



## **Nyström M-Hilbert-Schmidt independence criterion**

**LSE Research Online URL for this paper:** <http://eprints.lse.ac.uk/118251/>

Version: Accepted Version

---

### **Conference or Workshop Item:**

Kalinke, Florian and Szabo, Zoltan (2023) Nyström M-Hilbert-Schmidt independence criterion. In: Conference on Uncertainty in Artificial Intelligence, 2024-07-31 - 2024-08-04, University Center at Carnegie Mellon University, Pittsburgh, United States, USA.

---

### **Reuse**

Items deposited in LSE Research Online are protected by copyright, with all rights reserved unless indicated otherwise. They may be downloaded and/or printed for private study, or other acts as permitted by national copyright laws. The publisher or other rights holders may allow further reproduction and re-use of the full text version. This is indicated by the licence information on the LSE Research Online record for the item.

---

# Nyström $M$ -Hilbert-Schmidt Independence Criterion (Supplementary Material)

---

Florian Kalinke<sup>1</sup>

Zoltán Szabó<sup>2</sup>

<sup>1</sup>Institute for Program Structures and Data Organization, Karlsruhe Institute of Technology, Karlsruhe, Germany

<sup>2</sup>Department of Statistics, London School of Economics, London, UK

## A APPENDIX

Section A.1 contains two external theorems and lemmas that we use. Section A.2 is about our proofs.

### A.1 EXTERNAL THEOREMS AND LEMMAS

In this section two theorems and lemmas are recalled for self-completeness, Theorem A.1 is about bounding the error of Nyström mean embeddings [Chatalic et al., 2022, Theorem 4.1], Theorem A.2 is a well-known result [Serfling, 1980, Section 5.6, Theorem A] for bounding the deviation of U-statistics. Lemma A.1 is about connection between U- and V-statistics. Lemma A.2 recalls Markov's inequality.

**Theorem A.1** (Bound on mean embeddings). *Let  $\mathcal{X}$  be a locally compact second-countable topological space,  $X$  a random variable supported on  $\mathcal{X}$  with Borel probability measure  $\mathbb{P}$ , and let  $\mathcal{H}_k$  be a RKHS on  $\mathcal{X}$  with kernel  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ , and feature map  $\phi_k$ . Assume that there exists a constant  $K \in (0, \infty)$  such that  $\sup_{x \in \mathcal{X}} \sqrt{k(x, x)} \leq K$ . Let  $C_k = \mathbb{E}[\phi_k(X) \otimes \phi_k(X)]$ . Furthermore, assume that the data points  $\hat{\mathbb{P}}_n = \{x_1, \dots, x_n\}$  are drawn i.i.d. from the distribution  $\mathbb{P}$  and that  $n'$  subsamples  $\tilde{\mathbb{P}}_{n'} = \{\tilde{x}_1, \dots, \tilde{x}_{n'}\}$  are drawn uniformly with replacement from the dataset  $\hat{\mathbb{P}}_n$ . Then for any  $\delta \in (0, 1)$  it holds that*

$$\left\| \mu_k(\mathbb{P}) - \mu_k(\tilde{\mathbb{P}}_{n'}) \right\|_{\mathcal{H}_k} \leq \frac{c_1}{\sqrt{n}} + \frac{c_2}{n'} + \frac{c_3 \sqrt{\log(n'/\delta)}}{n'} \sqrt{\mathcal{N}_X \left( \frac{12K^2 \log(n'/\delta)}{n'} \right)},$$

with probability at least  $1 - \delta$  provided that

$$n' \geq \max \left( 67, 12K^2 \|C_k\|_{\text{op}}^{-1} \right) \log \left( \frac{n'}{\delta} \right),$$

where  $c_1 = 2K \sqrt{2 \log(6/\delta)}$ ,  $c_2 = 4\sqrt{3}K \log(12/\delta)$ , and  $c_3 = 12\sqrt{3 \log(12/\delta)}K$ .

Recall that a U-statistic is the average of a (symmetric) core function  $h = h(x_1, \dots, x_m)$  over the observations  $X_1, \dots, X_n \sim \mathbb{P}$  ( $n \geq m$ ) with form

$$U_n = U(X_1, \dots, X_m) = \frac{1}{\binom{n}{m}} \sum_c h(X_{i_1}, \dots, X_{i_m}), \quad (1)$$

where  $c$  is the set of the  $\binom{n}{m}$  combinations of  $m$  distinct elements  $\{i_1, \dots, i_m\}$  from  $\{1, \dots, n\}$ .  $U_n$  is an unbiased estimator of  $\theta = \theta(\mathbb{P}) = \mathbb{E}_{\mathbb{P}}[h(X_1, \dots, X_m)]$ .

**Theorem A.2** (Hoeffding's inequality for U-statistics). *Let  $h = h(x_1, \dots, x_m)$  be a core function for  $\theta = \theta(\mathbb{P}) = \mathbb{E}_{\mathbb{P}}[h(X_1, \dots, X_m)]$  with  $a \leq h(x_1, \dots, x_m) \leq b$ . Then, for any  $u > 0$  and  $n \geq m$ ,*

$$\mathbb{P}(U_n - \theta \geq u) \leq \exp \left( -\frac{2nu^2}{m(b-a)^2} \right).$$

Similar to (1) one can consider an alternative (slightly biased) estimator of  $\theta$ , which is called V-statistic:

$$V_n = V(X_1, \dots, X_m) = \frac{1}{n^m} \sum_{(i_1, \dots, i_m) \in T_m(n)} h(X_{i_1}, \dots, X_{i_m}), \quad (2)$$

where  $T_m(n)$  is the  $m$ -fold Cartesian product of the set  $[n]$ .

There is a close relation between U- and V-statistics, as it is made explicit by the following lemma [Serfling, 1980, Lemma, Section 5.7.3].

**Lemma A.1** (Connection between U- and V-statistics). *Let  $\mathbb{P}$  be a probability measure on a metric space  $\mathcal{X}$ . Let  $(X_i)_{i \in [n]} \stackrel{i.i.d.}{\sim} \mathbb{P}$ . Let  $m$  denote any element of  $[n]$ . Let  $h$  be a core function satisfying  $\mathbb{E}[|h(X_1, \dots, X_m)|^r] < \infty$  with some  $r \in \mathbb{Z}_+$ . Let  $U_n$  and  $V_n$  denote the U and V-statistic associated to  $h$  as defined in (1) and (2), respectively. Then it holds that*

$$\mathbb{E}[|U_n - V_n|^r] = \mathcal{O}(n^{-r}).$$

**Lemma A.2** (Markov inequality). *For a real-valued random variable  $X$  with probability distribution  $\mathbb{P}$  and a  $a > 0$ , it holds that*

$$\mathbb{P}(|X| \geq a) \leq \frac{\mathbb{E}(|X|)}{a}.$$

## A.2 PROOFS

This section is dedicated to proofs. Lemma 4.2 is derived in Section A.2.1. Proposition 4.1 is proved in Section A.2.3 relying on two lemmas shown in Section A.2.2. Lemma 4.4 is proved in Section A.2.5, with an auxiliary result in Section A.2.4.

### A.2.1 Proof of Lemma 4.2

Let  $\mu_k(\tilde{\mathbb{P}}_{n'}) = \sum_{i=1}^{n'} \alpha_k^i \otimes_{m=1}^M \phi_{k_m}(x_m^i)$ , and let  $\mu_{k_m}(\tilde{\mathbb{P}}_{m, n'}) = \sum_{i=1}^{n'} \alpha_{k_m}^i \phi_{k_m}(x_m^i)$  for  $m \in [M]$ . We write

$$\begin{aligned} \text{HSIC}_{k, N}^2(\hat{\mathbb{P}}_n) &= \left\| \mu_k(\tilde{\mathbb{P}}_{n'}) - \otimes_{m=1}^M \mu_{k_m}(\tilde{\mathbb{P}}_{m, n'}) \right\|_{\mathcal{H}_k}^2 \\ &= \underbrace{\left\| \mu_k(\tilde{\mathbb{P}}_{n'}) \right\|_{\mathcal{H}_k}^2}_{=: A} - 2 \cdot \underbrace{\left\langle \mu_k(\tilde{\mathbb{P}}_{n'}), \otimes_{m=1}^M \mu_{k_m}(\tilde{\mathbb{P}}_{m, n'}) \right\rangle_{\mathcal{H}_k}}_{=: C} + \underbrace{\left\| \otimes_{m=1}^M \mu_{k_m}(\tilde{\mathbb{P}}_{m, n'}) \right\|_{\mathcal{H}_k}^2}_{=: B}, \end{aligned}$$

and continue term-by-term. Using the definition of the tensor product, we have for term  $A$  that

$$\begin{aligned} A &= \left\langle \mu_k(\tilde{\mathbb{P}}_{n'}), \mu_k(\tilde{\mathbb{P}}_{n'}) \right\rangle_{\mathcal{H}_k} = \sum_{i=1}^{n'} \sum_{j=1}^{n'} \alpha_k^i \alpha_k^j \left\langle \otimes_{m=1}^M \phi_{k_m}(x_m^i), \otimes_{m=1}^M \phi_{k_m}(x_m^j) \right\rangle_{\mathcal{H}_k} = \sum_{i=1}^{n'} \sum_{j=1}^{n'} \alpha_k^i \alpha_k^j \prod_{m=1}^M k_m(x_m^i, x_m^j) \\ &= \alpha_k^\top (\circ_{m=1}^M \mathbf{K}_{k_m}) \alpha_k. \end{aligned}$$

Similarly, we obtain for term  $B$  that

$$\begin{aligned} B &= \left\langle \otimes_{m=1}^M \mu_{k_m}(\tilde{\mathbb{P}}_{m, n'}), \otimes_{m=1}^M \mu_{k_m}(\tilde{\mathbb{P}}_{m, n'}) \right\rangle_{\mathcal{H}_k} \\ &= \left\langle \otimes_{m=1}^M \sum_{i^{(m)}=1}^{n'} \alpha_{k_m}^{i^{(m)}} \phi_{k_m}(x_m^{i^{(m)}}), \otimes_{m=1}^M \sum_{j^{(m)}=1}^{n'} \alpha_{k_m}^{j^{(m)}} \phi_{k_m}(x_m^{j^{(m)}}) \right\rangle_{\mathcal{H}_k} \\ &\stackrel{(*)}{=} \prod_{m=1}^M \sum_{i^{(m)}=1}^{n'} \sum_{j^{(m)}=1}^{n'} \alpha_{k_m}^{i^{(m)}} \alpha_{k_m}^{j^{(m)}} k_m(x_m^{i^{(m)}}, x_m^{j^{(m)}}) = \prod_{m=1}^M \alpha_{k_m}^\top \mathbf{K}_{k_m} \alpha_{k_m}, \end{aligned}$$

where in  $(*)$  we used (1), the linearity of the inner product, and the reproducing property.

Last, we express term  $C$  as

$$\begin{aligned}
C &= \left\langle \sum_{i=1}^{n'} \alpha_k^i \otimes_{m=1}^M \phi_{k_m}(x_m^i), \otimes_{m=1}^M \sum_{j^{(m)}=1}^{n'} \alpha_{k_m}^{j^{(m)}} \phi_{k_m}(x_m^{j^{(m)}}) \right\rangle_{\mathcal{H}_k} \\
&\stackrel{(a)}{=} \sum_{i=1}^{n'} \alpha_k^i \left\langle \otimes_{m=1}^M \phi_{k_m}(x_m^i), \otimes_{m=1}^M \sum_{j^{(m)}=1}^{n'} \alpha_{k_m}^{j^{(m)}} \phi_{k_m}(x_m^{j^{(m)}}) \right\rangle_{\mathcal{H}_k} \\
&\stackrel{(b)}{=} \sum_{i=1}^{n'} \alpha_k^i \prod_{m \in [M]} \left\langle \phi_{k_m}(x_m^i), \sum_{j^{(m)}=1}^{n'} \alpha_{k_m}^{j^{(m)}} \phi_{k_m}(x_m^{j^{(m)}}) \right\rangle_{\mathcal{H}_k} \\
&\stackrel{(c)}{=} \sum_{i=1}^{n'} \alpha_k^i \prod_{m \in [M]} \sum_{j^{(m)}=1}^{n'} \alpha_{k_m}^{j^{(m)}} \langle \phi_{k_m}(x_m^i), \phi_{k_m}(x_m^{j^{(m)}}) \rangle_{\mathcal{H}_k} \\
&\stackrel{(d)}{=} \sum_{i=1}^{n'} \alpha_k^i \prod_{m \in [M]} \underbrace{\sum_{j^{(m)}=1}^{n'} \alpha_{k_m}^{j^{(m)}} k_m(x_m^i, x_m^{j^{(m)}})}_{(\mathbf{K}_{k_m})_i \alpha_{k_m}} = \alpha_k^\top (\circ_{m=1}^M \mathbf{K}_{k_m} \alpha_{k_m}),
\end{aligned}$$

where (a) follows from the linearity of the inner product, (b) holds by (1), (c) is implied by the linearity of the inner product, (d) is valid by the reproducing property, and we refer to the  $i$ -th row of  $\mathbf{K}_{k_m}$  as  $(\mathbf{K}_{k_m})_i$ .

Substituting terms  $A$ ,  $B$ , and  $C$  concludes the proof.

## A.2.2 Two Lemmas to the Proof of Proposition 4.1

Our main result relies on two lemmas.

**Lemma A.3** (Error bound for Nyström mean embedding of tensor product kernel). *Let  $X = (X_m)_{m=1}^M \in \mathcal{X} = \times_{m=1}^M \mathcal{X}_m$ ,  $X \sim \mathbb{P} \in \mathcal{M}_1^+(\mathcal{X})$ , and  $(\mathcal{X}_m)_{m \in [M]}$  locally compact, second-countable topological spaces. Let  $k_m : \mathcal{X}_m \times \mathcal{X}_m \rightarrow \mathbb{R}$  be a bounded kernel, i.e. there exists  $a_{k_m} \in (0, \infty)$  such that  $\sup_{x_m \in \mathcal{X}_m} \sqrt{k_m(x_m, x_m)} \leq a_{k_m}$  for  $m \in [M]$ . Let  $a_k = \prod_{m=1}^M a_{k_m}$ ,  $k = \otimes_{m=1}^M k_m$ ,  $\mathcal{H}_k$  the RKHS associated to  $k$ ,  $\phi_k = \otimes_{m=1}^M \phi_{k_m}$ ,  $C_k = \mathbb{E}[\phi_k(X) \otimes \phi_k(X)]$ ,  $n' \leq n$ , and  $\tilde{\mathbb{P}}_{n'}$  defined according to (14). Then for any  $\delta \in (0, 1)$  it holds that*

$$\left\| \mu_k(\mathbb{P}) - \mu_k(\tilde{\mathbb{P}}_{n'}) \right\|_{\mathcal{H}_k} \leq \frac{c_{k,1}}{\sqrt{n}} + \frac{c_{k,2}}{n'} + \frac{c_{k,3} \sqrt{\log(n'/\delta)}}{n'} \sqrt{\mathcal{N}_X \left( \frac{12a_k^2 \log(n'/\delta)}{n'} \right)},$$

with probability at least  $1 - \delta$ , provided that

$$n' \geq \max \left( 67, 12a_k^2 \|C_k\|_{\text{op}}^{-1} \right) \log \left( \frac{n'}{\delta} \right),$$

where  $c_{k,1} = 2a_k \sqrt{2 \log(6/\delta)}$ ,  $c_{k,2} = 4\sqrt{3}a_k \log(12/\delta)$ , and  $c_{k,3} = 12\sqrt{3 \log(12/\delta)}a_k$ .

*Proof.* With  $\mathcal{X} = \times_{m \in [M]} \mathcal{X}_m$ , noticing that  $\mathcal{X}$  is locally compact second-countable iff.  $(\mathcal{X}_m)_{m \in [M]}$  are so [Willard, 1970, Theorem 16.2(c), Theorem 18.6],  $\mathcal{H}_k = \otimes_{m=1}^M \mathcal{H}_{k_m}$ ,  $\phi_k = \otimes_{m=1}^M \phi_{k_m}$ , and  $\sqrt{k(x, x)} = \prod_{m=1}^M \sqrt{k_m(x_m, x_m)} \leq a_k$ , the statement is implied by Theorem A.1.  $\square$

*Proof of Lemma 4.3.* To simplify notation, let  $\mu_{k_m} = \mu_{k_m}(\mathbb{P}_m)$ ,  $\tilde{\mu}_{k_m} = \mu_{k_m}(\tilde{\mathbb{P}}_{m, n'})$ ,  $\mathcal{H}_k = \otimes_{m=1}^M \mathcal{H}_{k_m}$ , and  $d_{k_m} = \|\mu_{k_m} - \tilde{\mu}_{k_m}\|_{\mathcal{H}_{k_m}}$ . The proof proceeds by induction on  $M$ : For  $M = 1$  the l.h.s. = r.h.s. =  $\left\| \mu_{k_1}(\mathbb{P}_1) - \mu_{k_1}(\tilde{\mathbb{P}}_{1, n'}) \right\|_{\mathcal{H}_{k_1}}$  is

satisfied, and we assume that the statement holds for  $M = M - 1$ , to obtain

$$\begin{aligned}
& \left\| \otimes_{m=1}^M \mu_{k_m} - \otimes_{m=1}^M \tilde{\mu}_{k_m} \right\|_{\mathcal{H}_k} = \left\| \otimes_{m=1}^{M-1} \mu_{k_m} - \otimes_{m=1}^{M-1} \mu_{k_m} \otimes \tilde{\mu}_{k_M} + \otimes_{m=1}^{M-1} \mu_{k_m} \otimes \tilde{\mu}_{k_M} - \otimes_{m=1}^M \tilde{\mu}_{k_m} \right\|_{\mathcal{H}_k} \\
& = \left\| \otimes_{m=1}^{M-1} \mu_{k_m} \otimes (\mu_{k_M} - \tilde{\mu}_{k_M}) + (\otimes_{m=1}^{M-1} \mu_{k_m} - \otimes_{m=1}^{M-1} \tilde{\mu}_{k_m}) \otimes \tilde{\mu}_{k_M} \right\|_{\mathcal{H}_k} \\
& \stackrel{(a)}{\leq} \left\| \otimes_{m=1}^{M-1} \mu_{k_m} \otimes (\mu_{k_M} - \tilde{\mu}_{k_M}) \right\|_{\mathcal{H}_k} + \left\| (\otimes_{m=1}^{M-1} \mu_{k_m} - \otimes_{m=1}^{M-1} \tilde{\mu}_{k_m}) \otimes \tilde{\mu}_{k_M} \right\|_{\mathcal{H}_k} \\
& \stackrel{(b)}{=} \left( \prod_{m \in [M-1]} \|\mu_{k_m}\|_{\mathcal{H}_{k_m}} \right) d_{k_M} + \left\| \otimes_{m=1}^{M-1} \mu_{k_m} - \otimes_{m=1}^{M-1} \tilde{\mu}_{k_m} \right\|_{\otimes_{m=1}^{M-1} \mathcal{H}_{k_m}} \|\tilde{\mu}_{k_M}\|_{\mathcal{H}_{k_M}} \\
& \stackrel{(c)}{\leq} d_{k_M} \prod_{m \in [M-1]} a_{k_m} + \left\| \otimes_{m=1}^{M-1} \mu_{k_m} - \otimes_{m=1}^{M-1} \tilde{\mu}_{k_m} \right\|_{\otimes_{m=1}^{M-1} \mathcal{H}_{k_m}} (a_{k_M} + d_{k_M}) \\
& \stackrel{(d)}{\leq} d_{k_M} \prod_{m \in [M-1]} a_{k_m} + \left\{ \prod_{m \in [M-1]} (a_{k_m} + d_{k_m}) - \prod_{m \in [M-1]} a_{k_m} \right\} (a_{k_M} + d_{k_M}) \\
& = d_{k_M} \prod_{m \in [M-1]} a_{k_m} + \prod_{m \in [M]} (a_{k_m} + d_{k_m}) - \prod_{m \in [M]} a_{k_m} - d_{k_M} \prod_{m \in [M-1]} a_{k_m} \\
& = \prod_{m \in [M]} (a_{k_m} + d_{k_m}) - \prod_{m \in [M]} a_{k_m},
\end{aligned}$$

where (a) holds by the triangle inequality, (b) is implied by (2) and the definition of  $d_{k_M}$ , (c) follows from

$$\begin{aligned}
\|\mu_{k_m}\|_{\mathcal{H}_{k_m}} &= \left\| \int_{\mathcal{X}_m} k_m(\cdot, x_m) d\mathbb{P}_m(x_m) \right\|_{\mathcal{H}_{k_m}} \stackrel{(e)}{\leq} \int_{\mathcal{X}_m} \underbrace{\|k_m(\cdot, x_m)\|_{\mathcal{H}_{k_m}}}_{\stackrel{(f)}{=} \sqrt{k_m(x_m, x_m)} \stackrel{(g)}{\leq} a_{k_m}} d\mathbb{P}_m(x_m) \leq a_{k_m}, \quad (3) \\
\|\tilde{\mu}_{k_M}\|_{\mathcal{H}_{k_M}} &= \|\tilde{\mu}_{k_M} - \mu_{k_M} + \mu_{k_M}\|_{\mathcal{H}_{k_M}} \stackrel{(h)}{\leq} \|\tilde{\mu}_{k_M} - \mu_{k_M}\|_{\mathcal{H}_{k_M}} + \|\mu_{k_M}\|_{\mathcal{H}_{k_M}} \stackrel{(i)}{\leq} d_{k_M} + a_{k_M},
\end{aligned}$$

(d) is valid by the induction statement holding for  $M - 1$ , (e) is a property of Bochner integrals, (f) is implied by the reproducing property, (g) comes from the definition of  $a_{k_m}$ , the triangle inequality implies (h), (i) follows from (3) and the definition of  $d_{k_M}$ .  $\square$

### A.2.3 Proof of Proposition 4.1

Let  $k = \otimes_{m=1}^M k_m$ , and let  $\mathcal{H}_k = \otimes_{m=1}^M \mathcal{H}_{k_m}$ . We note that  $\mathcal{X} = \times_{m \in [M]} \mathcal{X}_m$  is locally compact second-countable as  $(\mathcal{X}_m)_{m \in [M]}$  are so [Willard, 1970, Theorem 16.2(c), Theorem 18.6].

We decompose the error of the Nyström approximation as

$$\begin{aligned}
\left| \text{HSIC}_k(\mathbb{P}) - \text{HSIC}_{k,N}(\hat{\mathbb{P}}_n) \right| &= \left| \left\| \mu_k(\mathbb{P}) - \otimes_{m=1}^M \mu_{k_m}(\mathbb{P}_m) \right\|_{\mathcal{H}_k} - \left\| \mu_k(\tilde{\mathbb{P}}_{n'}) - \otimes_{m=1}^M \mu_{k_m}(\tilde{\mathbb{P}}_{m,n'}) \right\|_{\mathcal{H}_k} \right| \\
&\stackrel{(a)}{\leq} \left\| \mu_k(\mathbb{P}) - \otimes_{m=1}^M \mu_{k_m}(\mathbb{P}_m) - \mu_k(\tilde{\mathbb{P}}_{n'}) + \otimes_{m=1}^M \mu_{k_m}(\tilde{\mathbb{P}}_{m,n'}) \right\|_{\mathcal{H}_k} \\
&\stackrel{(b)}{\leq} \underbrace{\left\| \mu_k(\mathbb{P}) - \mu_k(\tilde{\mathbb{P}}_{n'}) \right\|_{\mathcal{H}_k}}_{t_1} + \underbrace{\left\| \otimes_{m=1}^M \mu_{k_m}(\mathbb{P}_m) - \otimes_{m=1}^M \mu_{k_m}(\tilde{\mathbb{P}}_{m,n'}) \right\|_{\mathcal{H}_k}}_{t_2},
\end{aligned}$$

where (a) holds by the reverse triangle inequality, and (b) follows from the triangle inequality.

**First term ( $t_1$ ):** One can bound the error of the first term by Lemma A.3; in other words, for any  $\delta \in (0, 1)$  with probability at least  $(1 - \delta)$  it holds that

$$\left\| \mu_k(\mathbb{P}) - \mu_k(\tilde{\mathbb{P}}_{n'}) \right\|_{\mathcal{H}_k} \leq \frac{c_{k,1}}{\sqrt{n}} + \frac{c_{k,2}}{n'} + \frac{c_{k,3} \sqrt{\log(n'/\delta)}}{n'} \sqrt{\mathcal{N}_{X_m} \left( \frac{12a_{k_m}^2 \log(n'/\delta)}{n'} \right)}$$

provided that  $n' \geq \max\left(67, 12a_k^2 \|C_k\|_{\text{op}}^{-1}\right) \log\left(\frac{n'}{\delta}\right)$ , with the constants  $c_{k,1} = 2a_k\sqrt{2\log(6/\delta)}$ ,  $c_{k,2} = 4\sqrt{3}a_k \log(12/\delta)$ ,  $c_{k,3} = 12\sqrt{3}\log(12/\delta)a_k$ .

**Second term ( $t_2$ ):** Applying Lemma 4.3 to the second term gives

$$\left\| \otimes_{m=1}^M \mu_{k_m}(\mathbb{P}_m) - \otimes_{m=1}^M \mu_{k_m}(\tilde{\mathbb{P}}_{m,n'}) \right\|_{\mathcal{H}_k} \leq \prod_{m \in [M]} \left( a_{k_m} + \left\| \mu_{k_m}(\mathbb{P}_m) - \mu_{k_m}(\tilde{\mathbb{P}}_{m,n'}) \right\|_{\mathcal{H}_{k_m}} \right) - \prod_{m \in [M]} a_{k_m}.$$

We now bound the error of each of the  $M$  factors by Theorem A.1, i.e., for fixed  $m \in [M]$ ; particularly we get that for any  $\delta \in (0, 1)$  with probability at least  $1 - \delta$

$$\begin{aligned} \left\| \mu_{k_m}(\mathbb{P}_m) - \mu_{k_m}(\tilde{\mathbb{P}}_{m,n'}) \right\|_{\mathcal{H}_{k_m}} &\leq \frac{c_{k_m,1}}{\sqrt{n}} + \frac{c_{k_m,2}}{n'} + \frac{c_{k_m,3}\sqrt{\log(n'/\delta)}}{n'} \sqrt{\mathcal{N}_{X_m}\left(\frac{12a_{k_m}^2 \log(n'/\delta)}{n'}\right)}, \text{ hence} \\ a_{k_m} + \left\| \mu_{k_m}(\mathbb{P}_m) - \mu_{k_m}(\tilde{\mathbb{P}}_{m,n'}) \right\|_{\mathcal{H}_{k_m}} &\leq a_{k_m} + \frac{c_{k_m,1}}{\sqrt{n}} + \frac{c_{k_m,2}}{n'} + \frac{c_{k_m,3}\sqrt{\log(n'/\delta)}}{n'} \sqrt{\mathcal{N}_{X_m}\left(\frac{12a_{k_m}^2 \log(n'/\delta)}{n'}\right)}, \end{aligned}$$

and by union bound that their product is for any  $\delta \in (0, \frac{1}{M})$  with probability at least  $1 - M\delta$

$$\begin{aligned} \prod_{m \in [M]} \left[ a_{k_m} + \left\| \mu_{k_m}(\mathbb{P}_m) - \mu_{k_m}(\tilde{\mathbb{P}}_{m,n'}) \right\|_{\mathcal{H}_{k_m}} \right] &\leq \\ &\leq \prod_{m \in [M]} \left[ a_{k_m} + \frac{c_{k_m,1}}{\sqrt{n}} + \frac{c_{k_m,2}}{n'} + \frac{c_{k_m,3}\sqrt{\log(n'/\delta)}}{n'} \sqrt{\mathcal{N}_{X_m}\left(\frac{12a_{k_m}^2 \log(n'/\delta)}{n'}\right)} \right], \\ \prod_{m \in [M]} \left[ a_{k_m} + \left\| \mu_{k_m}(\mathbb{P}_m) - \mu_{k_m}(\tilde{\mathbb{P}}_{m,n'}) \right\|_{\mathcal{H}_{k_m}} \right] - \prod_{m \in [M]} a_{k_m} &\leq \\ &\leq \prod_{m \in [M]} \left[ a_{k_m} + \frac{c_{k_m,1}}{\sqrt{n}} + \frac{c_{k_m,2}}{n'} + \frac{c_{k_m,3}\sqrt{\log(n'/\delta)}}{n'} \sqrt{\mathcal{N}_{X_m}\left(\frac{12a_{k_m}^2 \log(n'/\delta)}{n'}\right)} \right] - \prod_{m \in [M]} a_{k_m}, \end{aligned}$$

provided that  $n' \geq \max\left(67, 12a_{k_m}^2 \|C_{k_m}\|_{\text{op}}^{-1}\right) \log\left(\frac{n'}{\delta}\right)$  for all  $m \in [M]$ , with  $C_{k_m} = \mathbb{E}[\phi_{k_m}(X_m) \otimes \phi_{k_m}(X_m)]$  and constants  $c_{k_m,1} = 2a_{k_m}\sqrt{2\log(6/\delta)}$ ,  $c_{k_m,2} = 4\sqrt{3}a_{k_m} \log(12/\delta)$ ,  $c_{k_m,3} = 12\sqrt{3}\log(12/\delta)a_{k_m}$ , with  $m \in [M]$ .

Combining the  $M + 1$  terms by union bound yields the stated result.

#### A.2.4 Lemma to the Proof of Lemma 4.4

**Lemma A.4** (Deviation bound for U-statistics based HSIC estimator). *It holds that*

$$\left| \text{HSIC}_{k,u}^2(\hat{\mathbb{P}}_n) - \text{HSIC}_k^2(\mathbb{P}) \right| = \mathcal{O}_P\left(\frac{1}{\sqrt{n}}\right),$$

where  $\text{HSIC}_{k,u}^2$  is the U-statistic based estimator of  $\text{HSIC}_k^2$ .

*Proof.* We show that (3) can be expressed as a sum of U-statistics and then bound the terms individually. First, square (3) to obtain

$$\begin{aligned} \text{HSIC}_k^2(\mathbb{P}) &= \underbrace{\mathbb{E}_{(x_1, \dots, x_M), (x'_1, \dots, x'_M) \sim \mathbb{P}} \left[ \prod_{m \in [M]} k_m(x_m, x'_m) \right]}_A + \underbrace{\mathbb{E}_{x_1, x'_1 \sim \mathbb{P}_1, \dots, x_M, x'_M \sim \mathbb{P}_M} \left[ \prod_{m \in [M]} k_m(x_m, x'_m) \right]}_B \\ &\quad - 2 \underbrace{\mathbb{E}_{(x_1, \dots, x_M) \sim \mathbb{P}, x'_1 \sim \mathbb{P}_1, \dots, x'_M \sim \mathbb{P}_M} \left[ \prod_{m \in [M]} k_m(x_m, x'_m) \right]}_C, \end{aligned}$$

where  $A$ ,  $B$ , and  $C$  can be estimated by U-statistics  $A'_n$ ,  $B'_n$ , and  $C'_n$ , respectively. Let  $\text{HSIC}_{k,u}^2(\hat{\mathbb{P}}_n) = A'_n + B'_n - 2C'_n$ , and split  $t$  as  $\alpha t + \beta t + (1 - \alpha - \beta)t$ , with  $\alpha, \beta > 0$  and  $\alpha + \beta < 1$ . One obtains

$$P\left(\left|\text{HSIC}_k^2(\mathbb{P}) - \text{HSIC}_{k,u}^2(\hat{\mathbb{P}}_n)\right| \geq t\right) \leq P(|A - A'_n| \geq \alpha t) + P(|B - B'_n| \geq \beta t) + P(2|C - C'_n| \geq (1 - \alpha - \beta)t).$$

Doubling and rewriting Theorem A.2, we have that for U-statistics and any  $\delta \in (0, 1)$

$$\mathbb{P}\left(|U_n - \theta| \geq \sqrt{\frac{m(b-a)^2 \ln(\frac{2}{\delta})}{2n}}\right) \leq \delta.$$

Now, choosing the  $(\theta, U_n, u)$  triplet to be  $(A, A'_n, \alpha t)$ ,  $(B, B'_n, \beta t)$ ,  $(C, C'_n, \frac{(1-\alpha-\beta)t}{2})$ , respectively, setting  $m = 2M$ , and observing that  $a \leq k(x, y) \leq b$  as  $k$  is bounded, we obtain that  $|A'_n - A|\sqrt{n}$ ,  $|B'_n - B|\sqrt{n}$ , and  $|C'_n - C|\sqrt{n}$  are bounded in probability and so is their sum.  $\square$

### A.2.5 Proof of Lemma 4.4

We consider the decomposition

$$\left|\text{HSIC}_k^2(\hat{\mathbb{P}}_n) - \text{HSIC}_k^2(\mathbb{P})\right| \leq \underbrace{\left|\text{HSIC}_k^2(\hat{\mathbb{P}}_n) - \text{HSIC}_{k,u}^2(\hat{\mathbb{P}}_n)\right|}_{t_1} + \underbrace{\left|\text{HSIC}_{k,u}^2(\hat{\mathbb{P}}_n) - \text{HSIC}_k^2(\mathbb{P})\right|}_{t_2}, \quad (4)$$

by using the triangle inequality, where  $\text{HSIC}_{k,u}$  is the U-statistic based HSIC estimator.

**Second term ( $t_2$ ):** Lemma A.4 establishes that  $t_2 = \mathcal{O}_P\left(\frac{1}{\sqrt{n}}\right)$ .

**First term ( $t_1$ ):** To bound  $t_1$ , first, by Markov's inequality (Lemma A.2) observe that

$$\begin{aligned} \mathbb{P}\left(\left|\text{HSIC}_k^2(\hat{\mathbb{P}}_n) - \text{HSIC}_{k,u}^2(\hat{\mathbb{P}}_n)\right| \geq a\right) &\leq \frac{\mathbb{E}\left(\left|\text{HSIC}_k^2(\hat{\mathbb{P}}_n) - \text{HSIC}_{k,u}^2(\hat{\mathbb{P}}_n)\right|\right)}{\underbrace{a}_{=: \epsilon}}, \\ \mathbb{P}\left(\left|\text{HSIC}_k^2(\hat{\mathbb{P}}_n) - \text{HSIC}_{k,u}^2(\hat{\mathbb{P}}_n)\right| \geq \frac{\mathbb{E}\left(\left|\text{HSIC}_k^2(\hat{\mathbb{P}}_n) - \text{HSIC}_{k,u}^2(\hat{\mathbb{P}}_n)\right|\right)}{\epsilon}\right) &\leq \epsilon, \\ \mathbb{P}\left(\left|\text{HSIC}_k^2(\hat{\mathbb{P}}_n) - \text{HSIC}_{k,u}^2(\hat{\mathbb{P}}_n)\right| < \frac{\mathbb{E}\left(\left|\text{HSIC}_k^2(\hat{\mathbb{P}}_n) - \text{HSIC}_{k,u}^2(\hat{\mathbb{P}}_n)\right|\right)}{\epsilon}\right) &\geq 1 - \epsilon, \\ \mathbb{P}\left(\left|\text{HSIC}_k^2(\hat{\mathbb{P}}_n) - \text{HSIC}_{k,u}^2(\hat{\mathbb{P}}_n)\right| < \frac{C}{n\epsilon}\right) &\stackrel{(*)}{\geq} 1 - \epsilon, \\ \mathbb{P}\left(\left|\text{HSIC}_k^2(\hat{\mathbb{P}}_n) - \text{HSIC}_{k,u}^2(\hat{\mathbb{P}}_n)\right| \geq \frac{C}{n\epsilon}\right) &\leq \epsilon, \end{aligned} \quad (5)$$

for constant  $C > 0$  and  $n$  large enough, where  $(*)$  follows from Lemma A.1 (with  $r = 1$ ). (5) implies that

$$\left|\text{HSIC}_k^2(\hat{\mathbb{P}}_n) - \text{HSIC}_{k,u}^2(\hat{\mathbb{P}}_n)\right| = \mathcal{O}_P\left(\frac{1}{n}\right).$$

**Combining the terms ( $t_1 + t_2$ ):** Combining the obtained results for the two terms, one gets that

$$\begin{aligned} \left|\text{HSIC}_k^2(\hat{\mathbb{P}}_n) - \text{HSIC}_k^2(\mathbb{P})\right| &\stackrel{(4)}{\leq} \left|\text{HSIC}_k^2(\hat{\mathbb{P}}_n) - \text{HSIC}_{k,u}^2(\hat{\mathbb{P}}_n)\right| + \left|\text{HSIC}_{k,u}^2(\hat{\mathbb{P}}_n) - \text{HSIC}_k^2(\mathbb{P})\right| \\ &= \mathcal{O}_P\left(\frac{1}{n}\right) + \mathcal{O}_P\left(\frac{1}{\sqrt{n}}\right) = \mathcal{O}_P\left(\frac{1}{\sqrt{n}}\right). \end{aligned} \quad (6)$$

Hence

$$\begin{aligned} \mathcal{O}_P\left(\frac{1}{\sqrt{n}}\right) &\stackrel{(6)}{\geq} \left| \text{HSIC}_k^2(\hat{\mathbb{P}}_n) - \text{HSIC}_k^2(\mathbb{P}) \right| = \left| \text{HSIC}_k(\hat{\mathbb{P}}_n) - \text{HSIC}_k(\mathbb{P}) \right| \left| \underbrace{\text{HSIC}_k(\hat{\mathbb{P}}_n)}_{\stackrel{(5)}{\geq 0}} + \underbrace{\text{HSIC}_k(\mathbb{P})}_{\geq 0} \right| \\ &\geq \left| \text{HSIC}_k(\hat{\mathbb{P}}_n) - \text{HSIC}_k(\mathbb{P}) \right| \text{HSIC}_k(\mathbb{P}), \end{aligned}$$

which by dividing with the constant  $\text{HSIC}_k(\mathbb{P}) > 0$  implies the statement.

## References

Antoine Chatalic, Nicolas Schreuder, Alessandro Rudi, and Lorenzo Rosasco. Nyström kernel mean embeddings. In *International Conference on Machine Learning (ICML)*, pages 3006–3024, 2022.

Robert J. Serfling. *Approximation theorems of mathematical statistics*. John Wiley & Sons, 1980.

Stephen Willard. *General Topology*. Addison-Wesley, 1970.