3. We consider two types of tests: instantaneous independence and independence at all time lags.

<sup>1</sup>The relation between different mixing coefficients is discussed in [8].

Table 1: Rejection rates for two-sample experiments. **MCMC**: sample size=500; a gaussian kernel with bandwidth  $\sigma = 1.7$  is used; every second Gibbs sample is kept (i.e., after a pass through both dimensions). **Pitch**: a gaussian kernel with bandwidth  $\sigma = 14$  is used. Both: wild bootstrap uses blocksize of  $l_n = 20$ ; averaged over at least 200 trials.

experiment \ method	vanilla	$\widehat{\text{MMD}}_{k,b}$	$V_{b1}$	$V_{b2}$
<b>MCMC</b> : i.i.d. vs i.i.d. ( $\mathbf{H}_0$ )	.040	.025	.012	.070
i.i.d. vs Gibbs ( $\mathbf{H}_0$ )	.528	.100	.052	.105
Gibbs vs Gibbs ( $\mathbf{H}_0$ )	.680	.110	.060	.100
<b>Pitch</b> : $n_x = 300, n_y = 200$ ( $\mathbf{H}_0$ )	.970	.145		
$n_x = 600, n_y = 400$ ( $\mathbf{H}_0$ )	.965	.104		
$n_x = 300, n_y = 200$ ( $\mathbf{H}_A$ )	1	.600		
$n_x = 600, n_y = 400$ ( $\mathbf{H}_A$ )	1	.898		

**Test of instantaneous independence** Here, the null hypothesis  $\mathbf{H}_0$  is that  $X_t$  and  $Y_t$  are independent at all times  $t$ , and the alternative hypothesis  $\mathbf{H}_1$  is that they are dependent.

**Proposition 2.** *Under the null hypothesis, if the stationary process  $Z_t = (X_t, Y_t)$  is  $\tau$ -dependent with a coefficient  $\tau(i) = i^{-6-\epsilon}$  for some  $\epsilon > 0$ , then  $\varphi(6nV_b(h), nV(h)) \rightarrow 0$  in probability, where  $\varphi$  is the Prokhorov metric.*

*Proof.* Since both  $k$  and  $l$  are bounded and Lipschitz continuous, the core  $h$  is bounded and Lipschitz continuous. One-degeneracy under the null hypothesis was stated in [14, Theorem 2] and the fact that  $h_2$  is a kernel was shown in [14, section A.2, following eq. 11]. The result then follows from Theorem 1.  $\square$

The following proposition holds by Theorem 2, since the core  $h$  is Lipschitz continuous, symmetric and bounded.

**Proposition 3.** *If the stationary process  $Z_t$  is  $\tau$ -dependent with a coefficient  $\tau(i) = i^{-6-\epsilon}$  for some  $\epsilon > 0$ , then under the alternative hypothesis  $nV_{b2}(h)$  converges in distribution to some random variable with a finite variance and  $V_{b1}$  converges to zero in probability.*

**Lag-HSIC** Propositions 2 and 3 also allow us to construct a test of time series independence that is similar to one designed by [4]. Here, we will be testing against a broader null hypothesis:  $X_t$  and  $Y_{t'}$  are independent for  $|t - t'| < M$  for an arbitrary large but fixed  $M$ . In the Appendix, we show how to construct a test when  $M \rightarrow \infty$ , although this requires an additional assumption about the uniform convergence of cumulative distribution functions.

Since the time series  $Z_t = (X_t, Y_t)$  is stationary, it suffices to check whether there exists a dependency between  $X_t$  and  $Y_{t+m}$  for  $-M \leq m \leq M$ . Since each lag corresponds to an individual hypothesis, we will require a Bonferroni correction to attain a desired test level  $\alpha$ . We therefore define  $q = 1 - \frac{\alpha}{2M+1}$ . The shifted time series will be denoted  $Z_t^m = (X_t, Y_{t+m})$ . Let  $S_{m,n} = nV(h, Z_t^m)$  denote the value of the normalized HSIC statistic calculated on the shifted process  $Z_t^m$ . Let  $F_{b,n}$  denote the empirical cumulative distribution function obtained by the bootstrap procedure using  $nV_b(h, Z)$ . The test will then reject the null hypothesis if the event  $\mathcal{A}_n = \left\{ \max_{-M \leq m \leq M} S_{m,n} > F_{b,n}^{-1}(q) \right\}$  occurs. By a simple application of the union bound, it is clear that the asymptotic probability of the Type I error will be  $\lim_{n \rightarrow \infty} P_{\mathbf{H}_0}(\mathcal{A}_n) \leq \alpha$ . On the other hand, if the alternative holds, there exists some  $m$  with  $|m| \leq M$  for which  $V(h, Z^m) = n^{-1}S_{m,n}$  converges to a non-zero constant. In this case

$$P_{\mathbf{H}_1}(\mathcal{A}_n) \geq P_{\mathbf{H}_1}(S_{m,n} > F_{b,n}^{-1}(q)) = P_{\mathbf{H}_1}(n^{-1}S_{m,n} > n^{-1}F_{b,n}^{-1}(q)) \rightarrow 1 \quad (8)$$

as long as  $n^{-1}F_{b,n}^{-1}(q) \rightarrow 0$ , which follows from the convergence of  $V_b$  zero in probability shown in Proposition 3. Therefore, the Type II error of the multiple lag test with  $V_{b2}$  is guaranteed to converge to zero as the sample size increases. Our experiments in the next Section demonstrate that while this procedure is defined over a finite range of lags, it results in tests more powerful than the procedure for an infinite number of lags proposed in [4]. We note that a procedure that works for an infinite

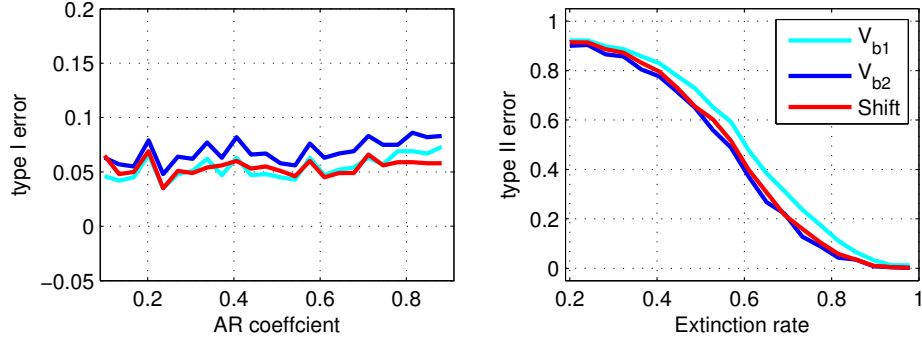


Figure 1: Comparison of Shift-HSIC and tests based on  $V_{b1}$  and  $V_{b2}$ . Note that the test based on  $V_{b2}$  has slightly higher type I error and that the test based on  $V_{b1}$  has slightly higher type II error.

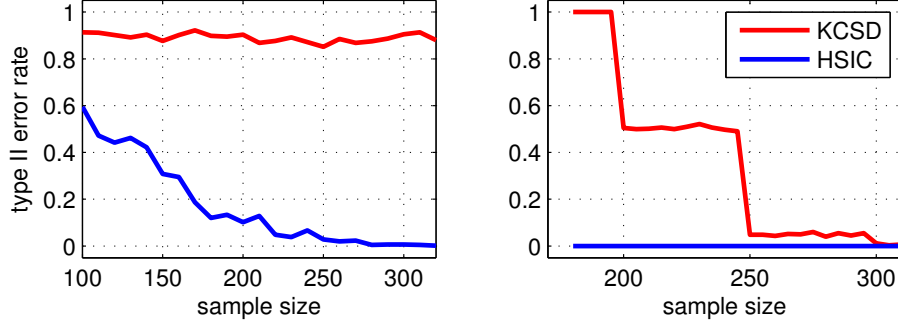


Figure 2: The left panel presents type I and II error of the lag-HSIC and KCSD algorithms for process following dynamics given by the equation (9), whereas the errors for process with dynamics given by equation (10) and (11) are shown in the right panel. X axis is indexed by the time series length i.e. sample size.

number of lags, although possible to construct, does not add much practical value under the present assumptions. Indeed, since the  $\tau$ -mixing assumption applies to the joint sequence  $Z_t = (X_t, Y_t)$ , dependence between  $X_t$  and  $Y_{t+m}$  is bound to disappear at a rate of  $o(m^{-6})$ , i.e., the variables both within and across the two series are assumed to become gradually independent.

## 5 Experiments

**The MCMC M.D.** It is natural to use MMD in order to measure how far the MCMC chain is from its stationary distribution [23, Section 5], by comparing the produced sample to a benchmark sample (e.g., one produced with a heavily thinned random walk, which is wasteful of samples and computationally burdensome). It is not clear how to construct a statistical test on whether the sampler has converged, however, due to the chain dependence—standard approximations of null distributions lead to too many rejections of the null hypothesis. Our experiments indicate that the wild bootstrap approach allows such testing as it attains a desired number of false positives.

To assess performance of the wild bootstrap in MCMC convergence diagnostics, we consider the situation where samples  $\{X_i\}$  and  $\{Y_i\}$  are bivariate and both have the identical marginal distribution given by an elongated normal  $P = \mathcal{N}\left(\begin{bmatrix} 0 & 0 \end{bmatrix}, \begin{bmatrix} 15.5 & 14.5 \\ 14.5 & 15.5 \end{bmatrix}\right)$ . However, they could have arisen either as independent samples or as outputs of the Gibbs sampler with stationary distribution  $P$ . Table 1 shows the Type I error under the significance level  $\alpha = 0.05$ . It is clear that in the case where at least one of the samples is a Gibbs chain, the permutation test has a Type I error much larger than  $\alpha$ , since the null distribution is misspecified. The wild bootstrap using  $V_{b1}$  (without artificial degeneration) yields the correct Type I error control in these cases. Consistent with findings in [18, Section 5],  $V_{b1}$  mimics the null distribution better than  $V_{b2}$ . Bootstrapped statistic  $\widehat{\text{MMD}}_{k,b}$  in

(7) with two auxiliary artificially degenerated bootstrap processes behaves similarly to  $V_{b2}$ . In the alternative scenario where  $\{Y_i\}$  was taken from a distribution with the same covariance structure but with the mean set to  $\mu = [2.5 \ 0]$ , Type II error for all tests was zero.

**Pitch-evoking sounds** The second experiment we consider is a two sample test on sounds studied in the field of pitch perception [16]. We synthesise the sounds with the fundamental frequency parameter of treble C, subsampled at 10.46kHz. Each  $i$ -th period of length  $\Omega$  contains  $d = 20$  audio samples at times  $0 = t_1 < \dots < t_d < \Omega$  – we treat this whole vector as a single observation  $X_i$  or  $Y_i$ , i.e., we are comparing distributions on  $\mathbb{R}^{20}$ . Sounds are generated based on the AR process  $a_i = \lambda a_{i-1} + \sqrt{1 - \lambda^2} \epsilon_i$ , where  $a_0, \epsilon_i \sim \mathcal{N}(0, I_d)$ , with  $X_{i,r} = \sum_j \sum_{s=1}^d a_{j,s} \exp\left(-\frac{(t_r - t_s - (j-i)\Omega)^2}{2\sigma^2}\right)$ , i.e., a given pattern – a smoothed version of  $a_0$  – slowly varies and thus the sound deviates from periodicity, but still evokes a pitch. We take  $X$  with  $\sigma = 0.1\Omega$  and  $\lambda = 0.8$ , and  $Y$  is either an independent copy of  $X$  (null scenario), or has  $\sigma = 0.05\Omega$  (alternative scenario)<sup>2</sup>.  $n_x$  is taken to be different from  $n_y$ . Results in Table 1 demonstrate how the approach using the wild bootstrapped statistic in (7) allows to control the Type I error while the permutation test virtually always rejects the null hypothesis.

**Instantaneous independence** To examine instantaneous independence test performance we compare it with the Shift-HSIC procedure [7] on the ‘Extinct Gaussian’ autoregressive process proposed in the [7, Section 4.1]. Using exactly the same setting we compute type I error as a function of the temporal dependence and type II error as a function of extinction rate<sup>3</sup>. The Figure 1 shows that all three tests, Shift-HSIC and tests based on  $V_{b1}$  and  $V_{b2}$ , perform similarly.

**Lag-HSIC** The KCSD test [4] is, to our knowledge, the only algorithm to conduct a long-range independence test for time series. In the experiments, we compare lag-HSIC with KCSD on two kinds of processes: one inspired by econometrics and one from [4].

In lag-HSIC, the amount of lags under examination was equal to  $\max\{10, \log n\}$ , where  $n$  is the sample size. We have used gaussian kernels with widths estimated by the median heuristic. The cumulative distribution of the  $V$ -statistics was approximated by samples from  $nV_{b2}$ . To model the tail of this distribution, we have fitted the generalized Pareto distribution to the bootstrapped samples ([20] shows that for a large class of underlying distribution functions such an approximation is valid). The first process is a pair of two time series which share a common variance,

$$X_t = \epsilon_{1,k} \sigma_k^2, \quad Y_t = \epsilon_{2,k} \sigma_k^2, \quad \sigma_k^2 = 1 + 0.45(X_t^2 + Y_t^2). \quad (9)$$

The above set of equations is an instance of the VEC dynamics [2] used in econometrics to model market volatility. The left panel of the Figure 2 presents the Type II error rate: for KCSD it remains at 90% while for lag-HSIC it gradually drops to zero. The Type I error, which we calculated by sampling two independent copies  $(X_t^{(1)}, Y_t^{(1)})$  and  $(X_t^{(2)}, Y_t^{(2)})$  of the process and performing the tests on the pair  $(X_t^{(1)}, Y_t^{(2)})$ , was around 5% for both of the tests.

Our next experiment is a process sampled according to the dynamics proposed by [4],

$$X_k = \cos(\phi_{k,1}) \quad \phi_{k,1} = \phi_{k-1,1} + 0.1\epsilon_{1,k} + 2\pi f_1 T_s \quad (10)$$

$$Y_k = [2 + C \sin(\phi_{k,1})] \cos(\phi_{k,2}) \quad \phi_{k,2} = \phi_{k-1,2} + 0.1\epsilon_{2,k} + 2\pi f_2 T_s \quad (11)$$

with parameters  $C = .4$ ,  $f_1 = 4Hz$ ,  $f_2 = 20Hz$ , and frequency  $\frac{1}{T_s} = 100Hz$ . We compared performance of the KCSD algorithm, with parameters set to vales recommended in [4], and the lag-HSIC algorithm. The Type II error of lag-HSIC, presented in the right panel of the Figure 2, is substantially lower than that of KCSD. Most oddly, KCSD error seems to converge to zero in steps. This may be due to the method relying on a spectral decomposition of the signals across a fixed set of bands. As the number of samples increases, the quality of the spectrogram will improve, and dependence will become apparent in bands where it was undetectable at shorter signal lengths.

## References

- [1] M.A. Arcones. The law of large numbers for u-statistics under absolute regularity. *Electron. Comm. Probab.*, 3:13–19, 1998.

<sup>2</sup>Such variation in the smoothness parameter changes the width of the spectral envelope, i.e., the brightness of the sound.

<sup>3</sup>larger extinction rate implies a greater dependence between processes; larger AR component implies a stronger temporal dependence



- 432 [2] L. Bauwens, S. Laurent, and J.V.K. Rombouts. Multivariate GARCH models: a survey. *J. Appl. Econ.*,  
433 21(1):79–109, January 2006.
- 434 [3] A. Berlinet and C. Thomas-Agnan. *Reproducing Kernel Hilbert Spaces in Probability and Statistics*.  
435 Kluwer, 2004.
- 436 [4] M. Besserve, N.K. Logothetis, and B. Schölkopf. Statistical analysis of coupled time series with kernel  
437 cross-spectral density operators. In *NIPS*, page 25352543, 2013.
- 438 [5] I.S. Borisov and N.V. Volodko. Orthogonal series and limit theorems for canonical u-and v-statistics of  
439 stationary connected observations. *Siberian Adv. Math.*, 18(4):242–257, 2008.
- 440 [6] S. Borovkova, R. Burton, and H. Dehling. Limit theorems for functionals of mixing processes with  
441 applications to U-statistics and dimension estimation. *Trans. Amer. Math. Soc.*, 353(11):4261–4318, 2001.
- 442 [7] K. Chwialkowski and A. Gretton. A kernel independence test for random processes. In *ICML*, 2014.
- 443 [8] J. Dedecker, P. Doukhan, G. Lang, S. Louhichi, and C. Prieur. *Weak dependence: with examples and*  
444 *applications*.
- 445 [9] P. Doukhan. *Mixing*. Springer, 1994.
- 446 [10] R.M. Dudley. *Real analysis and probability*, volume 74. Cambridge University Press, 2002.
- 447 [11] K. Fukumizu, A. Gretton, X. Sun, and B. Schölkopf. Kernel measures of conditional dependence. In  
448 *NIPS*, volume 20, pages 489–496, 2007.
- 449 [12] A. Gretton, K.M. Borgwardt, M.J. Rasch, B. Schölkopf, and A. Smola. A kernel two-sample test. *J.*  
450 *Mach. Learn. Res.*, 13:723–773, 2012.
- 451 [13] A. Gretton, O. Bousquet, A. Smola, and B. Schölkopf. Measuring statistical dependence with hilbert-  
452 schmidt norms. In *Algorithmic learning theory*, pages 63–77. Springer, 2005.
- 453 [14] A. Gretton, K. Fukumizu, C. Teo, L. Song, B. Schölkopf, and A. Smola. A kernel statistical test of  
454 independence. In *NIPS*, volume 20, pages 585–592, 2007.
- 455 [15] Z. Harchaoui, F. Bach, and E. Moulines. Testing for homogeneity with kernel Fisher discriminant analysis.  
456 In *NIPS*. 2008.
- 457 [16] P. Hehrmann. *Pitch Perception as Probabilistic Inference*. PhD thesis, Gatsby Computational Neuro-  
458 science Unit, University College London, 2011.
- 459 [17] A. Leucht. Degenerate U-and V-statistics under weak dependence: Asymptotic theory and bootstrap  
460 consistency. *Bernoulli*, 18(2):552–585, 2012.
- 461 [18] A. Leucht and M.H. Neumann. Dependent wild bootstrap for degenerate U- and V-statistics. *Journal of*  
462 *Multivariate Analysis*, 117:257–280, 2013.
- 463 [19] R. Lyons. Distance covariance in metric spaces. *The Annals of Probability*, 41(5):3051–3696, 2013.
- 464 [20] J. Pickands III. Statistical inference using extreme order statistics. *Ann. Statist.*, pages 119–131, 1975.
- 465 [21] D. Sejdinovic, A. Gretton, and W. Bergsma. A kernel test for three-variable interactions. In *NIPS*, pages  
466 1124–1132, 2013.
- 467 [22] D. Sejdinovic, B. Sriperumbudur, A. Gretton, and K. Fukumizu. Equivalence of distance-based and  
468 RKHS-based statistics in hypothesis testing. *Ann. Statist.*, 41(5):2263–2702, 2013.
- 469 [23] D. Sejdinovic, H. Strathmann, M. Lomeli Garcia, C. Andrieu, and A. Gretton. Kernel Adaptive  
470 Metropolis-Hastings. In *ICML*, 2014.
- 471 [24] R. Serfling. *Approximation Theorems of Mathematical Statistics*. Wiley, New York, 1980.
- 472 [25] X. Shao. The dependent wild bootstrap. *J. Amer. Statist. Assoc.*, 105(489):218–235, 2010.
- 473 [26] A. J Smola, A. Gretton, L. Song, and B. Schölkopf. A Hilbert space embedding for distributions. In *Al-*  
474 *gorithmic Learning Theory*, volume LNAI4754, pages 13–31, Berlin/Heidelberg, 2007. Springer-Verlag.
- 475 [27] B. Sriperumbudur, K. Fukumizu, and G. Lanckriet. Universality, characteristic kernels and RKHS em-  
476 bedding of measures. 12:2389–2410, 2011.
- 477 [28] B. Sriperumbudur, A. Gretton, K. Fukumizu, G. Lanckriet, and B. Schölkopf. Hilbert space embeddings  
478 and metrics on probability measures. *J. Mach. Learn. Res.*, 11:1517–1561, 2010.
- 479 [29] M. Sugiyama, T. Suzuki, Y. Itoh, T. Kanamori, and M. Kimura. Least-squares two-sample test. *Neural*  
480 *Networks*, 24(7):735–751, 2011.
- 481 [30] K. Zhang, J. Peters, D. Janzing, B., and B. Schölkopf. Kernel-based conditional independence test and  
482 application in causal discovery. In *UAI*, pages 804–813, 2011.
- 483 [31] X. Zhang, L. Song, A. Gretton, and A. Smola. Kernel measures of independence for non-iid data. In  
484 *NIPS*, volume 22, 2008.
- 485