

# EFFICIENT AGGREGATED KERNEL TESTS USING INCOMPLETE $U$ -STATISTICS

**ANTONIN SCHRAB**

Centre for Artificial Intelligence  
Gatsby Computational Neuroscience Unit  
University College London and Inria London  
a.schrab@ucl.ac.uk

**ILMUN KIM**

Department of Statistics & Data Science  
Department of Applied Statistics  
Yonsei University  
ilmun@yonsei.ac.kr

**BENJAMIN GUEDJ**

Centre for Artificial Intelligence  
University College London and Inria London  
b.guedj@ucl.ac.uk

**ARTHUR GRETTON**

Gatsby Computational Neuroscience Unit  
University College London  
arthur.gretton@gmail.com

## Abstract

We propose a series of computationally efficient, nonparametric tests for the two-sample, independence and goodness-of-fit problems, using the Maximum Mean Discrepancy (MMD), Hilbert Schmidt Independence Criterion (HSIC), and Kernel Stein Discrepancy (KSD), respectively. Our test statistics are incomplete  $U$ -statistics, with a computational cost that interpolates between linear time in the number of samples, and quadratic time, as associated with classical  $U$ -statistic tests. The three proposed tests aggregate over several kernel bandwidths to detect departures from the null on various scales: we call the resulting tests MMDAggInc, HSICAggInc and KSDAggInc. For the test thresholds, we derive a quantile bound for wild bootstrapped incomplete  $U$ -statistics, which is of independent interest. We derive uniform separation rates for MMDAggInc and HSICAggInc, and quantify exactly the trade-off between computational efficiency and the attainable rates: this result is novel for tests based on incomplete  $U$ -statistics, to our knowledge. We further show that in the quadratic-time case, the wild bootstrap incurs no penalty to test power over more widespread permutation-based approaches, since both attain the same minimax optimal rates (which in turn match the rates that use oracle quantiles). We support our claims with numerical experiments on the trade-off between computational efficiency and test power. In the three testing frameworks, we observe that our proposed linear-time aggregated tests obtain higher power than current state-of-the-art linear-time kernel tests.

## 1. Introduction

Nonparametric hypothesis testing is a fundamental field of statistics, and is widely used by the machine learning community and practitioners in numerous other fields, due to the increasing availability of huge amounts of data. When dealing with large-scale datasets, computational cost can quickly emerge as a major issue which might prevent from using expensive tests in practice; constructing efficient tests is therefore crucial for their real-world applications. In this paper, we construct kernel-based aggregated tests using incomplete  $U$ -statistics (Blom, 1976) for the **two-sample**, **independence** and **goodness-of-fit** problems (which we detail in Section 2). The quadratic-time aggregation procedure is known to lead to state-of-the-art powerful tests (Fromont et al., 2012, 2013; Albert et al., 2022; Schrab et al., 2021, 2022), and we propose efficient variants of these well-studied tests, with computational cost interpolating from the classical quadratic-time regime to the linear-time one.

**Related work: aggregated tests.** Kernel selection (or kernel bandwidth selection) is a fundamental problem in nonparametric hypothesis testing because it has a major influence on test power. Motivated by this problem, non-asymptotic aggregated tests, which combine tests with different kernel bandwidths, have been proposed for the two-sample (Fromont et al., 2012, 2013; Kim et al., 2022; Schrab et al., 2021), independence (Albert et al., 2022; Kim et al., 2022), and goodness-of-fit (Schrab et al., 2022) testing frameworks. Li and Yuan (2019) and Balasubramanian et al. (2021) construct similar aggregated tests for these three problems, with the difference that they work in the asymptotic regime. All the mentioned works study aggregated tests in terms of uniform separation rates and minimax rates of testing (Ingster, 1987, 1989, 1993a,b; Baraud, 2002; Tolstikhin et al., 2016). Those rates depend on the sample size and satisfy the following property: if the  $L^2$ -norm difference between the densities is greater than the uniform separation rate, then the test is guaranteed to have high power. All aggregated kernel-based tests in the existing literature have been studied using estimators which are  $U$ -statistics (Hoeffding, 1992) with tests running in quadratic time.

**Related work: linear-time kernel tests.** Several linear-time kernel tests have been proposed for those three testing frameworks. Those include tests using classical linear-time estimators with median bandwidth (Gretton et al., 2012a; Liu et al., 2016) or selecting an optimal bandwidth on held-out data to maximize power (Gretton et al., 2012b), tests using eigenspectrum approximation (Gretton et al., 2009), tests using post-selection inference for adaptive kernel selection, also using incomplete  $U$ -statistics (Yamada et al., 2018, 2019; Lim et al., 2019, 2020; Kübler et al., 2020; Freidling et al., 2021), tests which use a Nyström approximation of the asymptotic null distribution (Zhang et al., 2018; Cherfaoui et al., 2022), random Fourier features tests (Zhang et al., 2018; Zhao and Meng, 2015; Chwialkowski et al., 2015), the current state-of-the-art adaptive tests which use features selected on held-out data to maximize power (Jitkrittum et al., 2016, 2017a,b), as well as tests using neural networks to learn a discrepancy (Grathwohl et al., 2020). We also point out the very relevant works of Kübler et al. (2022) and Huggins and Mackey (2018) on quadratic-time tests, and of Ho and Shieh (2006), Zaremba et al. (2013) and Zhang et al. (2018) on the use of block  $U$ -statistics which have complexity  $\mathcal{O}(N^{1.5})$  for block size  $\sqrt{N}$  where  $N$  is the sample size.

**Contributions and outline.** In Section 2, we present the three testing problems with their associated well-known quadratic-time kernel-based estimators (MMD, HSIC, KSD) which are  $U$ -statistics. We introduce three associated incomplete  $U$ -statistics estimators, which can be computed in linear time, in Section 3. We then provide quantile and variance bounds for generic incomplete  $U$ -statistics using a wild bootstrap, in Section 4. We study the level and power guarantees of linear-time tests using incomplete  $U$ -statistics for a fixed kernel bandwidth, in Section 5. In particular, we obtain uniform separation rates for the two-sample and independence tests over a Sobolev ball, and show that these rates are minimax optimal up to the cost incurred for efficiency of the test. In Section 6, we propose our efficient aggregated tests which combine tests with multiple kernel bandwidths. We prove that the proposed tests are adaptive over Sobolev balls and achieve the same uniform separation rate (up to an iterated logarithmic term) as the tests with optimal bandwidths. As a result of our analysis, we have shown minimax optimality over Sobolev balls of the quadratic-time tests using quantiles estimated with a wild bootstrap. Whether this optimality result also holds for tests using the more general permutation-based procedure to approximate HSIC quantiles, was an open problem formulated by Kim et al. (2022), we prove that it indeed holds in Section 7. We close the paper with numerical experiments in Section 8, where we observe that MMDAggInc, HSICAggInc and KSDAggInc retain high power and outperform other state-of-the-art linear-time kernel tests. Our implementation of the tests and code for reproducibility of the experiments are available online under the MIT License: <https://github.com/antoninschrab/agginc-paper>.

## 2. Background

Here we briefly describe our main problems of interest, comprising the two-sample, independence and goodness-of-fit problems. We approach these problems from a nonparametric point of view using the kernel-based statistics: MMD, HSIC, and KSD. We briefly introduce original forms of these statistics, which can be computed in quadratic time, and also discuss ways of calibrating tests proposed in the literature.

**Two-sample testing.** In this problem, we are given independent samples  $\mathbb{X}_m := (X_i)_{1 \leq i \leq m}$  and  $\mathbb{Y}_n = (Y_j)_{1 \leq j \leq n}$ , consisting of i.i.d. random variables with respective probability density functions<sup>1</sup>  $p$  and  $q$  on  $\mathbb{R}^d$ . We assume we work with balanced sample sizes so that there exists a constant<sup>2</sup>  $C > 0$  such that  $\max(m, n) \leq C \min(m, n)$ . We are interested in testing the null hypothesis  $\mathcal{H}_0 : p = q$  against the alternative  $\mathcal{H}_1 : p \neq q$ ; that is, we want to know if the samples come from the same distribution. Gretton et al. (2012a) propose a non-parametric kernel test based on the *Maximum Mean Discrepancy* (MMD), a measure between probability distributions which uses a characteristic kernel  $k$  (Fukumizu et al., 2008; Sriperumbudur et al., 2011). It can be estimated using a quadratic-time estimator (Gretton et al., 2012a, Lemma 6) which, as noted by Kim et al. (2022), can be expressed as a two-sample  $U$ -statistic (both of second order) (Hoeffding, 1992),

$$\widehat{\text{MMD}}_k^2(\mathbb{X}_m, \mathbb{Y}_n) = \frac{1}{|\mathbf{i}_2^m| |\mathbf{i}_2^n|} \sum_{(i, i') \in \mathbf{i}_2^m} \sum_{(j, j') \in \mathbf{i}_2^n} h_k^{\text{MMD}}(X_i, X_{i'}; Y_j, Y_{j'}), \quad (1)$$

<sup>1</sup>All probability density functions in this paper are with respect to the Lebesgue measure.

<sup>2</sup>We use the convention that all constants are generically denoted by  $C$ , even though they are different.

where  $\mathbf{i}_a^b$  denotes the set of all  $a$ -tuples drawn without replacement from  $\{1, \dots, b\}$  so that  $|\mathbf{i}_a^b| = b \cdot \dots \cdot (b - a + 1)$ , and where, for  $x_1, x_2, y_1, y_2 \in \mathbb{R}^d$ , we let

$$h_k^{\text{MMD}}(x_1, x_2; y_1, y_2) := k(x_1, x_2) - k(x_1, y_2) - k(x_2, y_1) + k(y_1, y_2). \quad (2)$$

**Independence testing.** In this problem, we have access to i.i.d. pairs of samples  $\mathbb{Z}_N := (Z_i)_{1 \leq i \leq N} = ((X_i, Y_i))_{1 \leq i \leq N}$  with joint probability density  $p_{xy}$  on  $\mathbb{R}^{d_x} \times \mathbb{R}^{d_y}$  and marginals  $p_x$  on  $\mathbb{R}^{d_x}$  and  $p_y$  on  $\mathbb{R}^{d_y}$ . We are interested in testing  $\mathcal{H}_0 : p_{xy} = p_x \otimes p_y$  against  $\mathcal{H}_1 : p_{xy} \neq p_x \otimes p_y$ ; that is, we want to know if two components of the pairs of samples are independent or dependent. [Gretton et al. \(2005, 2008\)](#) propose a non-parametric kernel test based on the *Hilbert Schmidt Independence Criterion* (HSIC). It can be estimated using the quadratic-time estimator proposed by [Song et al. \(2012, Equation 5\)](#) which is a fourth-order one-sample  $U$ -statistic

$$\widehat{\text{HSIC}}_{k,\ell}(\mathbb{Z}_N) = \frac{1}{|\mathbf{i}_4^N|} \sum_{(i,j,r,s) \in \mathbf{i}_4^N} h_{k,\ell}^{\text{HSIC}}(Z_i, Z_j, Z_r, Z_s) \quad (3)$$

for characteristic kernels  $k$  on  $\mathbb{R}^{d_x}$  and  $\ell$  on  $\mathbb{R}^{d_y}$  ([Gretton, 2015](#)), and where for  $z_a = (x_a, y_a) \in \mathbb{R}^{d_x} \times \mathbb{R}^{d_y}$ ,  $a = 1, \dots, 4$ , we let

$$h_{k,\ell}^{\text{HSIC}}(z_1, z_2, z_3, z_4) := \frac{1}{4} h_k^{\text{MMD}}(x_1, x_2; x_3, x_4) h_\ell^{\text{MMD}}(y_1, y_2; y_3, y_4). \quad (4)$$

**Goodness-of-fit testing.** For this problem, we are given a model density  $p$  on  $\mathbb{R}^d$  and i.i.d. samples  $\mathbb{Z}_N := (Z_i)_{1 \leq i \leq N}$  drawn from a density  $q$  on  $\mathbb{R}^d$ . The aim is again to test  $\mathcal{H}_0 : p = q$  against  $\mathcal{H}_1 : p \neq q$ ; that is, we want to know if the samples have been drawn from the model. [Chwialkowski et al. \(2016\)](#) and [Liu et al. \(2016\)](#) both construct a non-parametric goodness-of-fit test using the *Kernel Stein Discrepancy* (KSD). A quadratic-time KSD estimator can be computed as the second-order one-sample  $U$ -statistic,

$$\widehat{\text{KSD}}_{p,k}^2(\mathbb{Z}_N) := \frac{1}{|\mathbf{i}_2^N|} \sum_{(i,j) \in \mathbf{i}_2^N} h_{k,p}^{\text{KSD}}(Z_i, Z_j), \quad (5)$$

where the *Stein kernel*  $h_{p,k} : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  is defined as

$$\begin{aligned} h_{k,p}^{\text{KSD}}(x, y) := & \left( \nabla \log p(x)^\top \nabla \log p(y) \right) k(x, y) + \nabla \log p(y)^\top \nabla_x k(x, y) \\ & + \nabla \log p(x)^\top \nabla_y k(x, y) + \sum_{i=1}^d \frac{\partial}{\partial x_i \partial y_i} k(x, y). \end{aligned} \quad (6)$$

In order to guarantee consistency of the Stein goodness-of-fit test ([Chwialkowski et al., 2016, Theorem 2.2](#)), we assume that the kernel  $k$  is  $C_0$ -universal ([Carmeli et al., 2010, Definition 4.1](#)) and that  $\mathbb{E}_q \left[ \left\| \nabla \log \frac{p(Z)}{q(Z)} \right\|_2^2 \right] < \infty$ .

**Quantile estimation.** Multiple strategies have been proposed to estimate the quantiles of test statistics under the null for these three tests. We primarily focus on the wild bootstrap approach ([Chwialkowski et al., 2014](#)), though our results also hold using a parametric bootstrap for the goodness-of-fit setting ([Schrab et al., 2022](#)). In Section 7, we show that the same uniform separation rates can be derived for HSIC quadratic-time tests using permutations instead of a wild bootstrap.

More details on MMD, HSIC, KSD, and on quantile estimation are provided in Appendix A.

### 3. Incomplete $U$ -statistics for MMD, HSIC and KSD

As presented above, the quadratic-time statistics for the two-sample (MMD), independence (HSIC) and goodness-of-fit (KSD) problems can be rewritten as  $U$ -statistics with kernels  $h_k^{\text{MMD}}$ ,  $h_{k,\ell}^{\text{HSIC}}$  and  $h_k^{\text{KSD}}$ , respectively. The computational cost of tests based on these  $U$ -statistics grows quadratically with the sample size. When working with very large sample sizes, as it is often the case in real-world uses of those tests, this quadratic cost can become very problematic, and faster alternative tests are better adapted to this ‘big data’ setting. Multiple linear-time kernel tests have been proposed in the three testing frameworks (see Section 1 for details). We construct linear-time variants of the aggregated kernel tests proposed by Fromont et al. (2013), Albert et al. (2022), Kim et al. (2022), and Schrab et al. (2021, 2022) for the three settings, with the aim of retaining the significant power advantages of the aggregation procedure observed for quadratic-time tests. To this end, we propose to replace the quadratic-time  $U$ -statistics presented in Equations (1), (3) and (5) with second order incomplete  $U$ -statistics (Blom, 1976; Janson, 1984; Lee, 1990),

$$\overline{\text{MMD}}_k^2(\mathbb{X}_m, \mathbb{Y}_n; \mathcal{D}_N) := \frac{1}{|\mathcal{D}_N|} \sum_{(i,j) \in \mathcal{D}_N} h_k^{\text{MMD}}(X_i, X_j; Y_i, Y_j), \quad (7)$$

$$\overline{\text{HSIC}}_{k,\ell}(Z_N; \mathcal{D}_{\lfloor N/2 \rfloor}) := \frac{1}{|\mathcal{D}_{\lfloor N/2 \rfloor}|} \sum_{(i,j) \in \mathcal{D}_{\lfloor N/2 \rfloor}} h_{k,\ell}^{\text{HSIC}}(Z_i, Z_j, Z_{i+\lfloor N/2 \rfloor}, Z_{j+\lfloor N/2 \rfloor}), \quad (8)$$

$$\overline{\text{KSD}}_{p,k}^2(Z_N; \mathcal{D}_N) := \frac{1}{|\mathcal{D}_N|} \sum_{(i,j) \in \mathcal{D}_N} h_{k,p}^{\text{KSD}}(Z_i, Z_j), \quad (9)$$

where for the two-sample problem we let  $N := \min(m, n)$ , and where the *design*  $\mathcal{D}_b$  is a subset of  $\mathbf{i}_2^b$  (the set of all 2-tuples drawn without replacement from  $\{1, \dots, b\}$ ). Note that  $\mathcal{D}_{\lfloor N/2 \rfloor} \subseteq \mathbf{i}_2^{\lfloor N/2 \rfloor} \subset \mathbf{i}_2^N$ . The design can be deterministic. For example, for the two-sample problem with equal even sample sizes  $m = n = N$ , the deterministic design  $\mathcal{D}_N = \{(2a-1, 2a) : a = 1, \dots, N/2\}$  corresponds to the MMD linear-time estimator proposed by Gretton et al. (2012a, Lemma 14). For fixed design size, the elements of the design can also be chosen at random without replacement, in which case the estimators in Equations (7) to (9) become random quantities given the data. The results presented in this paper hold for both deterministic and random (without replacement) design choices. By fixing the design sizes in Equations (7) to (9) to be

$$|\mathcal{D}_N| = |\mathcal{D}_{\lfloor N/2 \rfloor}| = cN \quad (10)$$

for some small constant  $c \in \mathbb{N} \setminus \{0\}$ , we obtain incomplete  $U$ -statistics which can be computed in linear time. Note that by pairing the samples  $Z_i := (X_i, Y_i)$ ,  $i = 1, \dots, N$  for the MMD case and  $\tilde{Z}_i := (Z_i, Z_{i+\lfloor N/2 \rfloor})$ ,  $i = 1, \dots, \lfloor N/2 \rfloor$  for the HSIC case, we observe that all three incomplete  $U$ -statistics of second order have the same form, with only the kernel functions and the design differing. The motivation for defining the estimators in Equations (7) to (9) as incomplete  $U$ -statistics of order 2 (rather than of higher order) derives from the reasoning of Kim et al. (2022, Section 6) using permuted complete  $U$ -statistics for the two-sample and independence problems.

## 4. Quantile and variance bounds for incomplete $U$ -statistics

Here we derive upper quantile and variance bounds for a second order incomplete degenerate  $U$ -statistic with a generic degenerate kernel  $h$ , for some design  $\mathcal{D} \subseteq \mathbf{i}_2^N$ , defined as

$$\bar{U}(\mathbb{Z}_N; \mathcal{D}) := \frac{1}{|\mathcal{D}|} \sum_{(i,j) \in \mathcal{D}} h(Z_i, Z_j).$$

We will use these results to bound the quantiles and variances of our three test statistics for our hypothesis tests in Section 5. The derived bounds are of independent interest.

In the following lemma, building on the results of Lee (1990), we directly derive an upper bound on the variance of the incomplete  $U$ -statistic in terms of the sample size  $N$  and of the design size  $|\mathcal{D}|$ .

**Lemma 4.1.** *The variance of the incomplete  $U$ -statistic can be upper bounded in terms of the quantities  $\sigma_1^2 := \text{var}(\mathbb{E}[h(Z, Z')|Z'])$  and  $\sigma_2^2 := \text{var}(h(Z, Z'))$  with different bounds depending on the design choice. For deterministic design  $\mathcal{D}^d$ , and for random design  $\mathcal{D}^r$ , we have*

$$\text{var}(\bar{U}) \leq C \left( \frac{N}{|\mathcal{D}^d|} \sigma_1^2 + \frac{1}{|\mathcal{D}^d|} \sigma_2^2 \right) \quad \text{and} \quad \text{var}(\bar{U}) \leq C \left( \frac{1}{N} \sigma_1^2 + \left( \frac{1}{|\mathcal{D}^r|} + \frac{1}{N^2} \right) \sigma_2^2 \right).$$

The proof of Lemma 4.1 is deferred to Appendix D. We emphasize the fact that this variance bound also holds for random design with replacement, as considered by Blom (1976) and Lee (1990). For random design, we observe that if  $|\mathcal{D}| \asymp N^2$  then the bound is  $\sigma_1^2/N + \sigma_2^2/N^2$  which is the variance bound of the complete  $U$ -statistic (Albert et al., 2022, Lemma 10). If  $N \leq |\mathcal{D}| \leq N^2$ , the variance bound is  $\sigma_1^2/N + \sigma_2^2/|\mathcal{D}|$ , and if  $|\mathcal{D}| \leq N$  it is  $\sigma_2^2/|\mathcal{D}|$  since  $\sigma_1^2 \leq \sigma_2^2/2$  (Blom, 1976, Equation 2.1).

Kim et al. (2022) develop exponential concentration bounds for permuted complete  $U$ -statistics, and Cl  men  on et al. (2013) study the uniform approximation of  $U$ -statistics by incomplete  $U$ -statistics. To the best of our knowledge, no quantile bounds have yet been obtained for incomplete  $U$ -statistics in the literature. While permutations are well-suited for complete  $U$ -statistics (Kim et al., 2022), using them with incomplete  $U$ -statistics results in having to compute new kernel values, and this comes at an extra computational cost we would like to avoid. Restricting the set of permutations to those for which the kernel values have already been computed for the original incomplete  $U$ -statistic corresponds exactly to using a wild bootstrap (Schrab et al., 2021, Appendix B). Hence, we consider the wild bootstrapped second order incomplete  $U$ -statistic

$$\bar{U}^\epsilon(\mathbb{Z}_N; \mathcal{D}) := \frac{1}{|\mathcal{D}|} \sum_{(i,j) \in \mathcal{D}} \epsilon_i \epsilon_j h(Z_i, Z_j) \tag{11}$$

for i.i.d. Rademacher random variables  $\epsilon_1, \dots, \epsilon_N$  with values in  $\{-1, 1\}$ , for which we derive an exponential concentration bound (quantile bound). We note the in-depth work of Chwialkowski et al. (2014) on the wild bootstrap procedure for kernel tests with applications to quadratic-time MMD and HSIC tests. We now provide exponential tail bounds for wild bootstrapped incomplete  $U$ -statistics.

**Lemma 4.2.** *There exists some constant  $C > 0$  such that, for every  $t \geq 0$ , we have*

$$\mathbb{P}_\epsilon(|\bar{U}^\epsilon| \geq t \mid \mathbb{Z}_N, \mathcal{D}) \leq 2 \exp\left(-C \frac{t}{A_{\text{inc}}}\right) \leq 2 \exp\left(-C \frac{t}{A}\right)$$

where  $A_{\text{inc}}^2 := |\mathcal{D}|^{-2} \sum_{(i,j) \in \mathcal{D}} h(Z_i, Z_j)^2$  and  $A^2 := |\mathcal{D}|^{-2} \sum_{(i,j) \in \mathbf{i}_2^N} h(Z_i, Z_j)^2$ .

Lemma 4.2 is proved in Appendix E. While the second bound in Lemma 4.2 is less tight, it has the benefit of not depending on the choice of design  $\mathcal{D}$  but only on the design size  $|\mathcal{D}|$  which is usually fixed.

## 5. Efficient kernel tests using incomplete $U$ -statistics

We now formally define the hypothesis tests obtained using the incomplete  $U$ -statistics with a wild bootstrap. This is done for fixed kernel bandwidths  $\lambda \in (0, \infty)^{d_x}, \mu \in (0, \infty)^{d_y}$ , for the kernels<sup>3</sup>

$$k_\lambda(x, y) := \prod_{i=1}^{d_x} \frac{1}{\lambda_i} K_i\left(\frac{x_i - y_i}{\lambda_i}\right), \quad \ell_\mu(x, y) := \prod_{i=1}^{d_y} \frac{1}{\mu_i} L_i\left(\frac{x_i - y_i}{\mu_i}\right), \quad (12)$$

for characteristic kernels  $(x, y) \mapsto K_i(x - y), (x, y) \mapsto L_i(x - y)$  on  $\mathbb{R} \times \mathbb{R}$  for functions  $K_i, L_i \in L^1(\mathbb{R}) \cap L^2(\mathbb{R})$  integrating to 1. We unify the notation for the three testing frameworks. For the two-sample and goodness-of-fit problems, we work only with  $k_\lambda$  and have  $d = d_x$ . For the independence problem, we work with the two kernels  $k_\lambda$  and  $\ell_\mu$ , and for ease of notation we let  $d := d_x + d_y$  and  $\lambda_{d_x+i} := \mu_i$  for  $i = 1, \dots, d_y$ . We also simply write  $p := p_{xy}$  and  $q := p_x \otimes p_y$ . We let  $\bar{U}_\lambda$  and  $h_\lambda$  denote either  $\overline{\text{MMD}}_{k_\lambda}^2$  and  $h_{k_\lambda}^{\text{MMD}}$ , or  $\overline{\text{HSIC}}_{k_\lambda, \ell_\mu}$  and  $h_{k_\lambda, \ell_\mu}^{\text{HSIC}}$ , or  $\overline{\text{KSD}}_{p, k_\lambda}^2$  and  $h_{k_\lambda, p}^{\text{KSD}}$ , respectively. We denote the design size of the incomplete  $U$ -statistics in Equations (7) to (9) by

$$L := |\mathcal{D}_N| = |\mathcal{D}_{\lfloor N/2 \rfloor}|.$$

For the three testing frameworks, we estimate the quantiles of the test statistics by simulating the null hypothesis using a wild bootstrap, as done in the case of complete  $U$ -statistics by Fromont et al. (2012), Schrab et al. (2021) for the two-sample problem, and by Schrab et al. (2022) for the goodness-of-fit problem. This is done by considering the original test statistic  $U_\lambda^{B_1+1} := \bar{U}_\lambda$  together with  $B_1$  wild bootstrapped incomplete  $U$ -statistics  $U_\lambda^1, \dots, U_\lambda^{B_1}$  computed as in Equation (11), and estimating the  $(1-\alpha)$ -quantile with a Monte Carlo approximation

$$\hat{q}_{1-\alpha}^\lambda := \inf \left\{ t \in \mathbb{R} : 1 - \alpha \leq \frac{1}{B_1 + 1} \sum_{b=1}^{B_1+1} \mathbf{1}(U_\lambda^b \leq t) \right\} = U_\lambda^{\bullet [B_1(1-\alpha)]}, \quad (13)$$

where  $U_\lambda^{\bullet 1} \leq \dots \leq U_\lambda^{\bullet B_1+1}$  are the sorted elements  $U_\lambda^1, \dots, U_\lambda^{B_1+1}$ . The test  $\Delta_\alpha^\lambda$  is defined as rejecting the null if the original test statistic  $\bar{U}_\lambda$  is greater than the estimated  $(1-\alpha)$ -quantile, that is,

$$\Delta_\alpha^\lambda(\mathbb{Z}_N) := \mathbf{1}\left(\bar{U}_\lambda(\mathbb{Z}_N) > \hat{q}_{1-\alpha}^\lambda\right).$$

We show in Proposition 5.1 that the test  $\Delta_\alpha^\lambda$  has well-calibrated asymptotic level for goodness-of-fit testing, and well-calibrated non-asymptotic level for two-sample and independence testing. The proof of the latter non-asymptotic guarantee is based on the exchangeability of  $U_\lambda^1, \dots, U_\lambda^{B_1+1}$  under the null hypothesis along with the result of Romano and Wolf (2005, Lemma 1). A similar proof strategy can be found in Fromont et al. (2012, Proposition 2), Albert et al. (2022, Proposition 1), and Schrab et al. (2021, Proposition 1). The

<sup>3</sup>Our results are presented for bandwidth selection, but they hold for more general kernel selection settings, as considered by Schrab et al. (2022). The goodness-of-fit results hold for a wider range of kernels including the IMQ (inverse multiquadric) kernel (Gorham and Mackey, 2017), as in Schrab et al. (2022).



exchangeability of wild bootstrapped incomplete  $U$ -statistics for independence testing does not follow directly from the mentioned works. We show this through an intriguing connection between the MMD kernel and the HSIC kernel (proof deferred to Appendix C).

**Proposition 5.1.** *The test  $\Delta_\alpha^\lambda$  has level  $\alpha \in (0, 1)$ , i.e.,  $\mathbb{P}_{\mathcal{H}_0}(\Delta_\alpha^\lambda(\mathbb{Z}_N) = 1) \leq \alpha$ . This holds non-asymptotically for the two-sample and independence cases, and asymptotically for goodness-of-fit.<sup>4</sup>*

Having established the validity of the test  $\Delta_\alpha^\lambda$ , we now study power guarantees for it in terms of the  $L^2$ -norm of the difference in densities  $\|p - q\|_2$ . In Theorem 5.2, we show for the three tests that, if  $\|p - q\|_2$  exceeds some threshold, we can guarantee high test power. For the two-sample and independence problems, we derive a uniform separation rate (Baraud, 2002) over Sobolev balls

$$\mathcal{S}_d^s(R) := \left\{ f \in L^1(\mathbb{R}^d) \cap L^2(\mathbb{R}^d) : \int_{\mathbb{R}^d} \|\xi\|_2^{2s} |\hat{f}(\xi)|^2 d\xi \leq (2\pi)^d R^2 \right\}, \quad (14)$$

with radius  $R > 0$  and smoothness parameter  $s > 0$ . This uniform separation rate is the smallest value of  $t$  such that for any alternative with  $\|p - q\|_2 > t$  and  $p - q \in \mathcal{S}_d^s(R)$  the probability of type II error of  $\Delta_\alpha^\lambda$  can be controlled by  $\beta \in (0, 1)$ . Before presenting Theorem 5.2, we need to introduce more notation unified over the three testing frameworks; we define the integral transform  $T_\lambda$  as

$$(T_\lambda f)(x) := \int_{\mathbb{R}^d} f(y) \mathcal{K}_\lambda(x, y) dy \quad (15)$$

for  $f \in L^2(\mathbb{R}^d)$ ,  $x \in \mathbb{R}^d$ , where  $\mathcal{K}_\lambda := k_\lambda$  for the two-sample problem,  $\mathcal{K}_\lambda := k_\lambda \otimes \ell_\mu$  for the independence problem, and  $\mathcal{K}_\lambda := h_{k_\lambda, p}^{\text{KSD}}$  for the goodness-of-fit problem. Note that, for the two-sample and independence testing frameworks, since  $\mathcal{K}_\lambda$  is translation-invariant, the integral transform corresponds to a convolution. However, this is not true for the goodness-of-fit framework as  $h_{k_\lambda, p}^{\text{KSD}}$  is not translation-invariant. We are now in a position to present our main contribution in Theorem 5.2: we derive a power guarantee condition for our tests using incomplete  $U$ -statistics, and a uniform separation rate over Sobolev balls for the two-sample and independence settings.

**Theorem 5.2.** *(i) Assume  $\|p\|_\infty \leq M$  and  $\|q\|_\infty \leq M$  for some  $M > 0$ ,  $\alpha \in (0, e^{-1})$ ,  $\beta \in (0, 1)$ ,  $B_1 \geq \frac{3}{\alpha^2} (\ln(\frac{8}{\beta}) + \alpha(1 - \alpha))$  and  $\lambda \in (0, \infty)^d$  with<sup>5</sup>  $\sigma_{2,\lambda}^2 := \mathbb{E}[h_\lambda(Z, Z')^2] \geq 1$ . If*

$$\|p - q\|_2^2 \geq \|(p - q) - T_\lambda(p - q)\|_2^2 + C \frac{N \ln(1/\alpha)}{L} \frac{1}{\beta} \sigma_{2,\lambda} \quad \text{for some constant } C > 0,$$

*then  $\mathbb{P}_{\mathcal{H}_1}(\Delta_\alpha^\lambda(\mathbb{Z}_N) = 0) \leq \beta$  (control of probability of type II error by  $\beta$ ).*

*(ii) Fix  $R > 0$  and  $s > 0$ , and consider the bandwidths  $\lambda_i^* := (N/L)^{2/(4s+d)}$  for  $i = 1, \dots, d$ . For MMD and HSIC, the uniform separation rate of  $\Delta_\alpha^{\lambda^*}$  over the Sobolev ball  $\mathcal{S}_d^s(R)$  is (up to a constant)*

$$(N/L)^{2s/(4s+d)}.$$

<sup>4</sup>Level is non-asymptotic for the goodness-of-fit case using a parametric bootstrap (Schrab et al., 2022).

<sup>5</sup>It is sufficient for this condition to hold up to a constant independent of  $\lambda$ , that is,  $\sigma_{2,\lambda}^2 \gtrsim 1$ . For MMD and HSIC, we have  $\sigma_{2,\lambda}^2 \lesssim 1/\lambda_1 \cdots \lambda_d$ , so in order to satisfy  $\sigma_{2,\lambda}^2 \gtrsim 1$  we must have  $\lambda_1 \cdots \lambda_d \lesssim 1$ , which is the condition required by Schrab et al. (2021) and Albert et al. (2022).



The proof of Theorem 5.2 relies on the variance and quantile bounds presented in Lemmas 4.1 and 4.2, and also uses results of Albert et al. (2022) and Schrab et al. (2021, 2022) on complete  $U$ -statistics. The details can be found in Appendix F. The power condition in Theorem 5.2 corresponds to a variance-bias decomposition; for large bandwidths the bias term (first term) dominates, while for small bandwidths the variance term (second term which also controls the quantile) dominates. We recall that the minimax (i.e. optimal) rate over the Sobolev ball  $\mathcal{S}_d^s(R)$  is  $(1/N)^{2s/(4s+d)}$  for the two-sample (Li and Yuan, 2019, Theorem 5 (ii)) and independence (Albert et al., 2022, Theorem 4) problems. We highlight that the rate for our incomplete  $U$ -statistic test has the same dependence in the exponent as the minimax rate; that is

$$(N/L)^{2s/(4s+d)} = (1/N)^{2s/(4s+d)} (N^2/L)^{2s/(4s+d)}$$

where we recall that  $L \leq N^2$  is the design size and  $N$  is the sample size. We reach the following conclusions.

- If  $L \asymp N^2$  then the test runs in quadratic time and the minimax rate is recovered.
- If  $N \lesssim L \lesssim N^2$  then the rate still converges to 0 but the cost  $(N^2/L)^{2s/(4s+d)}$  is incurred in the minimax rate (computational efficiency / rate convergence trade-off).
- If  $L \lesssim N$  then there is no guarantee that the rate converges to 0.

## 6. Efficient aggregated tests using incomplete $U$ -statistics

We now introduce our aggregated tests that combine single tests with different bandwidths. Our aggregation scheme is similar to those of Fromont et al. (2013), Albert et al. (2022) and Schrab et al. (2021, 2022), and can yield an adaptive test to the unknown smoothness parameter  $s$  of the Sobolev ball  $\mathcal{S}_d^s(R)$ , with relatively low price. Let  $\Lambda$  be a finite collection of bandwidths,  $(w_\lambda)_{\lambda \in \Lambda}$  be associated weights satisfying  $\sum_{\lambda \in \Lambda} w_\lambda \leq 1$  and  $u_\alpha$  be some correction term defined shortly in Equation (16). Then, using the incomplete  $U$ -statistic  $\bar{U}_\lambda$ , we define our aggregated test  $\Delta_\alpha^\Lambda$  as

$$\Delta_\alpha^\Lambda(\mathbb{Z}_N) := \mathbf{1}\left(\bar{U}_\lambda(\mathbb{Z}_N) > \hat{q}_{1-u_\alpha w_\lambda}^\lambda \text{ for some } \lambda \in \Lambda\right).$$

The levels of the single tests are weighted and adjusted with a correction term

$$u_\alpha := \sup_{B_3} \left\{ u \in \left(0, \min_{\lambda \in \Lambda} w_\lambda^{-1}\right) : \frac{1}{B_2} \sum_{b=1}^{B_2} \mathbf{1}\left(\max_{\lambda \in \Lambda} \left(\tilde{U}_\lambda^b - U_\lambda^{\bullet[B_1(1-uw_\lambda)]}\right) > 0\right) \leq \alpha \right\}, \quad (16)$$

where the wild bootstrapped incomplete  $U$ -statistics  $\tilde{U}_\lambda^1, \dots, \tilde{U}_\lambda^{B_2}$  computed as in Equation (11) are used to perform a Monte Carlo approximation of the probability under the null, and where the supremum is estimated using  $B_3$  steps of bisection method. Proposition 5.1, along with the reasoning of Schrab et al. (2021, Proposition 8), ensures that  $\Delta_\alpha^\Lambda$  has non-asymptotic level  $\alpha$  for the two-sample and independence cases, and asymptotic level  $\alpha$  for the goodness-of-fit case. We refer to the three aggregated test constructed using incomplete  $U$ -statistics as MMDAggInc, HSICAggInc and KSDAggInc. The computational complexity of those tests is  $\mathcal{O}(|\Lambda|(B_1 + B_2)L)$ , which means that if  $L \asymp N$  as in Equation (10), the tests run efficiently in linear time in the sample size.

We formally record error guarantees of  $\Delta_\alpha^\Lambda$  and derive uniform separation rates over Sobolev balls.

**Theorem 6.1.** (i) Assume  $\|p\|_\infty \leq M$  and  $\|q\|_\infty \leq M$  for some  $M > 0$ . Consider a collection  $\Lambda$  such that  $\sigma_{2,\lambda}^2 := \mathbb{E}[h_\lambda(Z, Z')^2] \geq 1$  for all  $\lambda \in \Lambda$ . For  $\alpha \in (0, e^{-1})$ ,  $B_1 \geq (\max_{\lambda \in \Lambda} w_\lambda^{-2}) \frac{12}{\alpha^2} (\ln(\frac{8}{\beta}) + \alpha(1 - \alpha))$ ,  $B_2 \geq \frac{8}{\alpha^2} \ln(\frac{2}{\beta})$ ,  $B_3 \geq \log_2(\frac{4}{\alpha} \min_{\lambda \in \Lambda} w_\lambda^{-1})$ , if

$$\|p - q\|_2^2 \geq \min_{\lambda \in \Lambda} \left( \|(p - q) - T_\lambda(p - q)\|_2^2 + C \frac{N \ln(1/(\alpha w_\lambda))}{L \beta} \sigma_{2,\lambda} \right) \text{ for some } C > 0,$$

then the probability of type II error is controlled as  $\mathbb{P}_{\mathcal{H}_1}(\Delta_\alpha^\Lambda(\mathbb{Z}_N) = 0) \leq \beta$ .

(ii) Consider the collections of bandwidths and weights (independent of  $R$  and  $s$ )

$$\Lambda := \left\{ (2^{-\ell}, \dots, 2^{-\ell}) \in (0, \infty)^d : \ell \in \left\{ 1, \dots, \left\lceil \frac{2}{d} \log_2 \left( \frac{L/N}{\ln(\ln(L/N))} \right) \right\rceil \right\} \right\}, \quad w_\lambda := \frac{6}{\pi^2 \ell^2}.$$

For two-sample and independence problems, the uniform separation rate of  $\Delta_\alpha^\Lambda$  over the Sobolev balls  $\{\mathcal{S}_d^s(R) : R > 0, s > 0\}$  is (up to a constant)

$$\left( \frac{\ln(\ln(L/N))}{L/N} \right)^{2s/(4s+d)}.$$

The extension from Theorem 5.2 to Theorem 6.1 has been proved for complete  $U$ -statistics in the two-sample (Fromont et al., 2013; Schrab et al., 2021), independence (Albert et al., 2022) and goodness-of-fit (Schrab et al., 2022) testing frameworks. The proof of Theorem 6.1 follows with the same reasoning by simply replacing  $N$  with  $L/N$  as we work with incomplete  $U$ -statistics; this ‘replacement’ is theoretically justified by Theorem 5.2. From Theorem 6.1, the aggregated test  $\Delta_\alpha^\Lambda$  is *adaptive* over Sobolev balls  $\{\mathcal{S}_d^s(R) : R > 0, s > 0\}$ : the test  $\Delta_\alpha^\Lambda$  does not depend on the unknown smoothness parameter  $s$  (unlike  $\Delta_\alpha^{\lambda^*}$  in Theorem 5.2) and achieves the minimax rate up to an iterated logarithmic factor and up to the cost incurred for efficiency of the test (i.e.  $L/N$  instead of  $N$ ).

## 7. Minimax optimal permuted quadratic-time aggregated independence test

Considering Theorem 6.1 with our incomplete  $U$ -statistic with full design  $\mathcal{D} = \mathbf{i}_2^N$  for which  $L \asymp N^2$ , we have proved that the quadratic-time two-sample and independence aggregated tests using a wild bootstrap achieve the rate  $(\ln(\ln(N))/N)^{2s/(4s+d)}$  over the Sobolev balls  $\{\mathcal{S}_d^s(R) : R > 0, s > 0\}$ . This is the minimax rate (Li and Yuan, 2019; Albert et al., 2022), up to some iterated logarithmic term. For the two-sample problem, Kim et al. (2022) and Schrab et al. (2021) show that this is also true using complete  $U$ -statistics with either a wild bootstrap or permutations. Whether the equivalent statement for independence test with permutations holds is unknown; the rate can be proved using theoretical (unknown) quantiles with a Gaussian kernel (Albert et al., 2022), but has not yet been proved using permutations. Kim et al. (2022, Proposition 8.7) consider this problem, again using a Gaussian kernel, but they do not obtain the correct dependence on  $\alpha$  (i.e.  $\ln(1/\alpha)$  is replaced with  $\alpha^{-1/2}$ ), hence they cannot recover the desired rate. As pointed out by Kim et al. (2022, Section 8): ‘It remains an open question as to whether [the power guarantee] continues to hold when  $\alpha^{-1/2}$  is replaced by  $\ln(1/\alpha)$ ’. We now prove that we can improve the  $\alpha$ -dependence to  $\ln(1/\alpha)^{3/2}$  for any bounded kernel of the form

presented in Equation (12), and that this allows us to obtain the desired rate over Sobolev balls  $\{\mathcal{S}_d^s(R) : R > 0, s > d/4\}$ . The assumption  $s > d/4$  imposes a stronger smoothness restriction on  $p - q \in \mathcal{S}_d^s(R)$ , which is similarly also considered by Li and Yuan (2019).

**Theorem 7.1.** *Consider the quadratic-time independence test using the complete  $U$ -statistic HSIC estimator with a quantile estimated using permutations as done by Kim et al. (2022, Proposition 8.7), with kernels as in (12) for bounded functions  $K_i$  and  $L_j$  for  $i = 1, \dots, d_x$ ,  $j = 1, \dots, d_y$ .*

(i) *Consider the assumptions of Theorem 5.2. For fixed  $R > 0$  and  $s > d/4$ , with the bandwidths  $\lambda_i^* := N^{-2/(4s+d)}$  for  $i = 1, \dots, d$ , the probability of type II error of the test is controlled by  $\beta$  when*

$$\|p - q\|_2^2 \geq \|(p - q) - T_{\lambda^*}(p - q)\|_2^2 + C \frac{1}{N} \frac{\ln(1/\alpha)^{3/2}}{\beta \sqrt{\lambda_1^* \cdots \lambda_d^*}} \quad \text{for some constant } C > 0.$$

The uniform separation rate over the Sobolev ball  $\mathcal{S}_d^s(R)$  is, up to a constant,  $(1/N)^{2s/(4s+d)}$ .

(ii) *Consider the assumptions of Theorem 6.1, the uniform separation rate over the Sobolev balls  $\{\mathcal{S}_d^s(R) : R > 0, s > d/4\}$  is  $(\ln(\ln(N))/N)^{2s/(4s+d)}$ , up to a constant, with the collections*

$$\Lambda := \left\{ (2^{-\ell}, \dots, 2^{-\ell}) \in (0, \infty)^d : \ell \in \left\{ 1, \dots, \left\lceil \frac{2}{d} \log_2 \left( \frac{N}{\ln(\ln(N))} \right) \right\rceil \right\} \right\}, \quad w_\lambda := \frac{6}{\pi^2 \ell^2}.$$

The proof of Theorem 7.1, in Appendix G, uses the exponential concentration bound of Kim et al. (2022, Theorem 6.3) for permuted complete  $U$ -statistics. As discussed by Kim et al. (2022, Section 8.3), their proposed sample-splitting method can also be used to obtain the correct dependency on  $\alpha$ .

## 8. Experiments

For the two-sample problem, we consider testing samples drawn from a uniform density on  $[0, 1]^d$  against samples drawn from a perturbed uniform density. For the independence problem, the joint density is a perturbed uniform density on  $[0, 1]^{d_x+d_y}$ , the marginals are then simply uniform densities. Those perturbed uniform densities can be shown to lie in Sobolev balls (Li and Yuan, 2019; Albert et al., 2022), to which our tests are adaptive. For the goodness-of-fit problem, we use a Gaussian-Bernoulli Restricted Boltzmann Machine as first considered by Liu et al. (2016) in this testing framework. Details on the experiments (e.g. model/test parameters) are presented in Appendix B.

We consider our incomplete aggregated tests MMDAggInc, HSICAggInc and KSDAggInc, with parameter  $R \in \{1, \dots, N - 1\}$  which fixes the deterministic design to consist of the first  $R$  sub-diagonals of the  $N \times N$  matrix, that is,  $\mathcal{D} := \{(i, i + r) : i = 1, \dots, N - r \text{ for } r = 1, \dots, R\}$  with size  $|\mathcal{D}| = RN - R(R - 1)/2$ . We run our incomplete tests with  $R \in \{1, 100, 200\}$  and also consider the complete test which uses the full design  $\mathcal{D} = \mathbf{i}_2^N$ . We compare their performances with current linear-time state-of-the-art tests: OST PSI (Kübler et al., 2020) which performs kernel selection using post selection inference, ME, SCF, FSIC and FSSD (Jitkrittum et al., 2016, 2017a,b) which evaluate the witness functions at a finite set of locations chosen to maximize the power, and LSD (Grathwohl et al., 2020) which uses a neural network to learn the Stein discrepancy (see Appendix B for details).

Similar trends are observed across all our experiments in Figure 1, in the three testing frameworks, when varying the sample size, the dimension, and the difficulty of the problem

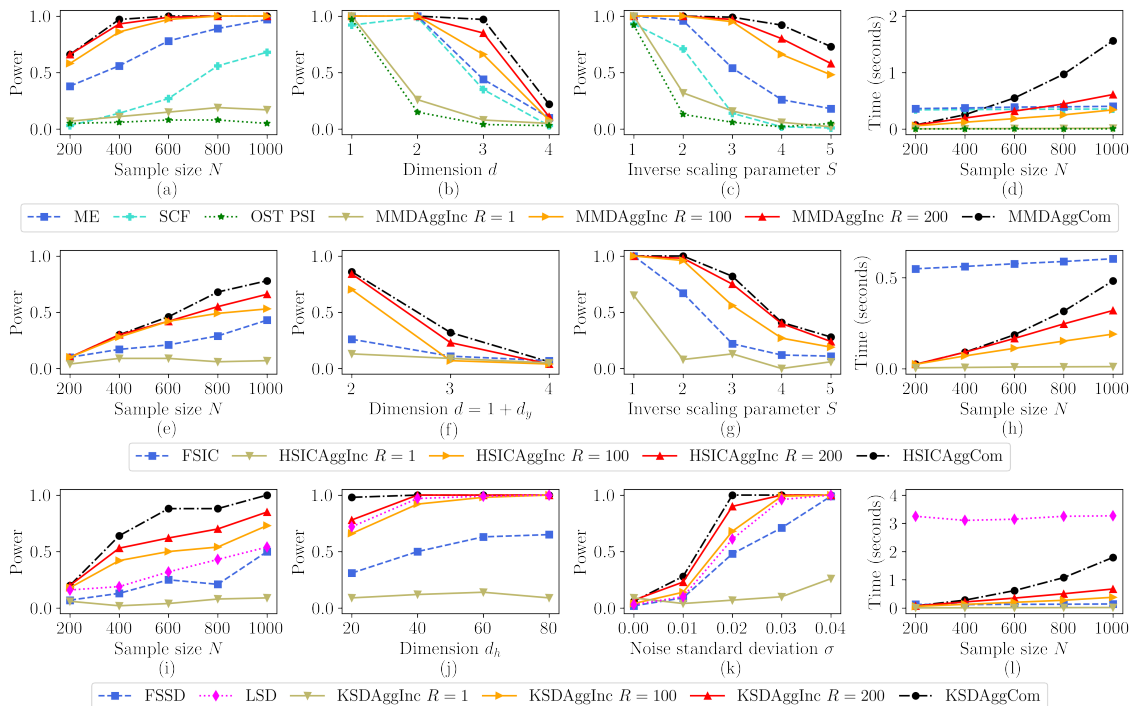


Figure 1: Two-sample (a–d) and independence (e–h) experiments using perturbed uniform densities. Goodness-of-fit (i–l) experiment using a Gaussian-Bernoulli Restricted Boltzmann Machine. The power results are averaged over 100 repetitions and the run times over 20 repetitions.

(scale of perturbations or noise level). The linear-time tests  $\text{AggInc } R = 200$  almost match the power obtained by the quadratic-time tests  $\text{AggCom}$  in all settings (except in Figure 1(i) where the difference is larger) while being computationally much more efficient as can be seen in Figure 1(d,h,l). The incomplete tests with  $R = 100$  has power only slightly below the one using  $R = 200$ , and runs roughly twice as fast (Figure 1(d,h,l)). In all experiments, those three tests ( $\text{AggInc } R = 100, 200$  and  $\text{AggCom}$ ) have significantly higher power than the linear-time tests which optimize test locations (ME, SCF, FSIC and FSSD); in the two-sample case the aggregated tests run faster for small sample size but slower for large sample size, in the independence case the aggregated tests run much faster, and in the goodness-of-fit case the tests optimizing test locations run faster. While both types of tests are linear, we note that the run times of the tests of [Jitkrittum et al. \(2016, 2017a,b\)](#) increase slower with the sample size than our aggregated tests with  $R = 100, 200$ , but a fixed computational cost is incurred for the optimization step, even for small sample sizes. In the goodness-of-fit framework, LSD matches the power of  $\text{KSDAggInc } R = 100$  when varying the noise level in Figure 1(k) ( $\text{KSDAggInc } R = 200$  has higher power), and matches the power of  $\text{KSDAggInc } R = 200$  when varying the hidden dimension in Figure 1(j) where  $d_x = 100$ . When varying the sample size in Figure 1(i), both  $\text{KSDAggInc}$  tests with  $R = 100, 200$  achieve much higher power than LSD. Unsurprisingly,  $\text{AggInc } R = 1$ , which runs much faster than all the aforementioned tests, has low power in every experiment. For the two-sample problem, it obtains slightly higher power than OST PSI which runs even faster.

## Acknowledgements

Antonin Schrab acknowledges support from the U.K. Research and Innovation under grant number EP/S021566/1. Ilmun Kim acknowledges support from the Yonsei University Research Fund of 2021-22-0332 as well as support from the Basic Science Research Program through the National Research Foundation of Korea (NRF) funded by the Ministry of Education (2022R1A4A1033384). Benjamin Guedj acknowledges partial support by the U.S. Army Research Laboratory and the U.S. Army Research Office, and by the U.K. Ministry of Defence and the U.K. Engineering and Physical Sciences Research Council (EPSRC) under grant number EP/R013616/1; Benjamin Guedj also acknowledges partial support from the French National Agency for Research, grants ANR-18-CE40-0016-01 and ANR-18-CE23-0015-02. Arthur Gretton acknowledges support from the Gatsby Charitable Foundation.

## Bibliography

- Albert, M., Laurent, B., Marrel, A., and Meynaoui, A. (2022). Adaptive test of independence based on HSIC measures. *The Annals of Statistics*, 50(2):858–879.
- Aronszajn, N. (1950). Theory of reproducing kernels. *Transactions of the American Mathematical Society*, 68(3):337–404.
- Balasubramanian, K., Li, T., and Yuan, M. (2021). On the optimality of kernel-embedding based goodness-of-fit tests. *Journal of Machine Learning Research*, 22(1).
- Baraud, Y. (2002). Non-asymptotic minimax rates of testing in signal detection. *Bernoulli*, 1(8(5):577–606).
- Blom, G. (1976). Some properties of incomplete U-statistics. *Biometrika*, 63(3):573–580.
- Carmeli, C., De Vito, E., Toigo, A., and Umanitá, V. (2010). Vector valued reproducing kernel Hilbert spaces and universality. *Analysis and Applications*, 8(01):19–61.
- Chebyshev, P. L. (1899). Oeuvres. *Commissionaires de l'Académie Impériale des Sciences*, 1.
- Cherfaoui, F., Kadri, H., Anthoine, S., and Ralaivola, L. (2022). A discrete RKHS standpoint for Nyström MMD. *HAL preprint hal-03651849*.
- Chwialkowski, K., Sejdinovic, D., and Gretton, A. (2014). A wild bootstrap for degenerate kernel tests. In *Advances in neural information processing systems*, pages 3608–3616.
- Chwialkowski, K., Strathmann, H., and Gretton, A. (2016). A kernel test of goodness of fit. In *International Conference on Machine Learning*, pages 2606–2615. PMLR.
- Chwialkowski, K. P., Ramdas, A., Sejdinovic, D., and Gretton, A. (2015). Fast two-sample testing with analytic representations of probability measures. In *Advances in Neural Information Processing Systems*, volume 28, pages 1981–1989.
- Cléménçon, S., Robbiano, S., and Tressou, J. (2013). Maximal deviations of incomplete U-statistics with applications to empirical risk sampling. In *Proceedings of the 2013 SIAM International Conference on Data Mining*, pages 19–27. SIAM.

- de la Peña, V. H. and Giné, E. (1999). *Decoupling: From Dependence to Independence*. Springer Science & Business Media.
- Duembgen, L. (1998). Symmetrization and decoupling of combinatorial random elements. *Statistics & probability letters*, 39(4):355–361.
- Dvoretzky, A., Kiefer, J., and Wolfowitz, J. (1956). Asymptotic minimax character of the sample distribution function and of the classical multinomial estimator. *The Annals of Mathematical Statistics*, pages 642–669.
- Freidling, T., Poinard, B., Climente-González, H., and Yamada, M. (2021). Post-selection inference with HSIC-Lasso. In *International Conference on Machine Learning*, pages 3439–3448. PMLR.
- Fromont, M., Laurent, B., Lerasle, M., and Reynaud-Bouret, P. (2012). Kernels based tests with non-asymptotic bootstrap approaches for two-sample problems. In *Conference on Learning Theory*, PMLR.
- Fromont, M., Laurent, B., and Reynaud-Bouret, P. (2013). The two-sample problem for Poisson processes: Adaptive tests with a nonasymptotic wild bootstrap approach. *The Annals of Statistics*, 41(3):1431–1461.
- Fukumizu, K., Gretton, A., Sun, X., and Schölkopf, B. (2008). Kernel measures of conditional dependence. In *Advances in Neural Information Processing Systems*, volume 1, pages 489–496.
- Gorham, J. and Mackey, L. (2017). Measuring sample quality with kernels. In *International Conference on Machine Learning*, pages 1292–1301. PMLR.
- Grathwohl, W., Wang, K.-C., Jacobsen, J.-H., Duvenaud, D., and Zemel, R. (2020). Learning the Stein discrepancy for training and evaluating energy-based models without sampling. In *International Conference on Machine Learning*, pages 3732–3747. PMLR.
- Gretton, A. (2015). A simpler condition for consistency of a kernel independence test. *arXiv preprint arXiv:1501.06103*.
- Gretton, A., Borgwardt, K. M., Rasch, M. J., Schölkopf, B., and Smola, A. (2012a). A kernel two-sample test. *Journal of Machine Learning Research*, 13:723–773.
- Gretton, A., Fukumizu, K., Harchaoui, Z., and Sriperumbudur, B. K. (2009). A fast, consistent kernel two-sample test. *Advances in Neural Information Processing Systems*, 22.
- Gretton, A., Fukumizu, K., Teo, C. H., Song, L., Schölkopf, B., and Smola, A. J. (2008). A kernel statistical test of independence. In *Advances in Neural Information Processing Systems*, volume 1, pages 585–592.
- Gretton, A., Herbrich, R., Smola, A., Bousquet, O., and Schölkopf, B. (2005). Kernel methods for measuring independence. *Journal of Machine Learning Research*, 6:2075–2129.
- Gretton, A., Sejdinovic, D., Strathmann, H., Balakrishnan, S., Pontil, M., Fukumizu, K., and Sriperumbudur, B. K. (2012b). Optimal kernel choice for large-scale two-sample tests. In *Advances in Neural Information Processing Systems*, volume 1, pages 1205–1213.

- Ho, H.-C. and Shieh, G. S. (2006). Two-stage U-statistics for hypothesis testing. *Scandinavian journal of statistics*, 33(4):861–873.
- Hoeffding, W. (1992). A class of statistics with asymptotically normal distribution. In *Breakthroughs in Statistics*, pages 308–334. Springer.
- Huggins, J. and Mackey, L. (2018). Random feature Stein discrepancies. *Advances in Neural Information Processing Systems*, 31.
- Ingster, Y. I. (1987). Minimax testing of nonparametric hypotheses on a distribution density in the  $L_p$  metrics. *Theory of Probability & its Applications*, 31(2):333–337.
- Ingster, Y. I. (1989). An asymptotically minimax test of the hypothesis of independence. *Journal of Soviet Mathematics*, 1(44:466–476).
- Ingster, Y. I. (1993a). Asymptotically minimax hypothesis testing for nonparametric alternatives. *Journal of Soviet Mathematics*, 1(44:466–476).
- Ingster, Y. I. (1993b). Minimax testing of the hypothesis of independence for ellipsoids in  $l_p$ . *Zapiski Nauchnykh Seminarov POMI*, 1(207:77–97).
- Janson, S. (1984). The asymptotic distributions of incomplete U-statistics. *Zeitschrift für Wahrscheinlichkeitstheorie und Verwandte Gebiete*, 66(4):495–505.
- Jitkrittum, W., Szabó, Z., Chwialkowski, K. P., and Gretton, A. (2016). Interpretable distribution features with maximum testing power. In *Advances in Neural Information Processing Systems*, volume 29, pages 181–189.
- Jitkrittum, W., Szabó, Z., and Gretton, A. (2017a). An adaptive test of independence with analytic kernel embeddings. In *International Conference on Machine Learning (ICML)*, pages 1742–1751.
- Jitkrittum, W., Xu, W., Szabó, Z., Fukumizu, K., and Gretton, A. (2017b). A linear-time kernel goodness-of-fit test. In *Advances in Neural Information Processing Systems*, pages 262–271.
- Key, O., Fernandez, T., Gretton, A., and Briol, F.-X. (2021). Composite goodness-of-fit tests with kernels. *arXiv preprint arXiv:2111.10275*.
- Kim, I., Balakrishnan, S., and Wasserman, L. (2022). Minimax optimality of permutation tests. *The Annals of Statistics*, 50(1):225–251.
- Kingma, D. P. and Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Kübler, J. M., Jitkrittum, W., Schölkopf, B., and Muandet, K. (2020). Learning kernel tests without data splitting. In *Advances in Neural Information Processing Systems 33*, pages 6245–6255. Curran Associates, Inc.
- Kübler, J. M., Jitkrittum, W., Schölkopf, B., and Muandet, K. (2022). A witness two-sample test. In *International Conference on Artificial Intelligence and Statistics*, pages 1403–1419. PMLR.
- Lee, J. (1990). *U-statistics: Theory and Practice*. Citeseer.



- Leucht, A. and Neumann, M. H. (2013). Dependent wild bootstrap for degenerate U- and V-statistics. *Journal of Multivariate Analysis*, 117:257–280.
- Li, T. and Yuan, M. (2019). On the optimality of gaussian kernel based nonparametric tests against smooth alternatives. *arXiv preprint arXiv:1909.03302*.
- Lim, J. N., Yamada, M., Jitkrittum, W., Terada, Y., Matsui, S., and Shimodaira, H. (2020). More powerful selective kernel tests for feature selection. In *International Conference on Artificial Intelligence and Statistics*, pages 820–830. PMLR.
- Lim, J. N., Yamada, M., Schölkopf, B., and Jitkrittum, W. (2019). Kernel Stein tests for multiple model comparison. In *Advances in Neural Information Processing Systems*, pages 2240–2250.
- Liu, Q., Lee, J., and Jordan, M. (2016). A kernelized Stein discrepancy for goodness-of-fit tests. In *International Conference on Machine Learning*, pages 276–284. PMLR.
- Massart, P. (1990). The tight constant in the Dvoretzky-Kiefer-Wolfowitz inequality. *The Annals of Probability*, 18(3):1269–1283.
- Müller, A. (1997). Integral probability metrics and their generating classes of functions. *Advances in Applied Probability*, 1:429–443.
- Ramachandran, P., Zoph, B., and Le, Q. V. (2017). Searching for activation functions. *arXiv preprint arXiv:1710.05941*.
- Romano, J. P. and Wolf, M. (2005). Exact and approximate stepdown methods for multiple hypothesis testing. *Journal of the American Statistical Association*, 100(469):94–108.
- Schrab, A., Guedj, B., and Gretton, A. (2022). KSD aggregated goodness-of-fit test. *arXiv preprint arXiv:2202.00824*.
- Schrab, A., Kim, I., Albert, M., Laurent, B., Guedj, B., and Gretton, A. (2021). MMD aggregated two-sample test. *arXiv preprint arXiv:2110.15073*.
- Shao, X. (2010). The dependent wild bootstrap. *Journal of the American Statistical Association*, 105(489):218–235.
- Song, L., Smola, A. J., Gretton, A., Bedo, J., and Borgwardt, K. M. (2012). Feature selection via dependence maximization. *Journal of Machine Learning Research*, 13:1393–1434.
- Sriperumbudur, B. K., Fukumizu, K., and Lanckriet, G. R. (2011). Universality, characteristic kernels and RKHS embedding of measures. *Journal of Machine Learning Research*, 12(7).
- Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. (2014). Dropout: a simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(1):1929–1958.
- Sutherland, D. J., Tung, H.-Y., Strathmann, H., De, S., Ramdas, A., Smola, A., and Gretton, A. (2017). Generative models and model criticism via optimized maximum mean discrepancy. In *International Conference on Learning Representations*.
- Tolstikhin, I., Sriperumbudur, B. K., and Schölkopf, B. (2016). Minimax estimation of maximum mean discrepancy with radial kernels. *Advances in Neural Information Processing Systems*, 29.

- Yamada, M., Umezū, Y., Fukumizu, K., and Takeuchi, I. (2018). Post selection inference with kernels. In *International Conference on Artificial Intelligence and Statistics*, pages 152–160. PMLR.
- Yamada, M., Wu, D., Tsai, Y. H., Ohta, H., Salakhutdinov, R., Takeuchi, I., and Fukumizu, K. (2019). Post selection inference with incomplete maximum mean discrepancy estimator. In *International Conference on Learning Representations*.
- Zaremba, W., Gretton, A., and Blaschko, M. (2013). B-test: A non-parametric, low variance kernel two-sample test. *Advances in neural information processing systems*, 26.
- Zhang, Q., Filippi, S., Gretton, A., and Sejdinovic, D. (2018). Large-scale kernel methods for independence testing. *Statistics and Computing*, 28(1):113–130.
- Zhao, J. and Meng, D. (2015). Fastmmd: Ensemble of circular discrepancy for efficient two-sample test. *Neural computation*, 27(6):1345–1372.

## A. Details on MMD, HSIC and KSD

In this section, we present more details than those presented in Section 2 on the Maximum Mean Discrepancy, on the Hilbert Schmidt Independence Criterion, and on the Kernel Stein Discrepancy.

**Maximum Mean Discrepancy.** Gretton et al. (2012a) introduce the *Maximum Mean Discrepancy* (MMD) which is a measure between probability densities  $p$  and  $q$  on  $\mathbb{R}^d$ . It is defined as the integral probability metric (IPM; Müller, 1997) over a reproducing kernel Hilbert space  $\mathcal{H}_k$  (RKHS; Aronszajn, 1950) with associated kernel  $k$ . Gretton et al. (2012a, Lemma 4) show that the MMD is equal to the  $\mathcal{H}_k$ -norm of the difference between the mean embeddings  $\mu_p(u) := \mathbb{E}_{X \sim p}[k(X, u)]$  and  $\mu_q(u) := \mathbb{E}_{Y \sim q}[k(Y, u)]$  for  $u \in \mathbb{R}^d$ . The square of the MMD is equal to

$$\begin{aligned} \text{MMD}_k^2(p, q) &:= \left( \sup_{f \in \mathcal{H}_k : \|f\|_{\mathcal{H}_k} \leq 1} |\mathbb{E}_p[f(X)] - \mathbb{E}_q[f(Y)]| \right)^2 \\ &= \|\mu_p - \mu_q\|_{\mathcal{H}_k}^2 \\ &= \mathbb{E}_{p,p}[k(X, X')] - 2 \mathbb{E}_{p,q}[k(X, Y)] + \mathbb{E}_{q,q}[k(Y, Y')] \end{aligned}$$

where  $X$  and  $X'$  (respectively  $Y$  and  $Y'$ ) are independent. Using a characteristic kernel (Fukumizu et al., 2008; Sriperumbudur et al., 2011) guarantees that  $\text{MMD}_k^2(p, q) = 0$  if and only if  $p = q$ , a crucial property for using the MMD to construct a two-sample test. With i.i.d. samples  $\mathbb{X}_m := (X_i)_{1 \leq i \leq m}$  from  $p$  and i.i.d. samples  $\mathbb{Y}_n = (Y_j)_{1 \leq j \leq n}$  from  $q$ , Gretton et al. (2012a, Lemma 6) propose to use the unbiased quadratic-time MMD estimator  $\widehat{\text{MMD}}_k^2(\mathbb{X}_m, \mathbb{Y}_n)$  defined as

$$\begin{aligned} &\frac{1}{m(m-1)} \sum_{(i,i') \in \mathbf{i}_2^m} k(X_i, X_{i'}) - \frac{2}{mn} \sum_{i=1}^m \sum_{j=1}^n k(X_i, Y_j) + \frac{1}{n(n-1)} \sum_{(j,j') \in \mathbf{i}_2^n} k(Y_j, Y_{j'}) \\ &= \frac{\mathbf{1}^\top \tilde{\mathbf{K}}_{\text{XX}} \mathbf{1}}{m(m-1)} - 2 \frac{\mathbf{1}^\top \tilde{\mathbf{K}}_{\text{XY}} \mathbf{1}}{mn} + \frac{\mathbf{1}^\top \tilde{\mathbf{K}}_{\text{YY}} \mathbf{1}}{n(n-1)} \end{aligned}$$

where  $\tilde{\mathbf{K}}_{\text{XX}}$  and  $\tilde{\mathbf{K}}_{\text{YY}}$  are the kernel matrices  $\mathbf{K}_{\text{XX}} := (k(X_i, X_j))_{1 \leq i, j \leq m}$  and  $\mathbf{K}_{\text{YY}} := (k(Y_i, Y_j))_{1 \leq i, j \leq n}$  with diagonal entries set to 0, where  $\mathbf{K}_{\text{XY}} := (k(X_i, Y_j))_{1 \leq i \leq m, 1 \leq j \leq n}$ , and where  $\mathbf{1}$  is a one-dimensional vector with all entries equal to 1 of variable length determined by the context<sup>6</sup>. As noted by Kim et al. (2022), this MMD estimator can be rewritten as a two-sample  $U$ -statistic (both of second order) (Hoeffding, 1992)

$$\widehat{\text{MMD}}_k^2(\mathbb{X}_m, \mathbb{Y}_n) = \frac{1}{|\mathbf{i}_2^m| |\mathbf{i}_2^n|} \sum_{(i,i') \in \mathbf{i}_2^m} \sum_{(j,j') \in \mathbf{i}_2^n} h_k^{\text{MMD}}(X_i, X_{i'}; Y_j, Y_{j'})$$

where  $\mathbf{i}_a^b$  denotes the set of all  $a$ -tuples drawn without replacement from  $\{1, \dots, b\}$  so that  $|\mathbf{i}_a^b| = b \cdots (b - a + 1)$ , for example  $|\mathbf{i}_2^m| = m(m-1)$ , and where, for  $x_1, x_2, y_1, y_2 \in \mathbb{R}^d$ , we let

$$h_k^{\text{MMD}}(x_1, x_2; y_1, y_2) := k(x_1, x_2) - k(x_1, y_2) - k(x_2, y_1) + k(y_1, y_2).$$

This kernel can easily be symmetrized (Kim et al., 2022) using a symmetrization trick (Duembgen, 1998), this corresponds to working with

$$\bar{h}_k^{\text{MMD}}(x_1, x_2; y_1, y_2) := \frac{1}{2!2!} \sum_{(i_1, i_2) \in \mathbf{i}_2^2} \sum_{(j_1, j_2) \in \mathbf{i}_2^2} h_k^{\text{MMD}}(x_1, x_2; y_1, y_2)$$

<sup>6</sup>We use this convention for the notation  $\mathbf{1}$  in this whole section.

and the MMD expression as a  $U$ -statistic still holds when replacing  $h_k^{\text{MMD}}$  with its symmetrized variant  $\bar{h}_k^{\text{MMD}}$ .

**Hilbert Schmidt Independence Criterion.** For a joint probability density  $p_{xy}$  on  $\mathbb{R}^{d_x} \times \mathbb{R}^{d_y}$  with marginals  $p_x$  on  $\mathbb{R}^{d_x}$  and  $p_y$  on  $\mathbb{R}^{d_y}$ , [Gretton et al. \(2005\)](#) introduce the *Hilbert Schmidt Independence Criterion* (HSIC) which is defined as

$$\begin{aligned} \text{HSIC}_{k,\ell}(p_{xy}) &:= \text{MMD}_{\kappa}^2(p_{xy}, p_x p_y) \\ &= \mathbb{E}_{p_{xy}, p_{xy}} \left[ k(X, X') \ell(Y, Y') \right] - 2 \mathbb{E}_{p_{xy}} \left[ \mathbb{E}_{p_x} [k(X, X')] \mathbb{E}_{p_y} [\ell(Y, Y')] \right] \\ &\quad + \mathbb{E}_{p_x, p_x} \left[ k(X, X') \right] \mathbb{E}_{p_y, p_y} \left[ \ell(Y, Y') \right]. \end{aligned}$$

with kernels  $k$  on  $\mathbb{R}^{d_x}$  and  $\ell$  on  $\mathbb{R}^{d_y}$  giving the product kernel  $\kappa((x, y), (x', y')) := k(x, x') \ell(y, y')$  on  $\mathbb{R}^{d_x} \times \mathbb{R}^{d_y}$ . With i.i.d. pairs of samples  $Z_N := (Z_i)_{1 \leq i \leq N} = ((X_i, Y_i))_{1 \leq i \leq N}$  drawn from  $p_{xy}$ , a natural unbiased HSIC estimator ([Gretton et al., 2008](#); [Song et al., 2012](#)) is then

$$\begin{aligned} \widehat{\text{HSIC}}_{k,\ell}(Z_N) &:= \frac{1}{|\mathbf{i}_2^N|} \sum_{(i,j) \in \mathbf{i}_2^N} k(X_i, X_j) \ell(Y_i, Y_j) - \frac{2}{|\mathbf{i}_3^N|} \sum_{(i,j,r) \in \mathbf{i}_3^N} k(X_i, X_j) \ell(Y_i, Y_r) \\ &\quad + \frac{1}{|\mathbf{i}_4^N|} \sum_{(i,j,r,s) \in \mathbf{i}_4^N} k(X_i, X_j) \ell(Y_r, Y_s) \\ &= \frac{1}{|\mathbf{i}_4^N|} \sum_{(i,j,r,s) \in \mathbf{i}_4^N} h_{k,\ell}^{\text{HSIC}}(Z_i, Z_j, Z_r, Z_s) \end{aligned}$$

which is a fourth-order one-sample  $U$ -statistic. For  $z_a = (x_a, y_a) \in \mathbb{R}^{d_x} \times \mathbb{R}^{d_y}$ ,  $a = 1, \dots, 4$ , we let

$$h_{k,\ell}^{\text{HSIC}}(z_1, z_2, z_3, z_4) := \frac{1}{4} h_k^{\text{MMD}}(x_1, x_2; x_3, x_4) h_\ell^{\text{MMD}}(y_1, y_2; y_3, y_4).$$

We stress the fact that this HSIC estimator can actually be computed in quadratic time as shown by [Song et al. \(2012, Equation 5\)](#) who provide the following closed-form expression

$$\widehat{\text{HSIC}}_{k,\ell}(Z_N) = \frac{1}{N(N-3)} \left( \text{tr}(\tilde{\mathbf{K}}\tilde{\mathbf{L}}) + \frac{\mathbf{1}^\top \tilde{\mathbf{K}} \mathbf{1} \mathbf{1}^\top \tilde{\mathbf{L}} \mathbf{1}}{(N-1)(N-2)} - \frac{2}{N-2} \mathbf{1}^\top \tilde{\mathbf{K}} \tilde{\mathbf{L}} \mathbf{1} \right)$$

where  $\tilde{\mathbf{K}}$  and  $\tilde{\mathbf{L}}$  are the kernel matrices  $\mathbf{K} := (k(X_i, X_j))_{1 \leq i, j \leq N}$  and  $\mathbf{L} := (\ell(Y_i, Y_j))_{1 \leq i, j \leq N}$  with diagonal entries set to 0. Again, this kernel can be symmetrized ([Song et al., 2012](#); [Kim et al., 2022](#)) using a symmetrization trick ([Duembgen, 1998](#)), and the HSIC expression as a  $U$ -statistic still holds when replacing  $h_k^{\text{HSIC}}$  with its symmetrized variant

$$\bar{h}_k^{\text{HSIC}}(z_1, z_2, z_3, z_4) := \frac{1}{4!} \sum_{(i_1, i_2, i_3, i_4) \in \mathbf{i}_4^4} h_k^{\text{HSIC}}(z_{i_1}, z_{i_2}, z_{i_3}, z_{i_4}).$$

**Kernel Stein Discrepancy.** For probability densities  $p$  and  $q$  on  $\mathbb{R}^d$ , [Chwialkowski et al. \(2016\)](#) and [Liu et al. \(2016\)](#) introduce the *Kernel Stein Discrepancy* (KSD) defined as

$$\begin{aligned} \text{KSD}_{p,k}^2(q) &:= \text{MMD}_{h_{k,p}^{\text{KSD}}}^2(q, p) \\ &= \mathbb{E}_{q,q} [h_{k,p}^{\text{KSD}}(Z, Z')] - 2 \mathbb{E}_{q,p} [h_{k,p}^{\text{KSD}}(Z, X)] + \mathbb{E}_{p,p} [h_{k,p}^{\text{KSD}}(X, X')] \\ &= \mathbb{E}_{q,q} [h_{k,p}^{\text{KSD}}(Z, Z')] \end{aligned}$$

where the *Stein kernel*  $h_{p,k}: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  is defined as

$$h_{k,p}^{\text{KSD}}(x, y) := \left( \nabla \log p(x)^\top \nabla \log p(y) \right) k(x, y) + \nabla \log p(y)^\top \nabla_x k(x, y) \\ + \nabla \log p(x)^\top \nabla_y k(x, y) + \sum_{i=1}^d \frac{\partial}{\partial x_i \partial y_i} k(x, y).$$

The Stein kernel satisfies the Stein identity  $\mathbb{E}_p[h_{k,p}^{\text{KSD}}(X, \cdot)] = 0$ . The KSD is particularly useful for the goodness-of-fit setting with a model density  $p$  and i.i.d. samples  $\mathbb{Z}_N := (Z_i)_{1 \leq i \leq N}$  drawn from a density  $q$  because it admits an estimator which does not require samples from the model  $p$ . The quadratic-time KSD estimator can be computed as the second-order one-sample  $U$ -statistic

$$\widehat{\text{KSD}}_{p,k}^2(\mathbb{Z}_N) := \frac{1}{\binom{N}{2}} \sum_{(i,j) \in \mathfrak{I}_2^N} h_{k,p}^{\text{KSD}}(Z_i, Z_j) = \frac{\mathbf{1}^\top \tilde{\mathbf{H}} \mathbf{1}}{N(N-1)}$$

where  $\tilde{\mathbf{H}}$  is the kernel matrix  $\mathbf{H} := (h_{k,p}^{\text{KSD}}(Z_i, Z_j))_{1 \leq i, j \leq N}$  with diagonal entries set to 0. The Stein kernel  $h_k^{\text{KSD}}$  is already symmetric, we can write  $\tilde{h}_k^{\text{KSD}}(x, y) := h_k^{\text{KSD}}(x, y)$  for all  $x, y \in \mathbb{R}^d$  for consistency of notation.

**Quantile estimation.** There exist many approaches to estimating the quantiles of the test statistics under the null hypothesis in the three frameworks: using the quantile of a known distribution-free asymptotic null distribution (Gretton et al., 2008, 2012a), sampling from an asymptotic null distribution with eigenspectrum approximation (Gretton et al., 2009), using permutations (Gretton et al., 2008; Albert et al., 2022; Kim et al., 2022; Schrab et al., 2021), using a wild bootstrap (Fromont et al., 2012; Chwialkowski et al., 2014, 2016; Schrab et al., 2021, 2022), using parametric bootstrap (Key et al., 2021; Schrab et al., 2022), using other bootstrap methods (Liu et al., 2016), to name but a few. Permutation-based tests have been shown to correctly control the non-asymptotic level for the two-sample (Schrab et al., 2021; Kim et al., 2022) and independence (Albert et al., 2022; Kim et al., 2022) problems. For the two-sample test, using a wild bootstrap also guarantees well-calibrated non-asymptotic level (Fromont et al., 2012; Schrab et al., 2021). For the goodness-of-fit setting, while a wild bootstrap guarantees only control of the asymptotic level (Chwialkowski et al., 2016), using a parametric bootstrap results in well-calibrated non-asymptotic level (Schrab et al., 2022). In this work, we focus on the wild bootstrap approach, though we point out that our results also hold using a parametric bootstrap for the goodness-of-fit setting as done by Schrab et al. (2022).

## B. Detailed experimental protocol

In this section, we present details on the experiments and on the tests considered. We also show empirically that all the tests have well-calibrated levels.

**Implementation and computational resources.** All experiments have been run on an AMD Ryzen Threadripper 3960X 24 Cores 128Gb RAM CPU at 3.8GHz, except the LSD test (Grathwohl et al., 2020) for which a neural network has been trained using a NVIDIA RTX A5000 24Gb Graphics Card. The overall runtime of all the experiments is of the order of a couple of hours (significant speedup can be obtained by using parallel computing). For the ME, SCF, FSIC and FSSD tests of Jitkrittum et al. (2016, 2017a,b), and for the LSD

test of [Grathwohl et al. \(2020\)](#), we use the implementations of the respective authors. The implementation of our computationally efficient aggregated tests, as well as the code for reproducibility of the experiments, are available [here](#) under the MIT license.

**Kernels.** For the two-sample and independence experiments, we use the Gaussian kernel<sup>7</sup> with equal bandwidths  $\lambda_1 = \dots = \lambda_d = \tilde{\lambda}$ , which is defined as

$$k_\lambda(x, y) := \exp\left(-\sum_{i=1}^d \frac{(x_i - y_i)^2}{\lambda_i^2}\right) = \exp\left(-\frac{\|x - y\|_2^2}{\tilde{\lambda}^2}\right),$$

and similarly for the kernel  $\ell_\mu$ . As shown by [Gorham and Mackey \(2017\)](#), a more appropriate kernel for goodness-of-fit testing is the IMQ (inverse multiquadric) kernel

$$k_\lambda(x, y) := \left(1 + \sum_{i=1}^d \frac{(x_i - y_i)^2}{\lambda_i^2}\right)^{-\beta_k} = \left(1 + \frac{\|x - y\|_2^2}{\tilde{\lambda}^2}\right)^{-\beta_k} \propto \left(\tilde{\lambda}^2 + \|x - y\|_2^2\right)^{-\beta_k}$$

for some  $\beta_k \in (0, 1)$ . In our goodness-of-fit experiments, we fix the parameter  $\beta_k = 0.5$ .

**Two-sample and independence experiments.** In our experiments, we consider perturbed uniform densities, those can be shown to lie in Sobolev balls and are used to derive the minimax rates over Sobolev balls for the two-sample and independence problems ([Albert et al., 2022](#); [Li and Yuan, 2019](#)). For the two-sample problem, we consider testing samples drawn from a uniform density against samples drawn from a perturbed uniform density, as considered by [Schrab et al. \(2021\)](#), see Equation 17 for formal definition and Figure 2 for illustrations). We scale the perturbations so that the perturbed density takes value in the whole interval  $[0, 2]$ , we then consider some inverse scaling parameter  $S \geq 1$  such that it takes value in the interval  $[1 - 1/S, 1 + 1/S]$ . Intuitively, as  $S$  increases, the perturbation is shrunk. In Figure 1(a,d), we consider 2 perturbations with inverse scaling parameter  $S = 2$  in dimension  $d = 1$  and vary the sample size  $N \in \{200, 400, 600, 800, 1000\}$ . In Figure 1(b), we vary the dimension  $d \in \{1, 2, 3, 4\}$  for 1 perturbation with  $S = 1$  and  $N = 1000$ . In Figure 1(c), we use 1 perturbation with  $d = 1$  and  $N = 1000$ , we vary the inverse scaling parameter  $S \in \{1, 2, 3, 4, 5\}$ . For the independence problem, we draw samples from the joint perturbed uniform density in dimension  $d_x + d_y$ , the marginals are simply uniform densities in dimensions  $d_x$  and  $d_y$ , respectively. We fix  $d_x = 1$  and vary  $d_y$  exactly as in the two-sample setting. The parameters for the independence experiments in Figure 1(e–h) are the same as those of the two-sample experiments in Figure 1(a–d) detailed above (with the only difference that for Figure 1(f) we consider  $d_y \in \{1, 2, 3\}$ ).

**Goodness-of-fit experiments.** In Figure 1(i–l), we use a Gaussian-Bernoulli Restricted Boltzmann Machine (GBRBM) with the same setting considered by [Liu et al. \(2016\)](#), [Grathwohl et al. \(2020\)](#) and [Schrab et al. \(2022\)](#). This is a hidden variable model with a continuous observable variable in  $\mathbb{R}^{d_x}$  and a hidden binary variable in  $\{-1, 1\}^{d_h}$ , the joint density is intractable but the score function admits a closed form. The GBRBM has parameters  $b \in \mathbb{R}^{d_x}$  and  $c \in \mathbb{R}^{d_h}$ , which are drawn from Gaussian standard distributions, and a parameter  $B \in \mathbb{R}^{d_x \times d_h}$ . For the model  $p$ , the elements of  $B$  are sampled uniformly from  $\{-1, 1\}$  (i.i.d. Rademacher variables). The samples come from a GBRBM  $q$  with the same parameters as the model  $p$  but where some Gaussian noise  $\mathcal{N}(0, \sigma)$  is injected into

<sup>7</sup>In practice, we do not need to normalize the kernels to integrate to 1 since our tests are invariant to multiplying the kernel by a scalar.

the elements of  $B$ . In Figure 1(i,l), we consider dimensions  $d_x = 50$  and  $d_h = 40$  with noise standard deviation  $\sigma = 0.02$  and we vary the sample size  $N \in \{200, 400, 600, 800, 1000\}$ . In Figure 1(j), we fix  $d_x = 100$ ,  $N = 1000$ ,  $\sigma = 0.03$  and we vary the hidden dimension  $d_h \in \{20, 40, 60, 80\}$ . For fixed observed dimension  $d_x$ , as the hidden dimension  $d_h$  increases the size of  $B \in \mathbb{R}^{d_x \times d_h}$  becomes larger, so there is more evidence of the noise being injected, which makes the problem easier. Hence, the test power increases as  $d_h$  increases for fixed  $d_x$ . In Figure 1(k), we consider dimensions  $d_x = 50$  and  $d_h = 40$  with sample size  $N = 1000$ , we vary the noise standard deviations  $\sigma \in \{0, 0.01, 0.02, 0.03, 0.04\}$ .

**AggInc tests.** For MMDAggInc and KSDAggInc, we use the collection of bandwidths

$$\Lambda := \left\{ 2^i \lambda_{\text{med}} \mathbf{1}_d : i \in \{-3, -2, -1, 0\} \right\} \quad \text{where} \quad \lambda_{\text{med}} := \text{median} \left\{ \|z_i - z_j\|_2 : (i, j) \in \mathbf{i}_2^N \right\}.$$

where  $\mathbf{1}_d$  is a  $d$ -dimensional vector with all entries equal to 1. For HSICAggInc, we work with the collection of pairs of bandwidths

$$\Lambda := \left\{ (2^i \lambda_{\text{med}} \mathbf{1}_{d_x}, 2^j \mu_{\text{med}} \mathbf{1}_{d_y}) : i, j \in \{-2, -1, 0\} \right\}$$

for the kernels  $k_\lambda$  and  $\ell_\mu$  defined in Equation (12), where

$$\lambda_{\text{med}} := \text{median} \left\{ \|x_i - x_j\|_2 : (i, j) \in \mathbf{i}_2^N \right\} \quad \text{and} \quad \mu_{\text{med}} := \text{median} \left\{ \|y_i - y_j\|_2 : (i, j) \in \mathbf{i}_2^N \right\}.$$

All aggregated tests are run with uniform weights defined as  $w_\lambda := 1/|\Lambda|$  for all  $\lambda \in \Lambda$ . The design choice consists of  $R$  sub-diagonals of the kernel matrix for  $R \in \{1, 100, 200\}$ , it is formally defined in Section 8. We use  $B_1 = 500$  and  $B_2 = 500$  wild bootstrapped statistics to estimate the quantiles and the probability under the null for the correction in Equation (16), respectively. For that correction term, we use  $B_3 = 50$  steps of bisection method to approximate the supremum.

**ME, SCF, FSIC and FSSD tests.** Jitkrittum et al. (2016) uses the two-sample tests ME and SCF proposed by Chwialkowski et al. (2015) with features which are chosen to maximise a lower bound on the test power. The ME test is based on analytic Mean Embeddings while the SCF test uses the difference in Smooth Characteristic Functions. For the independence problem, Jitkrittum et al. (2017a) constructs a FSIC test which uses their proposed normalised Finite Set Independence Criterion. Jitkrittum et al. (2017b) proposes a goodness-of-fit test based on the Finite Set Stein Discrepancy (FSSD). All those tests utilise test statistics which evaluate the witness function of either MMD, HSIC, or KSD, at some test locations (i.e. features) chosen on held-out data to maximise test power. For the two-sample SCF test, the test locations are in the frequency domain rather than in the spatial domain. All tests are used with 10 test locations which are chosen on half of the data, as done in the experiments of Jitkrittum et al. (2016, 2017a,b). The ME and SCF tests use the quantiles of the known chi-square asymptotic null distributions. The FSIC test uses 500 permutations to simulate the null hypothesis and compute the test threshold to ensure well-calibrated non-asymptotic level. The FSSD test simulates 2000 samples from the asymptotic null distribution (weighted sum of chi-squares) with the eigenvalues being computed from the covariance matrix with respect to the observed sample. For the two-sample and independence tests, the bandwidths of the Gaussian kernels are selected during the optimization procedure. For the goodness-of-fit test, the bandwidth of the IMQ (inverse multiquadric) kernel is set to one as done by Jitkrittum et al. (2017b), following the recommendation of Gorham and Mackey (2017).



**LSD test.** The Kernelised Stein Discrepancy (KSD) is a Stein Discrepancy (Gorham and Mackey, 2017) where the class of functions is taken to be the unit ball of a reproducing kernel Hilbert space (RKHS). Grathwohl et al. (2020) propose to instead consider some more expressive class of functions consisting of neural networks, resulting in the Learned Stein Discrepancy (LSD). For goodness-of-fit testing, they propose to split the data into training (80%), validation (10%) and testing (10%) sets. They construct a test statistic which is asymptotically normal under both  $\mathcal{H}_0$  and  $\mathcal{H}_1$ . Using the training set, they train the parametrised neural network to maximise the test power by optimizing a proxy for it which is derived following the reasonings of Gretton et al. (2012b), Sutherland et al. (2017) and Jitkrittum et al. (2017b). They perform model selection on the validation set. Finally, they run the test on the testing set using the quantile of the asymptotic normal distribution under the null. As in the experiments of Grathwohl et al. (2020), a 2-layer MLP with 300 units per layer and with Swish nonlinearity (Ramachandran et al., 2017) is used. Their model is trained using the Adam optimizer (Kingma and Ba, 2014) for 1000 iterations, with dropout (Srivastava et al., 2014), with weight decay of strength 0.0005, with learning rate  $10^{-3}$ , and with  $L^2$  regularising strength 0.5.

**Well-calibrated levels.** All tests are run with level  $\alpha = 0.05$ , it is verified in Tables 1 to 6 that all tests have well-calibrated levels for the three testing frameworks, when varying either the sample size or the dimension. The levels plotted are averages obtained across 200 repetitions, this explains the small fluctuations observed from the desired test level  $\alpha = 0.05$ . The settings of those six experiments correspond to the settings of the experiments presented in Figure 1(a,b,e,f,i,j) detailed above, with the difference that we are working under the null hypothesis (i.e. perturbed uniform densities are replaced with uniform densities, and the noise standard deviation for the GBRBB is set to  $\sigma = 0$ ).

Table 1: Two-sample level experiment using uniform densities varying the sample size.

Sample size	ME	SCF	OST PSI	MMDAggInc $R = 1$	MMDAggInc $R = 100$	MMDAggInc $R = 200$	MMDAggCom
200	0.055	0.005	0.045	0.04	0.05	0.055	0.055
400	0.08	0.01	0.04	0.035	0.06	0.03	0.03
600	0.08	0.005	0.105	0.085	0.04	0.04	0.07
800	0.05	0.005	0.055	0.075	0.03	0.035	0.055
1000	0.075	0.005	0.045	0.045	0.015	0.02	0.05

Table 2: Two-sample level experiment using uniform densities varying the dimension.

Dimension	ME	SCF	OST PSI	MMDAggInc $R = 1$	MMDAggInc $R = 100$	MMDAggInc $R = 200$	MMDAggCom
1	0.045	0	0.035	0.02	0.045	0.04	0.045
2	0.045	0.035	0.085	0.1	0.05	0.04	0.035
3	0.04	0.05	0.04	0.04	0.05	0.06	0.025
4	0.045	0.05	0.03	0.055	0.045	0.045	0.03

Table 3: Independence level experiment using uniform densities varying the sample size.

Sample size	FSIC	HSICAggInc $R = 1$	HSICAggInc $R = 100$	HSICAggInc $R = 200$	HSICAggCom
200	0.04	0.055	0.035	0.035	0.035
400	0.045	0.05	0.04	0.05	0.05
600	0.05	0.035	0.05	0.06	0.05
800	0.03	0.07	0.02	0.035	0.04
1000	0.07	0.02	0.085	0.035	0.04

Table 4: Independence level experiment using uniform densities varying the dimension.

Dimension	FSIC	HSICAggInc $R = 1$	HSICAggInc $R = 100$	HSICAggInc $R = 200$	HSICAggCom
2	0.035	0.065	0.08	0.055	0.07
3	0.065	0.055	0.035	0.02	0.025
4	0.04	0.035	0.045	0.055	0.055

Table 5: Goodness-of-fit level experiment using a Gaussian-Bernoulli Restricted Boltzmann Machine varying the sample size.

Sample size	FSSD	LSD	KSDAggInc $R = 1$	KSDAggInc $R = 100$	KSDAggInc $R = 200$	KSDAggCom
200	0.02	0.07	0.05	0.045	0.06	0.06
400	0.03	0.04	0.06	0.04	0.065	0.055
600	0.04	0.075	0.03	0.03	0.04	0.07
800	0.03	0.06	0.055	0.06	0.045	0.07
1000	0.025	0.05	0.045	0.035	0.045	0.065

Table 6: Goodness-of-fit level experiment using a Gaussian-Bernoulli Restricted Boltzmann Machine varying the dimension.

Dimension	FSSD	LSD	KSDAggInc $R = 1$	KSDAggInc $R = 100$	KSDAggInc $R = 200$	KSDAggCom
20	0.02	0.055	0.045	0.06	0.065	0.05
40	0.04	0.055	0.07	0.055	0.065	0.07
60	0.04	0.055	0.06	0.04	0.05	0.06
80	0.015	0.04	0.045	0.04	0.035	0.05

## C. Proof of Proposition 5.1

The asymptotic level of the goodness-of-fit test using a wild bootstrap follows from the results of [Shao \(2010\)](#), [Leucht and Neumann \(2013\)](#), [Chwialkowski et al. \(2014, 2016\)](#). As pointed out by [Schrab et al. \(2022\)](#) for the complete  $U$ -statistics, the KSD test statistic and the wild bootstrapped KSD statistics are not exchangeable under the null, and hence non-asymptotic level cannot be proved using the result of [Romano and Wolf \(2005, Lemma 1\)](#).

The non-asymptotic level for the two-sample test follows exactly from the reasoning of [Schrab et al. \(2021, Proposition 1\)](#). The fact that we work with incomplete  $U$ -statistics rather than with their complete counterparts does not affect the proof of exchangeability of  $U_\lambda^1, \dots, U_\lambda^{B_1+1}$ .

For the independence problem, [Albert et al. \(2022, Proposition 1\)](#) prove that the quadratic-time HSIC estimator and the permuted test statistics are exchangeable under the null hypothesis, it remains to be shown that this also holds in our incomplete setting using a wild bootstrap. Assuming that exchangeability under the null holds, the desired non-asymptotic level  $\alpha$  can then be guaranteed using the result of [Romano and Wolf \(2005, Lemma 1\)](#), exactly as done by [Albert et al. \(2022, Proposition 1\)](#).

We now prove that  $U_\lambda^1, \dots, U_\lambda^{B_1+1}$  for the independence problem are exchangeable under the null. Since  $U_\lambda^1, \dots, U_\lambda^{B_1}$  are i.i.d. given the data, they are exchangeable under the null. So, we need to prove that

$$\sum_{(i,j) \in \mathcal{D}_{\lfloor N/2 \rfloor}} h_{k,\ell}^{\text{HSIC}}(Z_i, Z_j, Z_{i+\lfloor N/2 \rfloor}, Z_{j+\lfloor N/2 \rfloor}) \quad (17)$$

is, under the null, distributed like

$$\sum_{(i,j) \in \mathcal{D}_{\lfloor N/2 \rfloor}} \epsilon_i \epsilon_j h_{k,\ell}^{\text{HSIC}}(Z_i, Z_j, Z_{i+\lfloor N/2 \rfloor}, Z_{j+\lfloor N/2 \rfloor}) \quad (18)$$

where  $\epsilon_1, \dots, \epsilon_N$  are i.i.d. Rademacher random variables. Using the result of [Schrab et al. \(2021, Appendix B, Proposition 11\)](#), considering the identity  $s_1(i) = i$  for  $i = 1, \dots, 2\lfloor N/2 \rfloor$ , and the swap function  $s_{-1}(i) = i + \lfloor N/2 \rfloor$  and  $s_{-1}(i + \lfloor N/2 \rfloor) = i$  for  $i = 1, \dots, 2\lfloor N/2 \rfloor$ , we have

$$\begin{aligned} & \sum_{(i,j) \in \mathcal{D}_{\lfloor N/2 \rfloor}} \epsilon_i \epsilon_j h_{k,\ell}^{\text{HSIC}}(Z_i, Z_j, Z_{i+\lfloor N/2 \rfloor}, Z_{j+\lfloor N/2 \rfloor}) \\ &= \frac{1}{4} \sum_{(i,j) \in \mathcal{D}_{\lfloor N/2 \rfloor}} h_k^{\text{MMD}}(X_i, X_j; X_{i+\lfloor N/2 \rfloor}, X_{j+\lfloor N/2 \rfloor}) \left( \epsilon_i \epsilon_j h_\ell^{\text{MMD}}(Y_i, Y_j; Y_{i+\lfloor N/2 \rfloor}, Y_{j+\lfloor N/2 \rfloor}) \right) \\ &= \frac{1}{4} \sum_{\substack{(i,j) \in \\ \mathcal{D}_{\lfloor N/2 \rfloor}}} h_k^{\text{MMD}}(X_i, X_j; X_{i+\lfloor N/2 \rfloor}, X_{j+\lfloor N/2 \rfloor}) h_\ell^{\text{MMD}}(Y_{s_{\epsilon_i}(i)}, Y_{s_{\epsilon_j}(j)}; Y_{s_{\epsilon_i}(i+\lfloor N/2 \rfloor)}, Y_{s_{\epsilon_j}(j+\lfloor N/2 \rfloor)}) \\ &= |\mathcal{D}_{\lfloor N/2 \rfloor}| \overline{\text{HSIC}}_{k,\ell}(\mathbb{Z}_N^\epsilon; \mathcal{D}_{\lfloor N/2 \rfloor}) \end{aligned}$$

where  $\mathbb{Z}_N^\epsilon := \left( (X_i, Y_{s_{\epsilon_i}(i)}) \right)_{1 \leq i \leq N}$  with  $Y_{s_{\epsilon_i}(i)} \in \{Y_i, Y_{i+\lfloor N/2 \rfloor}\}$  for  $i = 1, \dots, \lfloor N/2 \rfloor$  and  $Y_{s_{\epsilon_i}(i)} \in \{Y_i, Y_{i-\lfloor N/2 \rfloor}\}$  for  $i = 1 + \lfloor N/2 \rfloor, \dots, 2\lfloor N/2 \rfloor$ . Now, under the null the variables  $(X_i)_{1 \leq i \leq N}$  and  $(Y_i)_{1 \leq i \leq N}$  are independent, so  $\mathbb{Z}_N^\epsilon$  is distributed like  $\mathbb{Z}_N$ . We deduce that  $\overline{\text{HSIC}}_{k,\ell}(\mathbb{Z}_N; \mathcal{D}_{\lfloor N/2 \rfloor})$  and  $\overline{\text{HSIC}}_{k,\ell}(\mathbb{Z}_N^\epsilon; \mathcal{D}_{\lfloor N/2 \rfloor})$  have the same distribution under the null, and hence, that the terms in Equations (17) and (18) also have the same distribution under the null. We deduce that  $U_\lambda^1, \dots, U_\lambda^{B_1+1}$  for the independence problem are exchangeable under the null, which completes the proof.

## D. Proof of Lemma 4.1

Consider the case of fixed design. Using the variance expression of Lee (1990, Theorem 2, p. 190), we have

$$\text{var}(\bar{U}) = \frac{f_1\sigma_1^2 + f_2\sigma_2^2}{|\mathcal{D}|^2}$$

where  $f_i$  is the number of pairs of sets in the design  $\mathcal{D}$  that have  $i$  elements in common. The pairs of sets in  $\mathcal{D}$  with 2 elements in common are  $\{(\{i, j\}, \{i, j\}) : (i, j) \in \mathcal{D}\}$ , so  $f_2 = |\mathcal{D}|$ . We now calculate the number of pairs of sets in  $\mathcal{D}$  with 1 element in common. We start with a pair  $(i, j) \in \mathcal{D}$  (there are  $|\mathcal{D}|$  such pairs). The number of pairs in  $\mathcal{D}$  which have one element in common with  $(i, j)$  is upper bounded by the number of pairs in  $\mathbf{i}_2^N$  which have one element in common with  $(i, j)$ , those are  $\{\{i, r\} : 1 \leq r \leq N, r \neq i\} \cup \{\{j, r\} : 1 \leq r \leq N, r \neq j\}$  of size smaller than  $2N$ . We deduce that  $f_1 \leq 2N|\mathcal{D}|$ . Combining those results, we obtain

$$\text{var}(\bar{U}) \leq \frac{f_1\sigma_1^2 + f_2\sigma_2^2}{|\mathcal{D}|^2} \leq \frac{2N}{|\mathcal{D}|}\sigma_1^2 + \frac{1}{|\mathcal{D}|}\sigma_2^2$$

as desired.

Let us now consider the random design case. Recall that using the variance expression of the complete  $U$ -statistic  $U$  of Lee (1990, Theorem 3 p. 12), we can obtain that

$$\text{var}(U) \leq C \left( \frac{\sigma_1^2}{N} + \frac{\sigma_2^2}{N^2} \right)$$

as done by Kim et al. (2022, Appendix E) and Albert et al. (2022, Lemma 10). Using the result of Lee (1990, Theorem 4 p. 193), the variance of the incomplete  $U$ -statistic  $\bar{U}$  can be expressed in terms of the variance of the complete  $U$ -statistic  $U$ . For random design with replacement, we have

$$\text{var}(\bar{U}) = \frac{\sigma_2^2}{|\mathcal{D}|} + \left(1 - \frac{1}{|\mathcal{D}|}\right) \text{var}(U) \leq C \left( \frac{\sigma_1^2}{N} + \left(\frac{1}{|\mathcal{D}|} + \frac{1}{N^2}\right) \sigma_2^2 \right).$$

Letting  $S := N(N-1)/2$ , for random design without replacement, we have

$$\text{var}(\bar{U}) = \frac{S - |\mathcal{D}|}{|\mathcal{D}|(S-1)}\sigma_2^2 + \frac{S}{S-1} \left(1 - \frac{1}{|\mathcal{D}|}\right) \text{var}(U) \leq C \left( \frac{\sigma_1^2}{N} + \left(\frac{1}{|\mathcal{D}|} + \frac{1}{N^2}\right) \sigma_2^2 \right).$$

## E. Proof of Lemma 4.2

We rely on the concentration bound for i.i.d. Rademacher chaos of de la Peña and Giné (1999, Corollary 3.2.6) which as presented in Kim et al. (2022, Equation 39) takes the form

$$\mathbb{P}_\epsilon \left( \left| \sum_{(i,j) \in \mathbf{i}_2^N} \epsilon_i \epsilon_j a_{i,j} \right| \geq t \right) \leq 2 \exp \left( -Ct \left( \sum_{(i,j) \in \mathbf{i}_2^N} a_{i,j}^2 \right)^{-1} \right)$$

for some constant  $C > 0$  and for every  $t \geq 0$ , where  $\epsilon_1, \dots, \epsilon_N$  are i.i.d. Rademacher random variables taking values in  $\{-1, 1\}$ . Letting

$$a_{i,j} := \frac{h(Z_i, Z_j)}{|\mathcal{D}|} \mathbf{1}[(i, j) \in \mathcal{D}] \quad \text{for} \quad (i, j) \in \mathbf{i}_2^N,$$

we obtain

$$\begin{aligned} \mathbb{P}_\epsilon \left( \frac{1}{|\mathcal{D}|} \left| \sum_{(i,j) \in \mathcal{D}} \epsilon_i \epsilon_j h(Z_i, Z_j) \right| \geq t \mid \mathbb{Z}_N, \mathcal{D} \right) &\leq 2 \exp \left( -Ct \left( \frac{1}{|\mathcal{D}|^2} \sum_{(i,j) \in \mathcal{D}} h(Z_i, Z_j)^2 \right)^{-1} \right) \\ &\leq 2 \exp \left( -Ct \left( \frac{1}{|\mathcal{D}|^2} \sum_{(i,j) \in \mathfrak{I}_2^N} h(Z_i, Z_j)^2 \right)^{-1} \right) \end{aligned}$$

which concludes the proof.

## F. Proof of Theorem 5.2

We start by reviewing the steps of proofs of [Albert et al. \(2022\)](#) and [Schrab et al. \(2021, 2022\)](#) who prove that, for each of the three respective testing frameworks, a sufficient condition to ensure control of the probability of type II error for the quadratic-time test is

$$\|p - q\|_2^2 \geq \|(p - q) - T_\lambda(p - q)\|_2^2 + C \frac{1}{N} \frac{\ln(1/\alpha)}{\beta} \sigma_{2,\lambda}. \quad (19)$$

Those quadratic-time tests use the complete  $U$ -statistics defined in Equations (1), (3) and (5), which we denote as  $U_\lambda$ . The key results for their proofs rely on deriving variance and quantile bounds.

The variance bound is of the form

$$\text{var}(U_\lambda) \leq C \left( \frac{1}{N} \sigma_{1,\lambda}^2 + \frac{1}{N^2} \sigma_{2,\lambda}^2 \right) \quad (20)$$

where they show that for  $h_\lambda \in \{h_{k_\lambda}^{\text{MMD}}, h_{k_\lambda, \ell_\mu}^{\text{HSIC}}, h_{k_\lambda, p}^{\text{KSD}}\}$  defined in Equations (2), (4) and (6)

$$\sigma_{1,\lambda}^2 := \text{var}(\mathbb{E}[h_\lambda(Z, Z') \mid Z']) \leq C \|T_\lambda(p - q)\|_2^2$$

and

$$\sigma_{2,\lambda}^2 := \text{var}(h_\lambda(Z, Z')) = \mathbb{E}[h_\lambda(Z, Z')^2] \leq \frac{C}{\lambda_1 \cdots \lambda_d} \quad (21)$$

where the last inequality holds only for  $h_\lambda \in \{h_{k_\lambda}^{\text{MMD}}, h_{k_\lambda, \ell_\mu}^{\text{HSIC}}\}$ .

The quantile bound ([Schrab et al., 2021](#), Proposition 4) is of the form

$$\mathbb{P} \left( \widehat{q}_{1-\alpha}^{\lambda, U, B_1} \leq C \frac{1}{N} \frac{1}{\sqrt{\delta}} \ln \left( \frac{1}{\alpha} \right) \sigma_{2,\lambda} \right) \geq 1 - \delta$$

for  $\delta \in (0, 1)$ , where  $\widehat{q}_{1-\alpha}^{\lambda, U, B_1}$  is the quantile obtained using  $B_1$  wild bootstrapped similarly to the one defined in Equation (13) but using the complete  $U$ -statistic. Relying on Dvoretzky–Kiefer–Wolfowitz inequality ([Dvoretzky et al., 1956](#); [Massart, 1990](#)), [Schrab et al. \(2021, Proposition 4\)](#) show that it suffices to prove the bound

$$\mathbb{P} \left( \widehat{q}_{1-\alpha}^{\lambda, U, \infty} \leq C \frac{1}{N} \frac{1}{\sqrt{\delta}} \ln \left( \frac{1}{\alpha} \right) \sigma_{2,\lambda} \right) \geq 1 - \delta \quad (22)$$

for the true wild bootstrap quantile  $\widehat{q}_{1-\alpha}^{\lambda, U, \infty}$  without finite approximation.

Combining those variance and quantile bounds using Chebyshev's inequality (Chebyshev, 1899), they obtain a condition guaranteeing power in terms of the MMD, HSIC and KSD. By expressing these three measures as an RHKS inner product

$$\langle p - q, T_\lambda(p - q) \rangle = \frac{1}{2} \left( \|p - q\|_2^2 + \|T_\lambda(p - q)\|_2^2 - \|(p - q) - T_\lambda(p - q)\|_2^2 \right),$$

they obtain the condition in Equation (19) which guarantees high power in terms of  $\|p - q\|_2^2$ . Albert et al. (2022) and Schrab et al. (2021) then derive the minimax rate  $(1/N)^{2s/(4s+d)}$  over the Sobolev ball  $\mathcal{S}_d^s(R)$  for the independence and two-sample tests using the bandwidths  $\lambda_i^* := (1/N)^{2/(4s+d)}$  for  $i = 1, \dots, d$ .

For Theorem 5.2, we need to obtain the condition in Equation (19) with  $1/N$  replaced by  $N/L$ . Hence, following their reasoning, in order to prove Theorem 5.2 (i) & (ii), it suffices to derive variance and quantiles bounds for incomplete  $U$ -statistics which have the form of Equations (20) and (22) with  $1/N$  replaced by  $N/L$ , which we now do.

Using the variance bound for incomplete  $U$ -statistics  $\bar{U}_\lambda$  of Lemma 4.1, together with the fact that the design size  $L := |\mathcal{D}|$  is smaller than  $N^2$  so that  $1/L = L/L^2 \leq N^2/L^2$ , we obtain for fixed design that

$$\text{var}(\bar{U}_\lambda) \leq C \left( \frac{N}{L} \sigma_{1,\lambda}^2 + \frac{1}{L} \sigma_{2,\lambda}^2 \right) \leq C \left( \frac{N}{L} \sigma_{1,\lambda}^2 + \left( \frac{N}{L} \right)^2 \sigma_{2,\lambda}^2 \right),$$

we get the same bound for random design since

$$\text{var}(\bar{U}_\lambda) \leq C \left( \frac{1}{N} \sigma_{1,\lambda}^2 + \left( \frac{1}{L} + \frac{1}{N^2} \right) \sigma_{2,\lambda}^2 \right) \leq C \left( \frac{N}{L} \sigma_{1,\lambda}^2 + \left( \frac{N}{L} \right)^2 \sigma_{2,\lambda}^2 \right),$$

as desired.

For the quantile bound, we use Lemma 4.2 which, for  $A_\lambda^2 := L^{-2} \sum_{(i,j) \in \mathbf{i}_2^N} h_\lambda(Z_i, Z_j)^2$ , gives

$$\mathbb{P}_\epsilon(\bar{U}_\lambda^\epsilon \geq t \mid \mathcal{Z}_N, \mathcal{D}) \leq \mathbb{P}_\epsilon(|\bar{U}_\lambda^\epsilon| \geq t \mid \mathcal{Z}_N, \mathcal{D}) \leq 2 \exp\left(-C \frac{t}{A_\lambda}\right).$$

Setting  $\alpha := 2 \exp(-Ct/A_\lambda)$ , we obtain

$$\hat{q}_{1-\alpha}^{\lambda, \bar{U}, \infty} = t = \frac{A_\lambda \ln(2/\alpha)}{C} = CA_\lambda \ln(1/\alpha)$$

for a different constant  $C > 0$  since  $\alpha \in (0, e^{-1})$ . For  $\delta \in (0, 1)$ , using Markov's inequality, we obtain

$$\mathbb{P}\left(A_\lambda^2 \leq \frac{1}{\delta} \mathbb{E}[A_\lambda^2]\right) \geq 1 - \delta$$

where

$$\mathbb{E}[A_\lambda^2] = \mathbb{E}\left[\frac{1}{L^2} \sum_{(i,j) \in \mathbf{i}_2^N} h(Z_i, Z_j)^2\right] = \frac{N(N-1)}{L^2} \mathbb{E}[h_\lambda(Z, Z')^2] \leq C \frac{N^2}{L^2} \sigma_{2,\lambda}^2$$

using Equation (21). We deduce that

$$\begin{aligned} 1 - \delta &\leq \mathbb{P}\left(A_\lambda^2 \leq \frac{1}{\delta} \mathbb{E}[A_\lambda^2]\right) \\ &= \mathbb{P}\left(\hat{q}_{1-\alpha}^{\lambda, \bar{U}, \infty} \leq C \frac{1}{\sqrt{\delta}} \ln\left(\frac{1}{\alpha}\right) \sqrt{\mathbb{E}[A_\lambda^2]}\right) \\ &\leq \mathbb{P}\left(\hat{q}_{1-\alpha}^{\lambda, \bar{U}, \infty} \leq C \frac{1}{\sqrt{\delta}} \frac{N}{L} \ln\left(\frac{1}{\alpha}\right) \sigma_{2,\lambda}\right) \end{aligned}$$

as desired, which concludes the proof.

## G. Proof of Theorem 7.1

### G.1. Proof of Theorem 7.1 (i)

In this setting, we consider as proposed by [Kim et al. \(2022, Equation 32\)](#) the permuted HSIC complete  $U$ -statistic

$$U_N^\pi := \frac{1}{|\mathbf{i}_4^N|} \sum_{(i,j,r,s) \in \mathbf{i}_4^N} h_{k,\ell}^{\text{HSIC}}((X_i, Y_{\pi_i}), (X_j, Y_{\pi_j}), (X_r, Y_{\pi_r}), (X_s, Y_{\pi_s}))$$

for a permutation  $\pi$  of the indices  $\{1, \dots, N\}$ , and for  $h_{k,\ell}^{\text{HSIC}}$  as defined in Equation (4).

Applying the exponential concentration bound of [Kim et al. \(2022, Theorem 6.3\)](#), which uses the result of [de la Peña and Giné \(1999, Theorem 4.1.12\)](#), we obtain that there exist constants  $C_1, C_2 > 0$  such that

$$\mathbb{P}_\pi(U_N^\pi \geq t \mid \mathbb{Z}_N) \leq C_1 \exp\left(-C_2 \min\left(\frac{Nt}{\Lambda_N}, \frac{Nt^{2/3}}{M_N^{2/3}}\right)\right) \quad (23)$$

where

$$\Lambda_N^2 := \frac{1}{N^4} \sum_{i=1}^N \sum_{j=1}^N \sum_{r=1}^N \sum_{s=1}^N k_\lambda(X_i, X_j)^2 \ell_\mu(Y_r, Y_s)^2$$

and

$$\begin{aligned} M_N &:= \max_{1 \leq i,j,r,s \leq N} |k_\lambda(X_i, X_j) \ell_\mu(Y_r, Y_s)| \\ &= \max_{1 \leq i,j,r,s \leq N} \left| \prod_{a=1}^{d_x} \frac{1}{\lambda_a} K_a\left(\frac{(X_i)_a - (X_j)_a}{\lambda_a}\right) \prod_{b=1}^{d_y} \frac{1}{\lambda_b} L_b\left(\frac{(Y_r)_b - (Y_s)_b}{\lambda_b}\right) \right| \\ &\leq \frac{C}{\lambda_1 \dots \lambda_{d_x} \mu_1 \dots \mu_{d_y}} \\ &= \frac{C}{\lambda_1 \dots \lambda_d} \end{aligned}$$

for some constant  $C > 0$  since the functions  $K_1, \dots, K_{d_x}$  and  $L_1, \dots, L_{d_y}$  are bounded, and where we recall our notational convention that  $d := d_x + d_y$  and  $\lambda_{d_x+i} := \mu_i$  for  $i = 1, \dots, d_y$ .

Using the reasoning of [Schrab et al. \(2021, Proposition 3\)](#), we see that the results of [Albert et al. \(2022, Equations C.17, C.18 & C.19\)](#) hold not only for the Gaussian kernel but more generally for any kernels of the form of Equation (12). Those results give us that there exists a constant  $C > 0$  depending on  $M$  and  $d$  such that

$$\begin{aligned} \mathbb{E}[k_\lambda(X_1, X_2)^2 \ell_\mu(Y_1, Y_2)^2] &\leq \frac{C}{\lambda_1 \dots \lambda_{d_x} \mu_1 \dots \mu_{d_y}} = \frac{C}{\lambda_1 \dots \lambda_d}, \\ \mathbb{E}[k_\lambda(X_1, X_2)^2 \ell_\mu(Y_1, Y_3)^2] &\leq \frac{C}{\lambda_1 \dots \lambda_d}, \\ \mathbb{E}[k_\lambda(X_1, X_2)^2 \ell_\mu(Y_3, Y_4)^2] &\leq \frac{C}{\lambda_1 \dots \lambda_d}. \end{aligned}$$

We deduce that

$$\mathbb{E}[\Lambda_N^2] \leq \frac{C}{\lambda_1 \dots \lambda_d}.$$

As explained in Appendix F, by relying on Dvoretzky–Kiefer–Wolfowitz inequality ([Dvoretzky et al., 1956](#); [Massart, 1990](#)) as done by [Schrab et al. \(2021, Proposition 4\)](#), it



is sufficient to prove upper bounds for the true permutation quantile  $\widehat{q}_{1-\alpha}^{\lambda, \infty}$  without finite approximation. From Equation (23), we obtain that this quantile satisfies

$$\widehat{q}_{1-\alpha}^{\lambda, \infty} \leq C \max \left( \frac{\Lambda_N}{N} \ln \left( \frac{1}{\alpha} \right), \frac{M_N}{N^{3/2}} \ln \left( \frac{1}{\alpha} \right)^{3/2} \right).$$

Using Markov's inequality and bounds obtained above, we get that

$$\begin{aligned} \widehat{q}_{1-\alpha}^{\lambda, \infty} &\leq C \max \left( \frac{\sqrt{\mathbb{E}[\Lambda_N^2]}}{\sqrt{\delta}N} \ln \left( \frac{1}{\alpha} \right), \frac{M_N}{N^{3/2}} \ln \left( \frac{1}{\alpha} \right)^{3/2} \right) \\ &\leq C \max \left( \frac{\ln(1/\alpha)}{\sqrt{\delta}N\sqrt{\lambda_1 \cdots \lambda_d}}, \frac{\ln(1/\alpha)^{3/2}}{N^{3/2}\lambda_1 \cdots \lambda_d} \right) \end{aligned} \quad (24)$$

holds with probability at least  $1 - \delta$  where  $\delta \in (0, 1)$ , and where the constants are different on each line.

Now, recall that  $4s > d$  and  $\lambda_i^* = N^{-2/(4s+d)}$  for  $i = 1, \dots, d$ , so that

$$\frac{1}{\lambda_1^* \cdots \lambda_d^*} = N^{2d/(4s+d)} < N \quad \iff \quad N^{-1/2} < \sqrt{\lambda_1^* \cdots \lambda_d^*}$$

which gives

$$\begin{aligned} \widehat{q}_{1-\alpha}^{\lambda^*, \infty} &\leq C \max \left( \frac{\ln(1/\alpha)}{\sqrt{\delta}N\sqrt{\lambda_1^* \cdots \lambda_d^*}}, \frac{\ln(1/\alpha)^{3/2}}{N\sqrt{\lambda_1^* \cdots \lambda_d^*}} \right) \\ &\leq C \frac{\ln(1/\alpha)^{3/2}}{\sqrt{\delta}N\sqrt{\lambda_1^* \cdots \lambda_d^*}} \end{aligned}$$

holding with probability at least  $1 - \delta$ , since  $\alpha \in (0, e^{-1})$ . By combining this result with the reasoning of [Albert et al. \(2022\)](#) as explained in Appendix F, we obtain that the probability of type II error of the test is controlled by  $\beta \in (0, 1)$  when

$$\|p - q\|_2^2 \geq \|(p - q) - T_{\lambda^*}(p - q)\|_2^2 + C \frac{1}{N} \frac{\ln(1/\alpha)^{3/2}}{\beta \sqrt{\lambda_1^* \cdots \lambda_d^*}}.$$

We have recovered the correct dependency with respect to  $N$  and  $\lambda$  with an improved  $\alpha$ -dependency of  $\ln(1/\alpha)^{3/2}$  compared to the  $\alpha^{-1/2}$  dependency obtained by [Kim et al. \(2022, Proposition 8.7\)](#). The proof of minimax optimality of the quadratic-time test with fixed bandwidth  $\lambda^*$  does not depend on the  $\alpha$ -dependency and can be derived in both our setting and the one of [Kim et al. \(2022\)](#) using quantiles obtained from permutations by following the reasoning of [Albert et al. \(2022, Corollary 2\)](#). We obtain that the uniform separation rate over the Sobolev ball  $\mathcal{S}_d^s(R)$  is, up to a constant,  $(1/N)^{2s/(4s+d)}$ . The improved  $\alpha$ -dependency is crucial for deriving the rate of the aggregated quadratic-time test over Sobolev balls because the weights appear in the  $\alpha$ -term (i.e.  $\alpha$  is replaced by  $\alpha w_\lambda$  which depends on the sample size  $N$ ).

## G.2. Proof of Theorem 7.1 (ii)

Similarly to the proofs of [Albert et al. \(2022, Corollary 3\)](#) and [Schrab et al. \(2021, Corollary 10\)](#), consider

$$\ell^* := \left\lceil \frac{2}{4s+d} \log_2 \left( \frac{N}{\ln(\ln(N))} \right) \right\rceil \leq \left\lceil \frac{2}{d} \log_2 \left( \frac{N}{\ln(\ln(N))} \right) \right\rceil$$

and the bandwidth  $\lambda^* := (2^{-\ell^*}, \dots, 2^{-\ell^*}) \in \Lambda$  which satisfies

$$\ln\left(\frac{1}{w_{\lambda^*}}\right) \leq C \ln(\ell^*) \leq C \ln(\ln(N))$$

as  $w_{\lambda^*} := 6\pi^{-2}(\ell^*)^{-2}$ , and

$$\frac{1}{2} \left(\frac{\ln(\ln(N))}{N}\right)^{2/(4s+d)} \leq \lambda_i^* \leq \left(\frac{\ln(\ln(N))}{N}\right)^{2/(4s+d)}$$

for  $i = 1, \dots, d$ . Since  $4s > d$ , we have

$$\frac{1}{\lambda_1^* \cdots \lambda_d^*} \leq 2^d \left(\frac{N}{\ln(\ln(N))}\right)^{\frac{2d}{4s+d}} < C \frac{N}{\ln(\ln(N))} \iff N^{-1/2} < C \frac{\sqrt{\lambda_1^* \cdots \lambda_d^*}}{\sqrt{\ln(\ln(N))}},$$

all with different constants. By Equation (24), we get that

$$\widehat{q}_{1-\alpha w_{\lambda^*}}^{\lambda^*, \infty} \leq C \max\left(\frac{\ln(1/(\alpha w_{\lambda^*}))}{\sqrt{\delta} N \sqrt{\lambda_1^* \cdots \lambda_d^*}}, \frac{\ln(1/(\alpha w_{\lambda^*}))^{3/2}}{N^{3/2} \lambda_1^* \cdots \lambda_d^*}\right) \quad (25)$$

holds with probability at least  $1 - \delta$  for  $\delta \in (0, 1)$ . If the largest term is the first one, then we get

$$\begin{aligned} \widehat{q}_{1-\alpha w_{\lambda^*}}^{\lambda^*, \infty} &\leq C \frac{\ln(1/\alpha) + \ln(1/w_{\lambda^*})}{\sqrt{\delta} N \sqrt{\lambda_1^* \cdots \lambda_d^*}}, \\ \widehat{q}_{1-\alpha w_{\lambda^*}}^{\lambda^*, \infty} &\leq C \frac{\ln(\ln(N))}{\sqrt{\delta} N \sqrt{\lambda_1^* \cdots \lambda_d^*}}, \end{aligned} \quad (26)$$

where the constant also depends on  $\alpha$ . The result then follows exactly as in the proofs of [Albert et al. \(2022, Corollary 3\)](#) and [Schrab et al. \(2021, Corollary 10\)](#). So, we consider the case where the second term is the largest one, so that

$$\begin{aligned} \widehat{q}_{1-\alpha w_{\lambda^*}}^{\lambda^*, \infty} &\leq C \frac{\ln(1/(\alpha w_{\lambda^*}))^{3/2}}{N^{3/2} \lambda_1^* \cdots \lambda_d^*} \\ &\leq C \frac{\ln(1/w_{\lambda^*})^{3/2}}{N \lambda_1^* \cdots \lambda_d^*} N^{-1/2} \\ &\leq C \frac{\ln(\ln(N))^{3/2}}{N \lambda_1^* \cdots \lambda_d^*} \frac{\sqrt{\lambda_1^* \cdots \lambda_d^*}}{\sqrt{\ln(\ln(N))}} \\ &= C \frac{\ln(\ln(N))}{N \sqrt{\lambda_1^* \cdots \lambda_d^*}}. \end{aligned}$$

We have recovered the same dependency as in Equation (26) when considering the first term as the largest one, and the proof then follows exactly the ones of [Albert et al. \(2022, Corollary 3\)](#) and [Schrab et al. \(2021, Corollary 10\)](#). We have treated both cases in Equation (25), we conclude that the uniform separation rate over the Sobolev balls  $\{\mathcal{S}_d^s(R) : R > 0, s > d/4\}$  of the quadratic-time aggregated test using a quantile obtained with permutations is (up to a constant)

$$\left(\frac{\ln(\ln(N))}{N}\right)^{2s/(4s+d)}.$$