

An Efficient Permutation-Based Kernel Two-Sample Test

Antoine Chatalic^{1,2}, Marco Letizia^{1,3}, Nicolas Schreuder^{1,4}, and Lorenzo Rosasco^{1,5,6}

¹MaLGA Center - DIBRIS, Università di Genova, Genoa, Italy

²CNRS, Univ. Grenoble-Alpes, GIPSA-lab, France

³INFN - Sezione di Genova, Genoa, Italy

⁴CNRS, Laboratoire d'Informatique Gaspard Monge, Champs-sur-Marne, France

⁵Center for Brains, Minds and Machines, MIT, Cambridge, MA, USA

⁶Istituto Italiano di Tecnologia, Genoa, Italy

February 20, 2025

Abstract Two-sample hypothesis testing—determining whether two sets of data are drawn from the same distribution—is a fundamental problem in statistics and machine learning with broad scientific applications. In the context of nonparametric testing, maximum mean discrepancy (MMD) has gained popularity as a test statistic due to its flexibility and strong theoretical foundations. However, its use in large-scale scenarios is plagued by high computational costs. In this work, we use a Nyström approximation of the MMD to design a computationally efficient and practical testing algorithm while preserving statistical guarantees. Our main result is a finite-sample bound on the power of the proposed test for distributions that are sufficiently separated with respect to the MMD. The derived separation rate matches the known minimax optimal rate in this setting. We support our findings with a series of numerical experiments, emphasizing realistic scientific data.

Contents

1	Introduction	2
2	Background on testing	3
2.1	The two-sample testing problem	3
2.2	Kernel-based two-sample tests	4
2.3	Related work	5
3	An approximate MMD permutation test	6
3.1	Projection-based approximation of the MMD	6
3.2	Using permutations to determine the threshold	7
4	Theoretical guarantees	8
4.1	Level and power guarantees	9
4.2	Main elements for the proof	9
4.2.1	Quality of the MMD approximation	9

4.2.2	Bound on the empirical quantile threshold	11
5	Numerical studies	11
5.1	Datasets	12
5.2	Results	12
6	Conclusion	13
A	Quality of the MMD approximation: additional lemmas and proofs	18
B	Controlling the level and the power	19
B.1	Proof of Lemma 4.1	19
B.2	Proof of Lemma 4.3	19
B.3	Bounds on the quantiles	20
B.4	Proof of Theorem 4.2 (main result)	22

1 Introduction

Let P, Q be probability distributions over a space \mathcal{Z} . We consider the two-sample hypothesis testing problem, where one observes independent random samples

$$X_1, \dots, X_{n_X} \stackrel{i.i.d.}{\sim} P \quad \text{and} \quad Y_1, \dots, Y_{n_Y} \stackrel{i.i.d.}{\sim} Q, \quad (1)$$

and would like to test the hypothesis $\mathcal{H}_0 : P = Q$ against the alternative $\mathcal{H}_1 : P \neq Q$.

This problem is of paramount importance in many areas, particularly in precision sciences such as high-energy physics, where large volumes of multivariate data are commonly analyzed. A key application is the search for new physics, where small deviations between observed and predicted distributions may signal new phenomena (Chakravarti et al. 2023; D’Agnolo et al. 2021; Letizia et al. 2022). Another important example is the validation of generative models, which are emerging as promising surrogates to replace expensive Monte Carlo simulations while maintaining high fidelity (Amram et al. 2024; Grossi et al. 2024).

This paper focuses on kernel-based two-sample tests, which commonly use the maximum mean discrepancy (MMD; see Section 2.2) as a test statistic to determine whether to reject the null hypothesis \mathcal{H}_0 (Gretton et al. 2012). These tests are versatile and powerful tools for comparing probability distributions P and Q without requiring prior assumptions. However, they are often hindered by their computational cost, which scales quadratically with the total sample size $n := n_X + n_Y$.

In order to mitigate this computational drawback, several approximations of the maximum mean discrepancy (MMD) have been explored, each offering a trade-off between efficiency and accuracy. Notable approaches include linear-MMD (Gretton et al. 2012), incomplete U-statistics (Schrab et al. 2022; Yamada et al. 2018), block-MMD (Zaremba et al. 2013). A key limitation of these approaches is that quadratic time complexity is necessary to achieve statistical optimality (Domingo-Enrich et al. 2023, Proposition 2). More details on related works will be provided in Section 2.

To address this limitation, we propose a new procedure for two-sample testing based on a Nyström approximation of the MMD (Chatalic et al. 2022; Nyström 1930; Williams et al. 2001). From a theoretical perspective, we put forward a setting in which we demonstrate the minimax optimality (Baraud 2002) of our procedure while maintaining sub-quadratic computational complexity—a property shared with random feature approximations (Choi et al. 2024a; Zhao et al. 2015) and coresot-based

\mathcal{Z}	Data space
$\hat{\Psi}(X, Y)$	Test statistic
κ, ϕ	(Base) kernel and associated feature map
$\tilde{\kappa}, \varphi$	Feature map used to build the test statistic and associated kernel
$\hat{\Psi}_{(b_\alpha)}$	Threshold of the test
$n_X = X , n_Y = Y $	Number of samples (from P , from Q)
$n = n_X + n_Y$	Total number of samples
ℓ	dimension of the feature map ($\ell = \ell_X + \ell_Y$ for Nyström)
ℓ_X, ℓ_Y	Number of landmarks (from X , from Y)
\mathcal{P}	Number of permutations

Table 1: Main notations used throughout the paper

methods (Domingo-Enrich et al. 2023). Moreover, our method is simple to implement, with its approximation quality controlled by a single hyperparameter, making it both computationally efficient and easy to use in practice. Additionally, it applies to a broad class of kernels and does not require the input space to be a Euclidean space.

Organization of the paper We begin by reviewing the two-sample testing problem and kernel-based two-sample tests in Section 2. Section 3 introduces our MMD-based permutation test. Theoretical guarantees for this test are provided in Section 4. Finally, numerical studies presented in Section 5 demonstrate the practical effectiveness of our method.

Table 1 contains the main notations used throughout the paper.

2 Background on testing

In this section, we formally introduce the two-sample testing problem and provide an overview of MMD-based procedures to address it.

2.1 The two-sample testing problem

Let $(\mathcal{Z}, \mathcal{M})$ be a measurable space¹. Let $\mathcal{P}(\mathcal{Z})$ be the space of probability measures on \mathcal{Z} , and consider two probability distributions $P, Q \in \mathcal{P}(\mathcal{Z})$. Recall that we are given i.i.d. random samples $X_1, \dots, X_{n_X} \sim P$ and $Y_1, \dots, Y_{n_Y} \sim Q$ as described in Eq. (1). The goal of the two-sample testing problem is to test whether these samples are drawn from the same distribution. Formally, we aim to test the null hypothesis $\mathcal{H}_0 : P = Q$ against the alternative hypothesis $\mathcal{H}_1 : P \neq Q$.

A test is defined as a function of the data $T_{n_X, n_Y} : \mathcal{Z}^{n_X} \times \mathcal{Z}^{n_Y} \rightarrow \{0, 1\}$, designed to distinguish between the null and alternative hypotheses. Specifically, the test rejects the null hypothesis if and only if $T_{n_X, n_Y}(X, Y) = 1$. Most tests are built by considering a **test statistic** Ψ measuring some discrepancy between X and Y , and reject the null if and only if $\Psi(X, Y) > t$ for some threshold t , which can also depend on X, Y and whose choice is part of the definition of the test. To evaluate the effectiveness of such a test, we consider the errors it may incur in distinguishing between the null and alternative hypotheses.

Definition 1 (Type I/II errors) *A type I error occurs when the null hypothesis $P = Q$ is incorrectly rejected. Conversely, a type II error happens when the null hypothesis $P = Q$ is not rejected, despite $P \neq Q$*

¹One can consider for simplicity $\mathcal{Z} \subseteq \mathbb{R}^d$ with the Borel σ -algebra, however our results hold more generally for locally compact second-countable topological spaces.

being different.

Using these definitions, we can characterize the performance of a test in terms of its level and power.

Definition 2 (Level, power) A test T_{n_X, n_Y} is said to have level α if its type I error is uniformly bounded by α , and power $1 - \beta$ for a class of alternative distributions $\mathcal{P}_1 \subseteq \mathcal{P}(\mathcal{Z})^2$ if its type II error is uniformly bounded by β :

$$\sup_{P=Q \in \mathcal{P}(\mathcal{Z})} \mathbb{P}_{X \sim P^{\otimes n_X}, Y \sim Q^{\otimes n_Y}} [T_{n_X, n_Y}(X, Y) = 1] \leq \alpha, \quad (2)$$

$$\sup_{(P, Q) \in \mathcal{P}_1} \mathbb{P}_{X \sim P^{\otimes n_X}, Y \sim Q^{\otimes n_Y}} [T_{n_X, n_Y}(X, Y) = 0] \leq \beta. \quad (3)$$

Here the probabilities are taken not only over X and Y , but also over all other sources of randomness.

Beyond level and power, the **uniform separation rate** $\rho(T, \beta, \mathcal{P}_1)$ (Baraud 2002) of a test T provides a finer characterization of a test's performance. It quantifies the smallest separation between distributions that the test T can reliably distinguish and is defined as

$$\inf \left\{ \rho > 0 \mid \sup_{(P, Q) \in \mathcal{P}_1(\rho)} \mathbb{P}_{X \sim P^{\otimes n_X}, Y \sim Q^{\otimes n_Y}} [T(X, Y) = 0] \leq \beta \right\}.$$

Here $\rho > 0$ represents the separation, and $\mathcal{P}_1(\rho) := \{(P, Q) \in \mathcal{P}(\mathcal{Z})^2 : d(P, Q) > \rho\}$ denotes the class of distinct probability distributions P and Q that are separated by at least ρ with respect to a specified metric d (e.g., the MMD). The uniform separation rate determines the smallest separation ρ such that the test T reliably distinguishes ρ -separated distributions P and Q up to failure probability β .

Finally, for a fixed (non-asymptotic) level α , the **minimax rate of testing** is defined as the smallest uniform separation rate achievable by level- α tests:

$$\underline{\rho}(n_X, n_Y, \alpha, \beta, \mathcal{P}_1) := \inf_{T_\alpha} \rho(T_\alpha, \beta, \mathcal{P}_1), \quad (4)$$

where the infimum is taken over all level- α tests T_α .

Minimax rates for two-sample testing have been studied under different metrics and various classes of alternative hypotheses, see for instance Li et al. (2019). In this paper, we focus on distributions separated with respect to the MMD metric. For this setting, the minimax rate of testing is known to be lower-bounded by $\log(1/(\alpha + \beta))^{1/2} n^{-1/2}$ for translation-invariant kernels on \mathbb{R}^d (Kim et al. 2024, Th. 8). In Section 4, we will demonstrate that our proposed test achieves this minimax rate, establishing its optimality with respect to the MMD metric. Before introducing our testing procedure, we now provide a review of kernel-based two-sample tests and related works.

2.2 Kernel-based two-sample tests

Let $\delta(\cdot)$ denote the Dirac measure. We define the empirical probability distributions associated with the samples X and Y as $\hat{P} = \frac{1}{n_X} \sum_{i=1}^{n_X} \delta(X_i)$ and $\hat{Q} = \frac{1}{n_Y} \sum_{i=1}^{n_Y} \delta(Y_i)$, respectively. A common approach for performing a two-sample test is to assess whether a specific metric between \hat{P} and \hat{Q} exceeds a predefined threshold, which typically depends on the sample sizes. In the context of kernel methods, a standard choice of metric is the maximum mean discrepancy between \hat{P} and \hat{Q} , which we now introduce.

Let \mathcal{H} be a reproducing kernel Hilbert space (RKHS) with reproducing kernel $\kappa : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$ and canonical feature map $\phi(x) := \kappa(x, \cdot)$. We will impose the following assumption on the kernel κ .

Assumption 2.1: The kernel is bounded with $\sup_{x \in \mathcal{Z}} \kappa(x, x) = \|\kappa\|_\infty < \infty$ and measurable.

It is a common assumption in the kernel testing literature (see, e.g., Choi et al. (2024a) and Domingo-Enrich et al. (2023)). The kernel mean embedding (Berlinet et al. 2004) of a probability distribution π is defined as the Bochner integral (Diestel et al. 1977)

$$\mu(\pi) := \int \phi(x) d\pi(x),$$

and is well-defined under Assumption 2.1. Kernel mean embeddings allow probability distributions over arbitrary spaces to be represented as points in a Hilbert space. The maximum mean discrepancy (MMD) between two probability distributions P and Q is then defined as the distance between their respective kernel mean embeddings,

$$\text{MMD}(\pi_1, \pi_2) := \|\mu(\pi_1) - \mu(\pi_2)\|. \quad (5)$$

A kernel is said to be characteristic (Fukumizu et al. 2007) if and only if the mapping μ is injective, i.e. $P = Q \iff \|\mu(P) - \mu(Q)\| = 0$, in which case the MMD defines a metric on $\mathcal{P}(\mathcal{Z})$. Examples of characteristic kernels include Gaussian, Laplace and Matérn kernels. For general conditions under which kernels are characteristic, see Sriperumbudur et al. (2010), Simon-Gabriel et al. (2018) and related references.

Kernel two-sample testing is based on the principle that the MMD between two samples drawn from the same distribution should be small. The null hypothesis is rejected if the MMD (or a related statistic) exceeds a predefined threshold, indicating that the two samples are likely drawn from different distributions. Existing approaches mainly differ in their choice of test statistic and the method used to determine the test threshold.

2.3 Related work

Hypothesis testing and two-sample testing have been widely studied for a long time, and we refer the reader to Lehmann et al. (2022) for a general introduction.

Kernel-based test The introduction of two-sample tests using the MMD and its unbiased estimators as test statistics is due to Gretton et al. (2007, 2012). Based on either large deviation bounds or the asymptotic distribution of the unbiased test statistic, the authors derive test threshold values to achieve a target significance level α . Following this work, many variants have been proposed. For example, to address the issue that the standard kernel-MMD test statistic is a degenerate U-statistic under the null hypothesis, making its limiting distribution intractable, Shekhar et al. (2022) introduced cross-MMD. This quadratic-time MMD test statistic uses sample-splitting and studentization to ensure a limiting standard Gaussian distribution under the null. Departing from the MMD, other kernel-based test statistics have been explored. Since the MMD is an integral probability metric (Müller 1997), MMD-based tests can be interpreted as identifying the most discriminative test function from a set of witness functions belonging to a reproducing kernel Hilbert space. Inspired by this interpretation, tests based on optimized witness functions have been proposed by Kübler et al. (2022a) and Kübler et al. (2022b). Other kernel-based metrics include kernel Fisher discriminant analysis (Harchaoui et al. 2008) and its regularized variants (Hagrass et al. 2023), which can be viewed as kernelized versions of Hotelling’s T^2 test statistic. Additional approaches include the kernel density ratio (Kanamori et al. 2011) and kernel Stein discrepancies for goodness-of-fit tests (Huggins et al. 2018; Kalinke et al. 2024).

Efficient kernel-based tests The main disadvantage of kernel-based tests is that computing the MMD scales quadratically with the number of samples n . In their seminal paper, Gretton et al. (2012, Section

6) already introduced the linear-MMD, a statistic computable in $O((n+m)d)$ time, leveraging a partial evaluation of the terms appearing in the U-statistic estimator of the squared MMD. Variants of these incomplete U-statistics have subsequently been proposed by Schrab et al. (2022) and Yamada et al. (2018). Considering another partial evaluation of the MMD, Zaremba et al. (2013) introduced the block-MMD, a test statistic derived by (i) splitting the observations into disjoint blocks, (ii) computing the MMD for each block, and (iii) averaging the resulting statistics across all blocks. This approach has been further analyzed and refined by Ramdas et al. (2015) and Reddi et al. (2015). Chwialkowski et al. (2015) introduced a linear-time test statistic based on the average squared distance between empirical kernel embeddings evaluated at J randomly drawn points. Jitkrittum et al. (2016) proposed a variant of this statistic in which the J points are selected to maximize a lower bound on the power. A major limitation of these approaches is that either a quadratic time complexity is necessary to achieve an optimal power (Domingo-Enrich et al. 2023, Proposition 2) or the computation/power trade-off is yet to be characterized. Coreset-based approximation strategies have also been investigated, and proven to reach minimax separation rates at a subquadratic computational cost (Domingo-Enrich et al. 2023). Approximations of the MMD based on random Fourier features have been explored by Zhao et al. (2015), Zhao et al. (2021), and more recently by Choi et al. (2024a). Finally, Chatalic et al. (2022) proposed a Nyström approximation of the MMD, which is the base of our method and study.

Permutation tests There are several popular approaches to determining the test threshold. One common method is to use the quantile of the asymptotic distribution of the test statistic under the null hypothesis. However, this approach provides only asymptotic guarantees and may not perform well in finite samples. Another method relies on concentration inequalities, which, while theoretically sound, can be overly conservative, leading to thresholds that are too loose. Alternatively, permutation and bootstrap methods offer data-driven approaches that approximate the null distribution more accurately in practice, often resulting in improved empirical performance. The idea of comparing test statistics to their permuted replications dates back to the work of Hoeffding (1952). More recently, the combination of the MMD test statistic with permutation and bootstrap approaches has been explored in works such as Fromont et al. (2012), Kim et al. (2022), and Schrab et al. (2023).

Parameters selection The kernel function and its hyperparameters play an important role in the practical usability of kernel-based tests. Balasubramanian et al. (2021) and Li et al. (2024) established minimax optimality of carefully tuned MMD-based tests in specific settings. Multiple approaches have been investigated, such as aggregating multiple tests (Biggs et al. 2023; Fromont et al. 2012; Schrab et al. 2023, 2022) or using a Bayesian formalism (Zhang et al. 2022). In this work, we focus on the computational efficiency of the test, yet our approach could easily be combined with such ideas when adaptation is needed.

3 An approximate MMD permutation test

In order to avoid the quadratic cost of computing the MMD, we now introduce an efficient randomized approximation and employ it as a test statistic.

3.1 Projection-based approximation of the MMD

We approximate the MMD (5) between two empirical distributions \hat{P} and \hat{Q} using the Nyström method (Nyström 1930; Williams et al. 2001), that is by projecting their kernel mean embeddings of $\mu(\hat{P})$ and $\mu(\hat{Q})$ onto a data-dependent finite-dimensional subspace. Formally, given a set of landmark points $z_1, \dots, z_\ell \in \mathcal{Z}$, we define $\Phi_Z = [\phi(z_1), \dots, \phi(z_\ell)] : \mathbb{R}^\ell \rightarrow \mathcal{H}$ and $P_Z : \mathcal{H} \rightarrow \mathcal{H}$ the orthogonal projector onto $\text{span}(\Phi_Z)$. Using this projection, we approximate $\text{MMD}(\hat{P}, \hat{Q})$ by

$$\hat{\Psi}(X, Y) := \|P_Z \mu(\hat{P}) - P_Z \mu(\hat{Q})\|. \quad (6)$$

Denoting $\Phi_Z^* : h \mapsto [\kappa(h, z_1), \dots, \kappa(h, z_\ell)]^T$ the adjoint of the operator Φ_Z , the projection P_Z can be expressed as $P_Z = \Phi_Z(\Phi_Z^* \Phi_Z)^{\dagger} \Phi_Z^*$ and satisfies, as a projection, $P_Z = P_Z^2 = P_Z^*$. We also stress that $\Phi_Z^* \Phi_Z = K_Z$ corresponds to the kernel matrix K_ℓ of the chosen landmarks, i.e. defined by $(K_Z)_{ij} = \kappa(z_i, z_j)$. Hence, for any $v \in \mathcal{H}$ it holds $\|P_Z v\|^2 = \langle v, P_Z^2 v \rangle = \langle v, P_Z v \rangle = \|(K_Z^{\dagger})^{1/2} \Phi_Z^* v\|$. In particular, this implies that the approximation of the MMD in Eq. (6) can be computed efficiently as

$$\hat{\Psi}(X, Y) = \left\| \frac{1}{n_X} \sum_{i=1}^{n_X} \varphi(x_i) - \frac{1}{n_Y} \sum_{j=1}^{n_Y} \varphi(y_j) \right\| \quad (7)$$

$$\text{where } \varphi(x) := (K_Z^{\dagger})^{1/2} \begin{bmatrix} \kappa(z_1, x) \\ \dots \\ \kappa(z_\ell, x) \end{bmatrix}, \quad (8)$$

and K_Z^{\dagger} denotes the Moore-Penrose pseudo-inverse of K_Z . Note that $\hat{\Psi}(X, Y)$ can be computed in one pass over the data once Z has been chosen. We now detail how the landmarks are selected.

Choice of the landmarks Following Chatalic et al. (2022), we build our approximation by sampling points from both X and Y in order to obtain a good approximation of both kernel mean embeddings simultaneously. More precisely, we sample ℓ_X landmarks $\tilde{x}_1, \dots, \tilde{x}_{\ell_X}$ from X and ℓ_Y landmarks $\tilde{y}_1, \dots, \tilde{y}_{\ell_Y}$ from Y , and define $z_1 = \tilde{x}_1, \dots, z_{\ell_X} := \tilde{x}_{\ell_X}, z_{\ell_X+1} := \tilde{y}_1, \dots, z_\ell := \tilde{y}_{\ell_Y}$, so that $\ell = \ell_X + \ell_Y$.

Theoretical guarantees for approximating the MMD using a Nyström approximation with uniform sampling of the landmarks have been established in Chatalic et al. (2022). However, achieving an optimal convergence rate of order $O(n^{-1/2})$ with uniform sampling requires using a large number of landmarks. In this paper, we instead rely on leverage scores sampling, which is recognized as an optimal sampling strategy for compression (Chatalic et al. 2023).

Leverage scores quantify the relative importance of each point in a dataset and are closely related to the inverse of the Christoffel function in approximation theory (Fanuel et al. 2022; Pauwels et al. 2018). In this work, we focus on kernel ridge leverage scores (KRLS) (Alaoui et al. 2015), which are defined for a dataset of size n with an associated kernel matrix K as

$$\ell_\lambda(i) := \left(K(K + \lambda n I)^{-1} \right)_{ii}, \quad i = 1, \dots, n, \quad (9)$$

where $\lambda > 0$ is a regularization parameter. In the following, we sample landmarks from X and Y using leverage scores computed separately for each dataset. This corresponds to applying the definition in Eq. (9) to the kernel matrices of sizes $n_X \times n_X$ and $n_Y \times n_Y$, respectively. Since computing these scores exactly is typically expensive, we consider using multiplicative approximations of these scores.

Definition 3.1 (AKRLS): Let $\delta \in (0, 1]$, $\lambda_0 > 0$ and $z \in [1, \infty)$. The scores $(\hat{\ell}_\lambda(i))_{i \in [n]} \in \mathbb{R}^n$ are said to be (z, λ_0, δ) -approximate kernel ridge leverage scores (AKRLS) of X if with probability at least $1 - \delta$, for all $\lambda \geq \lambda_0, i \in [1, \dots, n]$ it holds $\frac{1}{z} \ell_\lambda(i) \leq \hat{\ell}_\lambda(i) \leq z \ell_\lambda(i)$.

Efficient algorithms have been proposed in the literature to efficiently sample from such approximate kernel ridge leverage scores, see for instance Chen et al. (2021), Musco et al. (2017), and Rudi et al. (2018).

3.2 Using permutations to determine the threshold

To construct a test based on the test statistic defined in Eq. (6), we need to specify how the threshold of the test is chosen. In this work, we adopt a permutation-based approach (Lehmann et al. 2022, Chapter 17). It exploits the fact that, under the null hypothesis (i.e., when $P = Q$), all observed data

points are exchangeable. Consequently, permuting the samples does not change the distribution of the test statistic.

We introduce the set of concatenated datasets $W = (W_i)_{1 \leq i \leq n}$, defined as $W_i = X_i$ for $1 \leq i \leq n_X$ and $W_{n_X+i} = Y_i$ for $1 \leq i \leq n_Y$. In what follows we consider a random variable σ that is uniformly distributed over the set of all permutations of $\{1, \dots, n\}$ and that is independent of all other sources of randomness. By sampling \mathcal{P} i.i.d. permutations $(\sigma_p)_{p=1}^{\mathcal{P}}$ and taking $\sigma_0 = \text{Id}$, we compute $(\mathcal{P} + 1)$ permuted test statistic

$$\hat{\Psi}_p := \hat{\Psi}((W_{\sigma_p(i)})_{i=1}^{n_X}, (W_{\sigma_p(j)})_{j=n_X+1}^n), \quad 0 \leq p \leq \mathcal{P}.$$

We use the notation $\hat{\Psi}_{(i)}$ to refer to the i -th value of the *ordered* list of these statistics, i.e. $\hat{\Psi}_{(0)} \leq \dots \leq \hat{\Psi}_{(\mathcal{P})}$. The threshold is then set as the empirical quantile $\hat{\Psi}_{(b_\alpha)}$ of the permuted test statistic, where $b_\alpha := \lceil (1 - \alpha)(\mathcal{P} + 1) \rceil$. Our testing procedure is detailed in Algorithm 3.1.

Algorithm 3.1: Permutation test based on a Nyström approximation of the MMD

Input: Feature map φ as in (8), data $W = (W_i)_{1 \leq i \leq n} \in \mathcal{Z}^n$, level $\alpha \in (0, 1)$, number of permutations \mathcal{P}

Output: Result of the test (boolean)

$w^{(0)} \leftarrow \left[\frac{1}{n_X}, \dots, \frac{1}{n_X}, -\frac{1}{n_Y}, \dots, -\frac{1}{n_Y} \right] \in \mathbb{R}^n$;

foreach $p = 1, \dots, \mathcal{P}$ **do** $w^{(p)} \leftarrow \text{shuffle}(w)$;

foreach $p = 0, \dots, \mathcal{P}$ **do** $v^{(p)} \leftarrow 0 \in \mathbb{R}^\ell$;

foreach $i = 1, \dots, n$ **do** // $O(n\mathcal{P}\ell)$ time, $O(\mathcal{P}\ell)$ space

foreach $p = 0, \dots, \mathcal{P}$ **do** $v^{(p)} \leftarrow v^{(p)} + w_i^{(p)} \varphi(W_i)$;

foreach $p = 0, \dots, \mathcal{P}$ **do** $\hat{\Psi}_p \leftarrow \left\| v^{(p)} \right\|$;

$b_\alpha \leftarrow \lceil (1 - \alpha)(\mathcal{P} + 1) \rceil$;

if $\hat{\Psi}_0 > \hat{\Psi}_{(b_\alpha)}$ **then** // reject H_0

| return 1 ;

else if $\hat{\Psi}_0 = \hat{\Psi}_{(b_\alpha)}$ **then**

$\hat{\Psi}^> \leftarrow \#\{0 \leq b \leq \mathcal{P} : \hat{\Psi}_b > \hat{\Psi}_{(b_\alpha)}\}$;

$\hat{\Psi}^= \leftarrow \#\{0 \leq b \leq \mathcal{P} : \hat{\Psi}_b = \hat{\Psi}_{(b_\alpha)}\}$;

return 1 with probability $\frac{\alpha(\mathcal{P}+1) - \hat{\Psi}^>}{\hat{\Psi}^=}$;

else

| return 0 ; // fail to reject H_0

Computational and memory efficiency The space and time complexity of our algorithm is limited. Once the landmarks have been chosen, all test statistics are computed in a single pass over the data. Moreover, the \mathcal{P} permutations can be efficiently stored using $O(\mathcal{P}n)$ bits and only $O((\mathcal{P} + 1)m)$ space is required to store the mean embeddings corresponding to all permutations and the non-permuted test statistic. In particular, it is never needed to store the features of the whole dataset in memory. This enables us to easily permute all samples without resorting to batch compression strategies that are typically required in coresets-based methods (Domingo-Enrich et al. 2023).

4 Theoretical guarantees

In this section, we focus on controlling the level and power of the proposed test. Without loss of generality, we assume $r \leq n_X/n_Y \leq 1$.

4.1 Level and power guarantees

The level of the proposed test is established by the following lemma, while a bound on the power is provided in the rest of this section. This lemma holds under the assumption that the data are exchangeable under the null hypothesis. It is satisfied in particular when the data are i.i.d., as is the case here.

Lemma 4.1: For any $\alpha \in (0, 1)$, the test described in Algorithm 3.1 with input level α has exact level α . (→ Proof)

By “exact”, we mean here that the inequality in (2) is actually an equality. We now state our main result which bounds the power of the test, and whose proof is based on lemmas 4.3, 4.5 and 4.6.

Theorem 4.2 (Main result: power of the test): Let $\beta \in (0, 1)$. Assume that the hypotheses of Lemma 4.5 are satisfied for $\delta := \beta/2$. Let $c_\alpha = \lfloor \alpha(\mathcal{P} + 1) \rfloor$. Then there exists a universal constant c s.t. the test of Algorithm 3.1 has power at least $1 - \beta$ provided that

$$\text{MMD}(P, Q)^2 \geq \frac{c \|\kappa\|_\infty}{n_X} \left[\frac{(r+1)}{r} \log \left(\frac{2e}{\alpha} \left(\frac{2}{\beta} \right)^{\frac{1}{c_\alpha}} \right) + \log \left(\frac{n_X}{\beta} \right) \right].$$

(→ Proof)

A detailed bound, with constants that are made comparatively more explicit is provided in the proof of Theorem 4.2. Note that the dependence of the separation rate (with respect to the MMD) on the smallest sample size n_X is $n_X^{-1/2}$. This matches the known minimax optimal rate of testing in this setting as established in Section 2.1 (Kim et al. 2024, Th. 8).

4.2 Main elements for the proof

Our proof relies on the following lemma, which is similar in spirit to Schrab et al. (2023, Lemma 4).

Lemma 4.3: Let $0 < \beta \leq 1$. Assume that there exists a function \mathcal{E}_{MMD} such that our estimator of the MMD satisfies

$$\mathbb{P}[|\text{MMD}(P, Q) - \hat{\Psi}(X, Y)| \geq \mathcal{E}_{\text{MMD}}(n_X, n_Y, \beta/2)] \leq \beta/2. \quad (10)$$

If the distributions P and Q are such that

$$\mathbb{P}[\text{MMD}(P, Q) \geq \mathcal{E}_{\text{MMD}}(n_X, n_Y, \beta/2) + \hat{\Psi}_{(b_\alpha)}] > 1 - \beta/2,$$

then $\mathbb{P}[\hat{\Psi}(X, Y) \leq \hat{\Psi}_{(b_\alpha)}] \leq \beta$.

(→ Proof)

In Section 4.2.1, we will show that Assumption (10) is satisfied for the proposed Nyström-based estimator. Then, following Domingo-Enrich et al. (2023, Lemma 6) we will show in Section 4.2.2 how the (random) empirical quantile threshold $\hat{\Psi}_{(b_\alpha)}$ can be bounded by the (deterministic) quantile of the permuted test statistic conditionally on X, Y .

4.2.1 Quality of the MMD approximation

The test statistic defined in Equation (6) is an estimator of the MMD between P and Q . In this section, we provide a high-probability bound for this approximation. This result is of independent interest,

as it improves upon Chatalic et al. (2022) by using approximate leverage score sampling instead of uniform sampling.

For $\square \in \{P, Q\}$, let $C_\square := \mathbb{E}_{x \sim \square} \phi(x) \otimes \phi(x) : \mathcal{H} \rightarrow \mathcal{H}$ denote the (uncentered) covariance operators associated with P and Q . Both C_P and C_Q are self-adjoint trace-class operators under Assumption 2.1, and we denote $(\lambda_i(C))_{i \in \mathbb{N}^*}$ the i -th eigenvalue of C . We will make the following assumption on the decay of the spectra of the considered covariance operators.

Assumption 4.4 (Polynomial spectral decay): There exist $\gamma \in (0, 1]$ and $a_\gamma > 0$ such that $\max(\lambda_i(C_P), \lambda_i(C_Q)) \leq a_\gamma i^{-1/\gamma}$.

Assuming that a covariance operator has such a polynomial decay is relatively standard in the literature, and corresponds for instance to the case of Sobolev spaces as studied by Widom (1964).

A key quantity in the analysis is the so-called effective dimension. This quantity depends both on the choice of the kernel as well as on a probability distribution, and is defined for any $\lambda > 0$ as $d_{\text{eff}}(P, \lambda) := \text{tr}(C_P(C_P + \lambda I)^{-1}) = \sum_{i \in \mathbb{N}^*} \frac{\lambda_i(C_P)}{\lambda_i(C_P) + \lambda}$. Notably, the effective dimension depends on the distribution P only via its covariance operator, and can be interpreted as a smooth estimate of the number of eigenvalues of C_P that are greater than λ . Under Assumptions 2.1 and 4.4, there exists a constant c_γ depending on a_γ, γ and $\|\kappa\|_\infty$ such that it holds $\max(d_{\text{eff}}(P, \lambda), d_{\text{eff}}(Q, \lambda)) \leq c_\gamma \lambda^{-\gamma}$ for any $\lambda > 0$ (see Chatalic et al. 2023, Lemma F.1).

We now have the following result regarding the error induced by the Nyström approximation of the MMD.

Lemma 4.5 (Nyström MMD approximation): Let $\ell = \ell_X + \ell_Y, \ell_X \geq L(n_X), \ell_Y \geq L(n_Y)$ where

$$L(n) := n^\gamma \left(\log \frac{32n}{\delta} \right)^{1-\gamma} \frac{78c_\gamma z^2}{(19\|\kappa\|_\infty)^\gamma}. \quad (11)$$

Let $\tilde{x}_1, \dots, \tilde{x}_{\ell_X}$ and $\tilde{y}_1, \dots, \tilde{y}_{\ell_Y}$ be drawn with replacement respectively from the datasets X and Y , and respectively proportionally to $(z, \lambda_0, \delta/8)$ -AKRLS for X and $(z, \lambda_0, \delta/8)$ -AKRLS for Y for some $z \geq 1$ and $\lambda_0 > 0$ satisfying $\lambda_0 \leq 19\|\kappa\|_\infty \min\left(\frac{\log(\frac{32n_X}{\delta})}{n_X}, \frac{\log(\frac{32n_Y}{\delta})}{n_Y}\right)$. Let

$$\mathcal{E}_{\text{MMD}}(n_X, n_Y, \delta) := \mathcal{E}_{\text{KME}}(n_X, \delta/2) + \mathcal{E}_{\text{KME}}(n_Y, \delta/2),$$

where

$$\mathcal{E}_{\text{KME}}(n, \delta) := \frac{\|\kappa\|_\infty^{1/2}}{\sqrt{n}} \left(2\sqrt{2\log(4/\delta)} + \sqrt{57\log(32n/\delta)} \right).$$

Under Assumptions 2.1 and 4.4, it holds with probability at least $1 - \delta$:

$$|\text{MMD}(P, Q) - \hat{\Psi}(X, Y)| \leq \mathcal{E}_{\text{MMD}}(n_X, n_Y, \delta),$$

provided that n_X, n_Y are large enough, namely:

- $\min(n_X, n_Y) \geq (1655 + 233 \log(8\|\kappa\|_\infty/\delta))\|\kappa\|_\infty$
- $\frac{\log(\frac{64n_X}{\delta})}{n_X} \leq C(\|\kappa\|_\infty, \|C_P\|, \gamma, z, c_\gamma)$ and $\frac{\log(\frac{64n_Y}{\delta})}{n_Y} \leq C(\|\kappa\|_\infty, \|C_Q\|, \gamma, z, c_\gamma)$.

for some constant C made explicit in the proof. (\rightarrow Proof)

4.2.2 Bound on the empirical quantile threshold

We now derive an upper-bound on the quantile of the permuted test statistic conditionally on X, Y . It will serve as an upper-bound on the random empirical quantile threshold $\hat{\Psi}_{(b_\alpha)}$ in our main result. For this, we show (cf. Appendix B.3) that the randomly permuted squared test statistic can be decomposed as the sum of a (weighted) U -statistic and a remainder term:

$$\hat{\Psi}^\sigma(X, Y)^2 = \frac{(n_X-1)(n_Y-1)}{n_X n_Y} \hat{U}^\sigma(X, Y) + \hat{R}^\sigma(X, Y).$$

We bound the quantile of $\hat{\Psi}^\sigma(X, Y)$ given (X, Y) , leveraging a result from Kim et al. (2022) to control the quantile of the U -statistic.

Lemma 4.6 (Quantile bound for $\hat{\Psi}^\sigma(X, Y)|X, Y$): Let $0 < \alpha \leq e^{-1}$. Under Assumption 2.1, there is a universal constant C' such that the test statistic associated to the Nyström kernel approximation satisfies, with probability at least $1 - \alpha$ conditioned on X and Y ,

$$\hat{\Psi}^\sigma(X, Y)^2 \leq C' \frac{\sqrt{n(n-1)} \|\kappa\|_\infty}{n_X^2} \log(1/\alpha) + \left(\frac{1}{n_X} + \frac{1}{n_Y}\right) 4 \|\kappa\|_\infty$$

(→ Proof)

We use this bound to establish a deterministic upper bound on the (random) empirical quantile threshold $\hat{\Psi}_{(b_\alpha)}$ using a result from Domingo-Enrich et al. (2023).

Lemma 4.7 (Bound on the empirical quantile): Let $\beta > 0$, and $1/(2e) \geq \alpha \geq 1/(\mathcal{P} + 1)$. Let $c_\alpha = \lfloor \alpha(\mathcal{P} + 1) \rfloor$. Then, with probability $1 - \beta/2$ conditioned on X and Y ,

$$\hat{\Psi}_{(b_\alpha)} \leq \|\kappa\|_\infty^{1/2} \left(\sqrt{C'} \frac{\sqrt{n}}{n_X} \sqrt{\log\left(\frac{2e}{\alpha(\beta/2)^{1/c_\alpha}}\right)} + 2\left(\frac{1}{\sqrt{n_X}} + \frac{1}{\sqrt{n_Y}}\right) \right)$$

(→ Proof)

Using Lemma 4.7, we are able to complete the proof of our main result, which is deferred to Appendix B.4.

5 Numerical studies

In this section, the empirical power and computational trade-offs of our method are explored and compared against the random Fourier features (RFF) approach presented in Choi et al. (2024b). We consider both AKRLS computed using the method of Musco et al. (2017) and uniform sampling for the selection of the Nyström landmarks. In all our experiments, we consider a Gaussian kernel $\kappa(x, y) = \exp(-\|x - y\|^2 / (2\sigma^2))$, whose bandwidth is fixed as the median of the inter-points Euclidean distance, estimated on a subset of 2000 random instances. The level of the test is set at $\alpha = 0.05$. The size is determined with $\mathcal{P} = 200$ permutations and the results are averaged over 400 repetitions. The error is estimated using Wilson score intervals at the 95% confidence level (Newcombe 1998). All our code is open-source and written in Python².

²<https://anonymous.4open.science/r/nystrom-mmd-0f7b/>

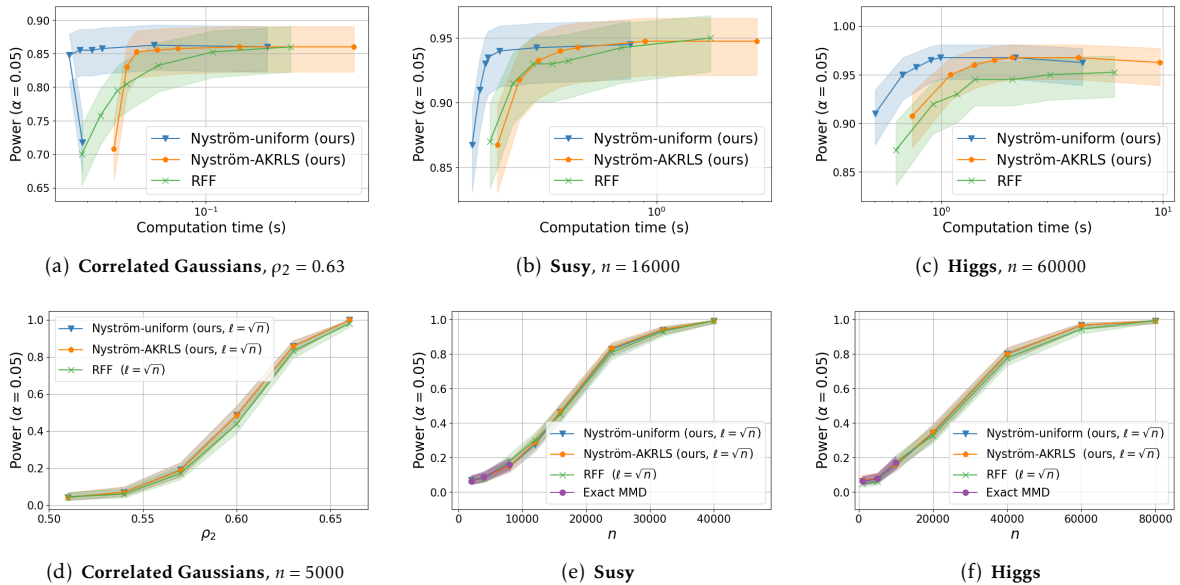


Figure 1: Power against computation time (top row) and power against correlation coefficient ρ_2 and total number of samples (bottom row). Each column represents a different dataset.

5.1 Datasets

Correlated Gaussians We consider a synthetic dataset based on 3-dimensional correlated Gaussian distributions. The first sample is drawn from $P = \mathcal{N}_3(\mathbf{0}, \Sigma(\rho_1))$, with covariance matrix $\Sigma(\rho_1) = \text{diag}(1, 1, 1) + \rho_1(1 - \text{diag}(1, 1, 1))$, and $\rho_1 = 0.5$. The second sample follows the same distribution, but with a varying correlation coefficient $0.51 \leq \rho_2 \leq 0.66$. The sample size is fixed at $n_x = n_y = 2500$.

Higgs and Susy These datasets consist of Monte Carlo simulations of collider data from high-energy physics (HEP) (Baldi et al. 2014). They are characterised by two classes, background data (P_b) and signal data (P_s). The latter is composed of processes producing a Higgs boson and supersymmetric particles, for the Higgs dataset and Susy dataset respectively. The Higgs dataset has 21 features. The last 7 features are functions of the first 21 and can be considered more discriminative *high-level* features. Consequently we will refer to the first 21 as *low-level* features. The dataset is composed of a total of 11M examples. We will only consider low-level features ($d = 14$). To mimic a typical scenario in HEP and make the test more challenging, we compare data from the distribution $P = P_b$ against a mixture of background and signal processes $Q = (1 - \alpha_{\text{mix}})P_b + \alpha_{\text{mix}}P_s$, with $\alpha_{\text{mix}} = 0.2$. The Susy dataset has 8 low-level features and 10 high-level features (that we do not consider). Again, we test background processes against a mixture of background and signal processes, with $\alpha_{\text{mix}} = 0.05$.

5.2 Results

In Figure 1 (first row), we report the power against computation time of the selected methods on the 3 datasets described in Section 5.1. We observe that our approach always matches, and in multiple cases outperforms, the results obtained using random Fourier features (RFF). The latter are themselves known to empirically outperform other linear-time test statistics existing in the literature (Choi et al. 2024a). In the first row, we fix the compression level at $\ell = \sqrt{n}$ and compare the same methods while varying the parameter ρ_2 for the synthetic dataset and the number of samples n for Susy and Higgs. We also report the result of the exact MMD estimators for the smaller values of n at which computations remain reasonable: in these settings, one can clearly observe that the efficiency of our

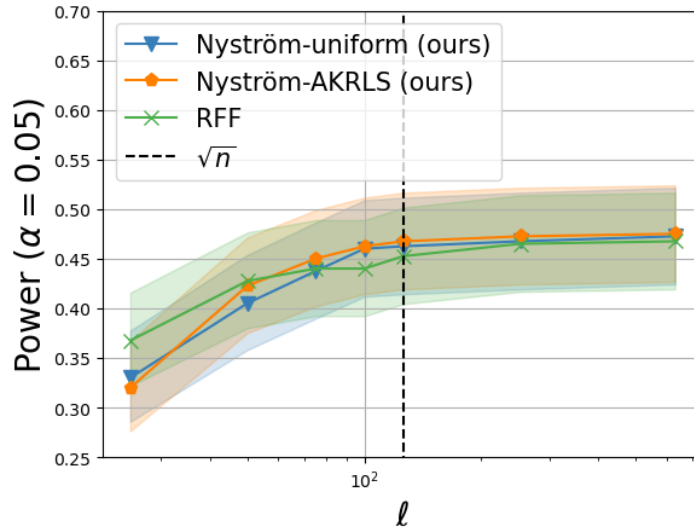


Figure 2: **Susy**, power against number of features, $n = 16000$.

approximate test procedure is achieved without compromising the power of the test.

In Figure 2, we report for the Susy dataset the power as a function of the compression level, which for our approach corresponds to the number of Nyström landmarks, and for RFF corresponds to the number of (real) features used in the approximation. Our approach performs similarly to random Fourier features, as well as Nyström with uniform sampling.

6 Conclusion

We introduced a computationally efficient procedure for two-sample testing, leveraging a data-adaptive Nyström approximation of the MMD as the test statistic. We provided a bound on the (non-asymptotic) power of the resulting test and demonstrated that it matches the minimax MMD-separation rates. Our procedure is simple to implement and compares favorably with existing state-of-the-art approaches, both theoretically and empirically. We leave open the possibility of using the proposed test statistic to design aggregated tests or to handle weighted samples.

Acknowledgments: All the authors acknowledge the financial support of the European Research Council (grant SLING 819789). L. R. acknowledges the financial support of the European Commission (Horizon Europe grant ELIAS 101120237), the US Air Force Office of Scientific Research (FA8655-22-1-7034), and the Ministry of Education, University and Research (FARE grant ML4IP R205T7J2KP; grant BAC FAIR PE00000013 funded by the EU - NGEU). This work represents only the view of the authors. The European Commission and the other organizations are not responsible for any use that may be made of the information it contains.

References

- Alaoui, Ahmed and Michael W. Mahoney (2015). “Fast Randomized Kernel Ridge Regression with Statistical Guarantees”. In: *Advances in Neural Information Processing Systems*, pp. 775–783.
- Amram, Oz et al. (2024). “CaloChallenge 2022: A Community Challenge for Fast Calorimeter Simulation”. In: ed. by Claudius Krause, Michele Fucci Giannelli, Gregor Kasieczka, Benjamin Nachman, Dalila Salamani, David Shih, and Anna Zaborowska. arXiv: [2410.21611](https://arxiv.org/abs/2410.21611) [[physics.ins-det](#)].
- Balasubramanian, Krishnakumar, Tong Li, and Ming Yuan (2021). “On the Optimality of Kernel-Embedding Based Goodness-of-Fit Tests”. In: *Journal of machine learning research* 22.1, pp. 1–45.
- Baldi, Pierre, Peter Sadowski, and Daniel Whiteson (2014). “Searching for Exotic Particles in High-Energy Physics with Deep Learning”. In: *Nature Commun.* 5, p. 4308. arXiv: [1402.4735](https://arxiv.org/abs/1402.4735) [[hep-ph](#)].
- Baraud, Yannick (2002). “Non-asymptotic minimax rates of testing in signal detection”. In: *Bernoulli* 8.5, pp. 577–606.
- Berlinet, Alain and Christine Thomas-Agnan (2004). *Reproducing Kernel Hilbert Spaces in Probability and Statistics*. Springer. ISBN: 978-1-4419-9096-9.
- Biggs, Felix, Antonin Schrab, and Arthur Gretton (2023). *MMD-FUSE: Learning and Combining Kernels for Two-Sample Testing Without Data Splitting*. Version 2. URL: <https://arxiv.org/abs/2306.08777>. Pre-published.
- Chakravarti, Purvasha, Mikael Kuusela, Jing Lei, and Larry Wasserman (2023). “Model-independent detection of new physics signals using interpretable SemiSupervised classifier tests”. In: *The Annals of Applied Statistics* 17.4, pp. 2759–2795.
- Chatalic, Antoine, Nicolas Schreuder, Ernesto De Vito, and Lorenzo Rosasco (2023). *Efficient Numerical Integration in Reproducing Kernel Hilbert Spaces via Leverage Scores Sampling*. arXiv: [2311.13548](https://arxiv.org/abs/2311.13548) [[cs](#), [math](#), [stat](#)]. URL: <http://arxiv.org/abs/2311.13548>. Pre-published.
- Chatalic, Antoine, Nicolas Schreuder, Lorenzo Rosasco, and Alessandro Rudi (2022). “Nyström Kernel Mean Embeddings”. In: *Proceedings of the 39th International Conference on Machine Learning*. Vol. 162. Proceedings of Machine Learning Research. PMLR, pp. 3006–3024.
- Chen, Yifan and Yun Yang (2021). “Fast Statistical Leverage Score Approximation in Kernel Ridge Regression”. In: *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*. International Conference on Artificial Intelligence and Statistics. PMLR, pp. 2935–2943.
- Choi, Ikjun and Ilmun Kim (2024a). *Computational-Statistical Trade-off in Kernel Two-Sample Testing with Random Fourier Features*. arXiv: [2407.08976](https://arxiv.org/abs/2407.08976) [[cs](#), [math](#), [stat](#)]. URL: <http://arxiv.org/abs/2407.08976>. Pre-published.
- (2024b). “Computational-Statistical Trade-off in Kernel Two-Sample Testing with Random Fourier Features”. In: *arXiv preprint arXiv:2407.08976*.
- Chwialkowski, Kacper P., Aaditya Ramdas, Dino Sejdinovic, and Arthur Gretton (2015). “Fast Two-Sample Testing with Analytic Representations of Probability Measures”. In: *Advances in Neural Information Processing Systems* 28.
- D’Agnolo, Raffaele Tito, Gaia Grosso, Maurizio Pierini, Andrea Wulzer, and Marco Zanetti (2021). “Learning multivariate new physics”. In: *Eur. Phys. J. C* 81.1, p. 89. arXiv: [1912.12155](https://arxiv.org/abs/1912.12155) [[hep-ph](#)].
- Diestel, Joseph and John Jerry Uhl (1977). *Vector Measures*. Mathematical Surveys and Monographs 15. American Mathematical Soc. ISBN: 0-8218-1515-6 978-0-8218-1515-1.
- Domingo-Enrich, Carles, Raaz Dwivedi, and Lester Mackey (2023). *Compress Then Test: Powerful Kernel Testing in Near-linear Time*. arXiv: [2301.05974](https://arxiv.org/abs/2301.05974) [[cs](#), [math](#), [stat](#)]. URL: <http://arxiv.org/abs/2301.05974>. Pre-published.
- Fanuel, Michaël, Joachim Schreurs, and Johan A. K. Suykens (2022). “Nyström Landmark Sampling and Regularized Christoffel Functions”. In: *Machine Learning* 111.6, pp. 2213–2254. ISSN: 1573-0565.
- Fromont, Magalie, Béatrice Laurent, Matthieu Lerasle, and Patricia Reynaud-Bouret (2012). “Kernels Based Tests with Non-Asymptotic Bootstrap Approaches for Two-Sample Problems”. In: *Conference on Learning Theory*. JMLR Workshop and Conference Proceedings, pp. 23–1.

- Fukumizu, Kenji, Arthur Gretton, Xiaohai Sun, and Bernhard Schölkopf (2007). “Kernel measures of conditional dependence”. In: *Advances in neural information processing systems* 20.
- Gretton, Arthur, Karsten Borgwardt, Malte Rasch, Bernhard Schölkopf, and Alex J. Smola (2007). “A Kernel Method for the Two-Sample-Problem”. In: *Advances in Neural Information Processing Systems*, pp. 513–520.
- Gretton, Arthur, Karsten M. Borgwardt, Malte J. Rasch, Bernhard Schölkopf, and Alexander Smola (2012). “A Kernel Two-Sample Test”. In: *The Journal of Machine Learning Research* 13.1, pp. 723–773.
- Grossi, Samuele, Marco Letizia, and Riccardo Torre (2024). “Refereeing the Referees: Evaluating Two-Sample Tests for Validating Generators in Precision Sciences”. In: arXiv: [2409.16336](https://arxiv.org/abs/2409.16336) [stat, ML].
- Hagrass, Omar, Bharath K. Sriperumbudur, and Bing Li (2023). *Spectral Regularized Kernel Two-Sample Tests*. arXiv: [2212.09201](https://arxiv.org/abs/2212.09201) [cs, math, stat]. URL: <http://arxiv.org/abs/2212.09201>. Pre-published.
- Harchaoui, Zaid, Francis Bach, and Eric Moulines (2008). *Testing for Homogeneity with Kernel Fisher Discriminant Analysis*. arXiv: [0804.1026](https://arxiv.org/abs/0804.1026) [stat]. URL: <http://arxiv.org/abs/0804.1026>. Pre-published.
- Hemerik, Jesse and Jelle Goeman (2018). “Exact Testing with Random Permutations”. In: *TEST* 27.4, pp. 811–825. ISSN: 1133-0686, 1863-8260.
- Hoeffding, Wassily (1952). “The large-sample power of tests based on permutations of observations”. In: *The Collected Works of Wassily Hoeffding*. Springer, pp. 247–271.
- Huggins, Jonathan and Lester Mackey (2018). “Random Feature Stein Discrepancies”. In: *Advances in Neural Information Processing Systems*. Vol. 31. Curran Associates, Inc.
- Jitkrittum, Wittawat, Zoltán Szabó, Kacper P Chwiałkowski, and Arthur Gretton (2016). “Interpretable Distribution Features with Maximum Testing Power”. In: *Advances in Neural Information Processing Systems*. Vol. 29. Curran Associates, Inc.
- Kalinke, Florian, Zoltan Szabo, and Bharath K. Sriperumbudur (2024). *Nyström Kernel Stein Discrepancy*. arXiv: [2406.08401](https://arxiv.org/abs/2406.08401) [cs, math, stat]. URL: <http://arxiv.org/abs/2406.08401>. Pre-published.
- Kanamori, Takafumi, Taiji Suzuki, and Masashi Sugiyama (2011). “F-Divergence Estimation and Two-Sample Homogeneity Test Under Semiparametric Density-Ratio Models”. In: *IEEE Transactions on Information Theory* 58.2, pp. 708–720.
- Kim, Ilmun, Sivaraman Balakrishnan, and Larry Wasserman (2022). *Minimax Optimality of Permutation Tests*. arXiv: [2003.13208](https://arxiv.org/abs/2003.13208) [math, stat]. URL: <http://arxiv.org/abs/2003.13208>. Pre-published.
- Kim, Ilmun and Antonin Schrab (2024). *Differentially Private Permutation Tests: Applications to Kernel Methods*. arXiv: [2310.19043](https://arxiv.org/abs/2310.19043) [cs, math, stat]. URL: <http://arxiv.org/abs/2310.19043>. Pre-published.
- Kübler, Jonas M., Wittawat Jitkrittum, Bernhard Schölkopf, and Krikamol Muandet (2022a). “A Witness Two-Sample Test”. arXiv: [2102.05573](https://arxiv.org/abs/2102.05573) [cs, stat].
- Kübler, Jonas M., Vincent Stimper, Simon Buchholz, Krikamol Muandet, and Bernhard Schölkopf (2022b). *AutoML Two-Sample Test*. URL: <https://arxiv.org/abs/2206.08843v3>. Pre-published.
- Lehmann, E.L. and Joseph P. Romano (2022). *Testing Statistical Hypotheses*. en. Springer Texts in Statistics. Cham: Springer International Publishing. ISBN: 978-3-030-70577-0 978-3-030-70578-7.
- Letizia, Marco, Gianvito Losapio, Marco Rando, Gaia Grosso, Andrea Wulzer, Maurizio Pierini, Marco Zanetti, and Lorenzo Rosasco (2022). “Learning new physics efficiently with nonparametric methods”. In: *Eur. Phys. J. C* 82.10, p. 879. arXiv: [2204.02317](https://arxiv.org/abs/2204.02317) [hep-ph].
- Li, Tong and Ming Yuan (2019). *On the Optimality of Gaussian Kernel Based Nonparametric Tests against Smooth Alternatives*. arXiv: [1909.03302](https://arxiv.org/abs/1909.03302) [math, stat]. URL: <http://arxiv.org/abs/1909.03302>. Pre-published.
- (2024). “On the Optimality of Gaussian Kernel Based Nonparametric Tests against Smooth Alternatives”. In: *Journal of Machine Learning Research* 25.334, pp. 1–62.

- Müller, Alfred (1997). “Integral probability metrics and their generating classes of functions”. In: *Advances in applied probability* 29.2, pp. 429–443.
- Musco, Cameron and Christopher Musco (2017). “Recursive Sampling for the Nystrom Method”. In: *Advances in Neural Information Processing Systems*, pp. 3833–3845.
- Newcombe, Robert G (1998). “Two-sided confidence intervals for the single proportion: comparison of seven methods”. In: *Statistics in medicine* 17.8, pp. 857–872.
- Nyström, E. J. (1930). “Über Die Praktische Auflösung von Integralgleichungen mit Anwendungen auf Randwertaufgaben”. In: *Acta Mathematica* 54, pp. 185–204. ISSN: 0001-5962, 1871-2509.
- Pauwels, Edouard, Francis Bach, and Jean-Philippe Vert (2018). “Relating Leverage Scores and Density Using Regularized Christoffel Functions”. In: *Proceedings of the 32nd International Conference on Neural Information Processing Systems* (Montréal, Canada). NIPS’18. Red Hook, NY, USA: Curran Associates Inc., pp. 1670–1679.
- Pinelis, Iosif (1994). “Optimum Bounds for the Distributions of Martingales in Banach Spaces”. In: *The Annals of Probability*, pp. 1679–1706.
- Ramdas, Aaditya, Sashank J. Reddi, Barnabás Póczos, Aarti Singh, and Larry A. Wasserman (2015). “Adaptivity and Computation-Statistics Tradeoffs for Kernel and Distance based High Dimensional Two Sample Testing”. In: *CoRR abs/1508.00655*. arXiv: [1508.00655](https://arxiv.org/abs/1508.00655).
- Reddi, Sashank, Aaditya Ramdas, Barnabás Póczos, Aarti Singh, and Larry Wasserman (2015). “On the high dimensional power of a linear-time two sample test under mean-shift alternatives”. In: pp. 772–780.
- Rudi, Alessandro, Daniele Calandriello, Luigi Carratino, and Lorenzo Rosasco (2018). “On Fast Leverage Score Sampling and Optimal Learning”. In: *Advances in Neural Information Processing Systems*, pp. 5672–5682.
- Rudi, Alessandro, Raffaello Camoriano, and Lorenzo Rosasco (2015). “Less Is More: Nyström Computational Regularization”. In: *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 1*. NIPS’15. Cambridge, MA, USA: MIT Press, pp. 1657–1665.
- Schrab, Antonin, Ilmun Kim, Mélisande Albert, Béatrice Laurent, Benjamin Guedj, and Arthur Gretton (2023). “MMD Aggregated Two-Sample Test”. In: *Journal of Machine Learning Research* 24.194, pp. 1–81. ISSN: 1533-7928.
- Schrab, Antonin, Ilmun Kim, Benjamin Guedj, and Arthur Gretton (2022). “Efficient Aggregated Kernel Tests Using Incomplete U-statistics”. In: *Advances in Neural Information Processing Systems* 35, pp. 18793–18807.
- Shekhar, Shubhanshu, Ilmun Kim, and Aaditya Ramdas (2022). “A Permutation-Free Kernel Two-Sample Test”. In: *Advances in Neural Information Processing Systems* 35, pp. 18168–18180.
- Simon-Gabriel, Carl-Johann and Bernhard Schölkopf (2018). “Kernel distribution embeddings: Universal kernels, characteristic kernels and kernel metrics on distributions”. In: *Journal of Machine Learning Research* 19.44, pp. 1–29.
- Sriperumbudur, Bharath, Kenji Fukumizu, and Gert Lanckriet (2010). “On the Relation between Universality, Characteristic Kernels and RKHS Embedding of Measures”. In: *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*. JMLR Workshop and Conference Proceedings, pp. 773–780.
- Widom, Harold (1964). “Asymptotic Behavior of the Eigenvalues of Certain Integral Equations. II”. In: *Archive for Rational Mechanics and Analysis* 17.3, pp. 215–229.
- Williams, Christopher and Matthias Seeger (2001). “Using the Nyström Method to Speed up Kernel Machines”. In: *Advances in Neural Information Processing Systems*, pp. 682–688.
- Yamada, Makoto, Denny Wu, Yao-Hung Hubert Tsai, Hirofumi Ohta, Ruslan Salakhutdinov, Ichiro Takeuchi, and Kenji Fukumizu (2018). “Post Selection Inference with Incomplete Maximum Mean Discrepancy Estimator”. In: *International Conference on Learning Representations*.
- Zaremba, Wojciech, Arthur Gretton, and Matthew Blaschko (2013). “B-Test: A Non-parametric, Low Variance Kernel Two-sample Test”. In: *Advances in Neural Information Processing Systems*. Vol. 26. Curran Associates, Inc.

- Zhang, Qinyi, Veit Wild, Sarah Filippi, Seth Flaxman, and Dino Sejdinovic (2022). “Bayesian Kernel Two-Sample Testing”. In: *Journal of Computational and Graphical Statistics* 31.4, pp. 1164–1176. ISSN: 1061-8600.
- Zhao, Ji and Deyu Meng (2015). “FastMMD: Ensemble of Circular Discrepancy for Efficient Two-Sample Test”. In: *Neural Computation* 27.6, pp. 1345–1372. ISSN: 0899-7667.
- Zhao, Shengjia, Abhishek Sinha, Yutong He, Aidan Perreault, Jiaming Song, and Stefano Ermon (2021). “Comparing Distributions by Measuring Differences That Affect Decision Making”. In: *International Conference on Learning Representations*.

A Quality of the MMD approximation: additional lemmas and proofs

For any projector P , we denote $P^\perp = I - P$. For an event $B \in \mathcal{B}$ we denote by B^c the complement of the event B , $B^c := E \setminus B$.

We first provide a lemma quantifying the impact of the Nyström projections on the covariance operators in operator norm, and then prove the bound on the MMD approximation (Lemma 4.5).

Similarly to Φ_Z and P_Z , we define $\Phi_{\tilde{X}} = [\phi(\tilde{x}_1), \dots, \phi(\tilde{x}_{\ell_X})] : \mathbb{R}^{\ell_X} \rightarrow \mathcal{H}$, $\Phi_{\tilde{Y}} = [\phi(\tilde{y}_1), \dots, \phi(\tilde{y}_{\ell_Y})] : \mathbb{R}^{\ell_Y} \rightarrow \mathcal{H}$, and $P_{\tilde{X}}, P_{\tilde{Y}}$ to be the orthogonal projections respectively onto $\text{span}(\Phi_{\tilde{X}})$ and $\text{span}(\Phi_{\tilde{Y}})$. Note that for any $v \in \mathcal{H}$, it holds $\|P_{\tilde{X}}^\perp v\| \leq \|P_{\tilde{X}}^\perp v\|$ and $\|P_{\tilde{Y}}^\perp v\| \leq \|P_{\tilde{Y}}^\perp v\|$.

Lemma A.1: Let $\delta > 0$. Let $\ell = \ell_X + \ell_Y$ where $\ell_X = L(n_X)$, $\ell_Y = L(n_Y)$ and L is as defined in Equation (11). Let $\tilde{x}_1, \dots, \tilde{x}_{\ell_X}$ and $\tilde{y}_1, \dots, \tilde{y}_{\ell_Y}$ be drawn with replacement proportionally to $(z, \lambda_0, \delta/4)$ -approximate leverage scores respectively from the datasets X and Y for some $z \geq 1$ and $\lambda_0 > 0$ satisfying $\lambda_0 \leq 19\|\kappa\|_\infty \min\left(\frac{\log(\frac{32n_X}{\delta})}{n_X}, \frac{\log(\frac{32n_Y}{\delta})}{n_Y}\right)$. Under Assumption 4.4 it holds jointly with probability at least $1 - \delta$:

$$\|P_{\tilde{X}}^\perp C_P^{1/2}\| \leq \mathcal{E}_{\text{Cov}}(n_X) \quad \text{and} \quad \|P_{\tilde{Y}}^\perp C_Q^{1/2}\| \leq \mathcal{E}_{\text{Cov}}(n_Y) \quad \text{where} \quad \mathcal{E}_{\text{Cov}}(n) := \sqrt{57\|\kappa\|_\infty \frac{\log(\frac{32n}{\delta})}{n}} \quad (12)$$

provided n_X, n_Y are large enough, namely:

- $\min(n_X, n_Y) \geq (1655 + 233 \log(8\|\kappa\|_\infty/\delta))\|\kappa\|_\infty$
- $\frac{\log(\frac{32n_X}{\delta})}{n_X} \leq C(\|\kappa\|_\infty, \|C_P\|, \gamma, z, c_\gamma)$ and $\frac{\log(\frac{32n_Y}{\delta})}{n_Y} \leq C(\|\kappa\|_\infty, \|C_Q\|, \gamma, z, c_\gamma)$.

for some function C made explicit in the proof.

(\rightarrow Proof)

Proof of Lemma A.1: We apply Rudi et al. (2015, Lemma 7) for X with probability $\delta/4$. We get

$$\mathbb{P}\left[\|P_{\tilde{X}}^\perp C_P^{1/2}\| \leq \sqrt{3\lambda_X}\right] \geq 1 - \delta/2 \quad (13)$$

and pick $\lambda_X = \frac{19\|\kappa\|_\infty \log(32n_X/\delta)}{n_X}$ and $\ell_X = (n_X)^\gamma \frac{78c_\gamma z^2 (\log \frac{32n_X}{\delta})^{1-\gamma}}{(19\|\kappa\|_\infty)^\gamma}$. These parameters correspond to the choice made in (Chatalic et al. 2023, Th.4.6) and thus satisfy the requirements of (Chatalic et al. 2023, Lemma F.4). A similar argument for Y yields a similar bound for $\|P_{\tilde{Y}}^\perp C_Q^{1/2}\|$, and the claimed result follows via a union bound. The function C in the statement of the lemma is thus defined as

$$C(\|\kappa\|_\infty, \|C_P\|, \gamma, z, c_\gamma) := \frac{1}{19\|\kappa\|_\infty} \min\left(\|C_P\|, \left(\frac{c_\gamma z^2}{5}\right)^{1/\gamma}\right) \quad (14)$$

We now have everything to prove Lemma 4.5.

Proof of Lemma 4.5: Applying multiple times the standard and inverse triangle inequalities,

$$\begin{aligned}
|\text{MMD}(P, Q) - \hat{\Psi}(X, Y)| &= \left| \|\mu(P) - \mu(Q)\| - \|P_Z \mu(X) - P_Z \mu(Y)\| \right| \\
&\leq \|\mu(P) - \mu(Q) - P_Z \mu(X) + P_Z \mu(Y)\| \\
&\leq \|\mu(P) - P_Z \mu(X)\| + \|\mu(Q) - P_Z \mu(Y)\| \\
&\leq \|\mu(P) - \mu(X)\| + \|P_Z^\perp \mu(P)\| + \|\mu(Q) - \mu(Y)\| + \|P_Z^\perp \mu(Q)\| \\
&\leq \|\mu(P) - \mu(X)\| + \|P_X^\perp \mu(X)\| + \|\mu(Q) - \mu(Y)\| + \|P_Y^\perp \mu(Y)\| .
\end{aligned}$$

We apply Chatalic et al. (2023, Theorem G.1) which is based on the standard Pinelis' concentration inequality (see Pinelis 1994, Th. 3.5) both for X and Y with probability $\delta/4$, and Lemma A.1 with probability $\delta/2$. Using a union bound, we get jointly with probability at least $1 - \delta$:

$$|\text{MMD}(P, Q) - \hat{\Psi}(X, Y)| \leq \mathcal{E}_{\text{KME}}(n_X, \delta/2) + \mathcal{E}_{\text{KME}}(n_Y, \delta/2). \quad (15)$$

The function C is the one defined in the proof of Lemma A.1.

B Controlling the level and the power

B.1 Proof of Lemma 4.1

Proof of Lemma 4.1: Note that under the null hypothesis, the distribution of the data $(X_1, \dots, X_{n_X}, Y_1, \dots, Y_{n_Y})$ is invariant to permutation of its coordinates. The stated lemma is a consequence of Hemerik et al. (2018, Proposition 3).

B.2 Proof of Lemma 4.3

Proof of Lemma 4.3: Define the events $\mathcal{A} = \{\hat{\Psi}(X, Y) \leq \hat{\Psi}_{(b_\alpha)}\}$ and $\mathcal{B}_\beta = \{\text{MMD}(P, Q) \geq \mathcal{E}_{\text{MMD}}(n_X, n_Y, \beta/2) + \hat{\Psi}_{(b_\alpha)}\}$. Assume that $\mathbb{P}[\mathcal{B}_\beta] > 1 - \beta/2$. We will show that $\mathbb{P}[\mathcal{A}] \leq \beta$. By the law of total probability, it holds

$$\mathbb{P}[\mathcal{A}] = \mathbb{P}[\mathcal{A} \cap \mathcal{B}_\beta] + \mathbb{P}[\mathcal{A} | \mathcal{B}_\beta^c] \mathbb{P}[\mathcal{B}_\beta^c] \leq \mathbb{P}[\mathcal{A} \cap \mathcal{B}_\beta] + \mathbb{P}[\mathcal{B}_\beta^c].$$

Using Lemma 4.5, the first term in the above upper bound can be bounded as

$$\mathbb{P}[\mathcal{A} \cap \mathcal{B}_\beta] = \mathbb{P}[\hat{\Psi}(X, Y) \leq \hat{\Psi}_{(b_\alpha)}, \hat{\Psi}_{(b_\alpha)} \leq \text{MMD}(P, Q) - \mathcal{E}_{\text{MMD}}(n_X, n_Y, \beta/2)] \quad (16)$$

$$\leq \mathbb{P}[\hat{\Psi}(X, Y) \leq \text{MMD}(P, Q) - \mathcal{E}_{\text{MMD}}(n_X, n_Y, \beta/2)] \quad (17)$$

$$\leq \mathbb{P}[|\text{MMD}(P, Q) - \hat{\Psi}(X, Y)| \geq \mathcal{E}_{\text{MMD}}(n_X, n_Y, \beta/2)] \quad (18)$$

$$\leq \beta/2 . \quad (19)$$

We get the desired result by combining the last inequality and the fact that $\mathbb{P}[\mathcal{B}_\beta^c] \leq \beta/2$.

B.3 Bounds on the quantiles

In order to prove Lemma 4.6 and Lemma 4.7 below, we first show how the squared test statistic can be decomposed as the sum of a U -statistic and a remainder, and provide a bound on the quantile of this U -statistic.

The (randomly) permuted squared test statistic can be written

$$\hat{\Psi}^\sigma(X, Y)^2 := \left\| \frac{1}{n_X} \sum_{i=1}^{n_X} \varphi(W_{\sigma(i)}) - \frac{1}{n_Y} \sum_{j=1}^{n_Y} \varphi(W_{\sigma(n_X+j)}) \right\|^2 \quad (20)$$

$$:= \left\| \frac{1}{n_Y n_X} \sum_{i=1}^{n_X} \sum_{j=1}^{n_Y} (\varphi(W_{\sigma(i)}) - \varphi(W_{\sigma(n_X+j)})) \right\|^2 \quad (21)$$

$$= \frac{1}{n_X^2 n_Y^2} \sum_{i=1}^{n_X} \sum_{j=1}^{n_Y} \sum_{i'=1}^{n_X} \sum_{j'=1}^{n_Y} h^\sigma(i, i'; j, j') \quad (22)$$

with $h^\sigma(i, i'; j, j') := \tilde{\kappa}(W_{\sigma(i)}, W_{\sigma(i')}) - \tilde{\kappa}(W_{\sigma(i)}, W_{\sigma(n_Y+j')}) - \tilde{\kappa}(W_{\sigma(n_Y+j)}, W_{\sigma(i')}) + \tilde{\kappa}(W_{\sigma(n_Y+j)}, W_{\sigma(n_Y+j')})$.

It can be decomposed as the sum of a (weighted) U -statistic and a remainder term,

$$\hat{\Psi}^\sigma(X, Y)^2 = \frac{(n_X-1)(n_Y-1)}{n_X n_Y} \hat{U}^\sigma(X, Y) + \hat{R}^\sigma(X, Y), \quad (23)$$

with

$$\hat{U}^\sigma(X, Y) := \frac{1}{n_X(n_X-1)n_Y(n_Y-1)} \sum_{1 \leq i \neq i' \leq n_X} \sum_{1 \leq j \neq j' \leq n_Y} h^\sigma(i, i'; j, j'), \quad (24)$$

and

$$\hat{R}^\sigma(X, Y) = \frac{1}{n_X^2 n_Y^2} \left(\sum_{i=1}^{n_X} \sum_{1 \leq j \neq j' \leq n_Y} h^\sigma(i, i; j, j') + \sum_{1 \leq i \neq i' \leq n_X} \sum_{j=1}^{n_Y} h^\sigma(i, i'; j, j) + \sum_{i=1}^{n_X} \sum_{j=1}^{n_Y} h^\sigma(i, i; j, j) \right). \quad (25)$$

We control the U -statistic with high probability using a result from Kim et al. (2022) while we obtain an upper-bound for $\hat{R}^\sigma(X, Y)$ using the boundedness assumption.

Lemma B.1 (Bound on the quantile of $\hat{U}^\sigma(X, Y)$ given X, Y): Let $0 < \alpha \leq e^{-1}$. Then there is a universal constant C' such that the U -statistic (24) associated to the Nyström kernel approximation satisfies

$$\mathbb{P} \left[\hat{U}^\sigma(X, Y) \leq C' \frac{\sqrt{n(n-1)} \|\kappa\|_\infty}{n_X(n_X-1)} \log(1/\alpha) \mid X, Y \right] \geq 1 - \alpha, \quad (26)$$

i.e. the r.h.s. of the above equation is an upper bound for the quantile $q_{1-\alpha, \hat{U}^\sigma(X, Y)}$. (\rightarrow Proof)

Proof of Lemma B.1: Let

$$\Sigma_{n_X, n_Y}^2 := \frac{1}{n_X^2 (n_X-1)^2} \left(\sum_{1 \leq i \neq i' \leq n_X} \langle \varphi(X_i), \varphi(X_{i'}) \rangle^2 \right) \quad (27)$$

$$+ \sum_{1 \leq j \neq j' \leq n_Y} \langle \varphi(Y_j), \varphi(Y_{j'}) \rangle^2 + 2 \sum_{1 \leq i \leq n_X, 1 \leq j \leq n_Y} \langle \varphi(X_i), \varphi(Y_j) \rangle^2 \quad (28)$$

Given that $\tilde{\kappa}(x, y)^2 = \langle P_Z \phi(x), P_Z \phi(y) \rangle^2 \leq \|\phi(x)\|^2 \|\phi(y)\|^2$ as P_Z is a projector, it holds

$$\Sigma_{n_X, n_Y}^2 \leq \frac{n(n-1) \|\kappa\|_\infty^2}{n_X^2 (n_X-1)^2}. \quad (29)$$

By Kim et al. (2022, Th. 6.1 and (31)), it holds for some constant C' conditionally on X, Y with probability at least $1 - \alpha$:

$$\hat{U}^\sigma(X, Y) \leq \max\left(\sqrt{C^{-1} \Sigma_{n_X, n_Y}^2 \log(1/\alpha)}, C^{-1} \Sigma_{n_X, n_Y} \log(1/\alpha)\right) \quad (30)$$

$$\leq C' \Sigma_{n_X, n_Y} \log(1/\alpha) \quad (31)$$

where C refers to the absolute constant of Kim et al. (2022, Th. 6.1).

We can now prove Lemma 4.6.

Proof of Lemma 4.6: Starting from (23)

$$\hat{\Psi}(X, Y)^2 = \frac{(n_X-1)(n_Y-1)}{n_X n_Y} \hat{U}^\sigma(X, Y) + \hat{R}^\sigma(X, Y) \quad (32)$$

Given that $\tilde{\kappa}(x, y) = \langle P_Z \phi(x), P_Z \phi(y) \rangle \leq \|\phi(x)\| \|\phi(y)\|$ as P_Z is a projector, it holds almost surely

$$|\hat{R}^\sigma(X, Y)| \leq \frac{(n_Y+n_X-1)}{n_X n_Y} 4 \|\kappa\|_\infty \leq \left(\frac{1}{n_X} + \frac{1}{n_Y}\right) 4 \|\kappa\|_\infty \quad (33)$$

$$(34)$$

Moreover, by Lemma B.1, conditionally on X, Y , with probability at least $1 - \alpha$ it holds

$$\hat{U}^\sigma(X, Y) \leq C' \frac{\sqrt{n(n-1) \|\kappa\|_\infty}}{n_X (n_X-1)} \log(1/\alpha). \quad (35)$$

Therefore, with probability at least $1 - \alpha$,

$$\hat{\Psi}(X, Y)^2 = \frac{(n_X-1)(n_Y-1)}{n_X n_Y} \hat{U}^\sigma(X, Y) + \hat{R}^\sigma(X, Y) \quad (36)$$

$$\leq C' \frac{(n_X-1)(n_Y-1)}{n_X n_Y} \frac{\sqrt{n(n-1) \|\kappa\|_\infty}}{n_X (n_X-1)} \log(1/\alpha) + \left(\frac{1}{n_X} + \frac{1}{n_Y}\right) 4 \|\kappa\|_\infty \quad (37)$$

$$\leq C' \frac{\sqrt{n(n-1) \|\kappa\|_\infty}}{n_X^2} \log(1/\alpha) + \left(\frac{1}{n_X} + \frac{1}{n_Y}\right) 4 \|\kappa\|_\infty \quad (38)$$

We now prove our high-probability bound on the threshold.

Proof of Lemma 4.7: Using (Domingo-Enrich et al. 2023, Lemma 6), it holds

$$\mathbb{P}\left[\hat{\Psi}_{(b_\alpha)} \leq q_{1-\alpha_1}(X, Y) \mid X, Y\right] > 1 - \beta/2 \quad (39)$$

where $\alpha_1 := \left(\frac{\beta/2}{\lfloor \alpha^{(\mathcal{P}+1)} \rfloor}\right)^{1/\lfloor \alpha^{(\mathcal{P}+1)} \rfloor}$. Moreover, by Lemma 4.6, provided $\alpha_1 \leq e^{-1}$, it holds (almost

surely for X, Y)

$$q_{1-\alpha_1}(X, Y) \leq \|\kappa\|_\infty^{1/2} \sqrt{C' \frac{\sqrt{n(n-1)}}{n_X^2} \log(1/\alpha_1) + \left(\frac{4}{n_X} + \frac{4}{n_Y}\right)} \quad (40)$$

$$\leq \sqrt{C'} \|\kappa\|_\infty^{1/2} \frac{\sqrt{n}}{n_X} \sqrt{\log(1/\alpha_1) + \left(\frac{1}{\sqrt{n_X}} + \frac{1}{\sqrt{n_Y}}\right) 2} \|\kappa\|_\infty^{1/2} . \quad (41)$$

We will now upper bound $\log(1/\alpha_1)$ to make our bound more explicit. Denoting $c_\alpha = \lfloor \alpha(\mathcal{P} + 1) \rfloor$, it holds $c_\alpha \geq 1$ by assumption and using the inequalities $\binom{n}{k} \leq \left(\frac{en}{k}\right)^k$ and $\lfloor x \rfloor \geq x/2$, we obtain

$$\log\left(\left(\frac{\mathcal{P}}{c_\alpha}\right)^{1/c_\alpha}\right) \leq \log\left(\frac{e\mathcal{P}}{c_\alpha}\right) \leq \log\left(\frac{2e\mathcal{P}}{\alpha(\mathcal{P} + 1)}\right) \leq \log\left(\frac{2e}{\alpha}\right) .$$

The population quantile can now be upper bounded as

$$q_{1-\alpha_1}(X, Y) \leq \sqrt{C'} \|\kappa\|_\infty^{1/2} \frac{\sqrt{n}}{n_X} \sqrt{\log\left(\frac{2e}{\alpha(\beta/2)^{1/c_\alpha}}\right) + \left(\frac{1}{\sqrt{n_X}} + \frac{1}{\sqrt{n_Y}}\right) 2} \|\kappa\|_\infty^{1/2} . \quad (42)$$

We finish the proof by noting that the assumption $\alpha_1 \leq 1/e$ is not restrictive. Indeed

$$\alpha_1 = \frac{(\beta/2)^{1/c_\alpha}}{\binom{\mathcal{P}}{c_\alpha}^{1/c_\alpha}} \leq \frac{1}{\binom{\mathcal{P}}{c_\alpha}} \leq \frac{\alpha(\mathcal{P} + 1)}{\mathcal{P}} \leq 2\alpha \quad (43)$$

which is upper-bounded by e^{-1} under our assumptions.

B.4 Proof of Theorem 4.2 (main result)

Let $q_{1-\alpha}(X, Y)$ denote the $(1 - \alpha)$ -quantile (with respect to the random variable σ) for the test statistic $\hat{\Psi}^\sigma(X, Y)$.

Proof of Theorem 4.2: By Lemma 4.5 (applied with $\delta = \beta/2$), defining $\mathcal{E}_{\text{MMD}}(n_X, n_Y, \delta) := \mathcal{E}_{\text{KME}}(n_X, \delta/2) + \mathcal{E}_{\text{KME}}(n_Y, \delta/2)$ and $\mathcal{E}_{\text{KME}}(n, \delta)$ as in Lemma 4.5, it holds

$$\mathbb{P}[|\text{MMD}(P, Q) - \hat{\Psi}(X, Y)| \geq \mathcal{E}_{\text{MMD}}(n_X, n_Y, \beta/2)] \leq \beta/2. \quad (44)$$

Hence in order to upper bound the power of the test by $1 - \beta$, it is sufficient by Lemma 4.3 to find a condition on P and Q such that

$$\mathbb{P}[\hat{\Psi}_{(b_\alpha)} \leq \text{MMD}(P, Q) - \mathcal{E}_{\text{MMD}}(n_X, n_Y, \beta/2)] > 1 - \beta/2 . \quad (45)$$

By Lemma 4.7 we get

$$\mathbb{P}\left[\hat{\Psi}_{(b_\alpha)} \leq \sqrt{C'} \|\kappa\|_\infty^{1/2} \frac{\sqrt{n}}{n_X} \sqrt{\log\left(\frac{2e}{\alpha(\beta/2)^{1/c_\alpha}}\right) + \left(\frac{1}{\sqrt{n_X}} + \frac{1}{\sqrt{n_Y}}\right) 2} \|\kappa\|_\infty^{1/2} \mid X, Y\right] > 1 - \beta/2 \quad (46)$$

Hence we obtain the desired power bound provided

$$\sqrt{C'}\|\kappa\|_{\infty}^{1/2}\frac{\sqrt{n}}{n_X}\sqrt{\log\left(\frac{2e}{\alpha}(2/\beta)^{1/c_{\alpha}}\right)} + \left(\frac{1}{\sqrt{n_X}} + \frac{1}{\sqrt{n_Y}}\right)2\|\kappa\|_{\infty}^{1/2} \leq \text{MMD}(P, Q) \\ - (\mathcal{E}_{\text{KME}}(n_X, \beta/4) + \mathcal{E}_{\text{KME}}(n_Y, \beta/4)),$$

that is

$$\text{MMD}(P, Q) \geq \sqrt{C'}\|\kappa\|_{\infty}^{1/2}\frac{\sqrt{n}}{n_X}\sqrt{\log\left(\frac{2e}{\alpha}\left(\frac{2}{\beta}\right)^{1/\lceil\alpha(\mathcal{P}+1)\rceil}\right)} \\ + \left(\frac{1}{\sqrt{n_X}} + \frac{1}{\sqrt{n_Y}}\right)2\|\kappa\|_{\infty}^{1/2}\left(1 + \sqrt{2\log(16/\beta)}\right) \\ + \|\kappa\|_{\infty}^{1/2}\sqrt{57}\left(\frac{\sqrt{\log(128n_X/\beta)}}{\sqrt{n_X}} + \frac{\sqrt{\log(128n_Y/\beta)}}{\sqrt{n_Y}}\right)$$