# An Adaptive Test of Independence with Analytic Kernel Embeddings

**Wittawat Jitkrittum** [1]   **Zoltán Szabó** [2]   **Arthur Gretton** [1]

## Abstract

A new computationally efficient dependence measure, and an adaptive statistical test of independence, are proposed. The dependence measure is the difference between analytic embeddings of the joint distribution and the product of the marginals, evaluated at a finite set of locations (features). These features are chosen so as to maximize a lower bound on the test power, resulting in a test that is data-efficient, and that runs in linear time (with respect to the sample size $n$). The optimized features can be interpreted as evidence to reject the null hypothesis, indicating regions in the joint domain where the joint distribution and the product of the marginals differ most. Consistency of the independence test is established, for an appropriate choice of features. In real-world benchmarks, independence tests using the optimized features perform comparably to the state-of-the-art quadratic-time HSIC test, and outperform competing $\mathcal{O}(n)$ and $\mathcal{O}(n \log n)$ tests.

## 1. Introduction

We consider the design of adaptive, nonparametric statistical tests of dependence: that is, tests of whether a joint distribution $P_{xy}$ factorizes into the product of marginals $P_x P_y$ with the null hypothesis that $H_0 : X$ and $Y$ are independent. While classical tests of dependence, such as Pearson's correlation and Kendall's $\tau$, are able to detect monotonic relations between univariate variables, more modern tests can address complex interactions, for instance changes in variance of $X$ with the value of $Y$. Key to many recent tests is to examine covariance or correlation between data features. These interactions become significantly harder to detect, and the features are more difficult to design, when the data reside in high dimensions.

Zoltán Szabó's ORCID ID: 0000-0001-6183-7603. Arthur Gretton's ORCID ID: 0000-0003-3169-7624. [1]Gatsby Unit, University College London, UK. [2]CMAP, École Polytechnique, France. Correspondence to: Wittawat Jitkrittum <wittawatj@gmail.com>.

A basic nonlinear dependence measure is the Hilbert-Schmidt Independence Criterion (HSIC), which is the Hilbert-Schmidt norm of the covariance operator between feature mappings of the random variables (Gretton et al., 2005; 2008). Each random variable $X$ and $Y$ is mapped to a respective reproducing kernel Hilbert space $\mathcal{H}_k$ and $\mathcal{H}_l$. For sufficiently rich mappings, the covariance operator norm is zero if and only if the variables are independent. A second basic nonlinear dependence measure is the smoothed difference between the characteristic function of the joint distribution, and that of the product of marginals. When a particular smoothing function is used, the statistic corresponds to the covariance between distances of $X$ and $Y$ variable pairs (Feuerverger, 1993; Székely et al., 2007; Székely & Rizzo, 2009), yielding a simple test statistic based on pairwise distances. It has been shown by Sejdinovic et al. (2013) that the distance covariance (and its generalization to semi-metrics) is an instance of HSIC for an appropriate choice of kernels. A disadvantage of these feature covariance statistics, however, is that they require quadratic time to compute (besides in the special case of the distance covariance with univariate real-valued variables, where Huo & Székely (2016) achieve an $\mathcal{O}(n \log n)$ cost). Moreover, the feature covariance statistics have intractable null distributions, and either a permutation approach or the solution of an expensive eigenvalue problem (e.g. Zhang et al., 2011) is required for consistent estimation of the quantiles. Several approaches were proposed by Zhang et al. (2017) to obtain faster tests along the lines of HSIC. These include computing HSIC on finite-dimensional feature mappings chosen as random Fourier features (RFFs) (Rahimi & Recht, 2008), a block-averaged statistic, and a Nyström approximation to the statistic. Key to each of these approaches is a more efficient computation of the statistic and its threshold under the null distribution: for RFFs, the null distribution is a finite weighted sum of $\chi^2$ variables; for the block-averaged statistic, the null distribution is asymptotically normal; for Nyström, either a permutation approach is employed, or the spectrum of the Nyström approximation to the kernel matrix is used in approximating the null distribution. Each of these methods costs significantly less than the $\mathcal{O}(n^2)$ cost of the full HSIC (the cost is linear in $n$, but also depends quadratically on the number of features retained). A potential disadvantage of the Nyström and Fourier approaches is that the features are not optimized to maximize test power,

but are chosen randomly. The block statistic performs worse than both, due to the large variance of the statistic under the null (which can be mitigated by observing more data).

In addition to feature covariances, correlation measures have also been developed in infinite dimensional feature spaces: in particular, Bach & Jordan (2002); Fukumizu et al. (2008) proposed statistics on the correlation operator in a reproducing kernel Hilbert space. While convergence has been established for certain of these statistics, their computational cost is high at $\mathcal{O}(n^3)$, and test thresholds have relied on permutation. A number of much faster approaches to testing based on feature correlations have been proposed, however. For instance, Dauxois & Nkiet (1998) compute statistics of the correlation between finite sets of basis functions, chosen for instance to be step functions or low order B-splines. The cost of this approach is $\mathcal{O}(n)$. This idea was extended by Lopez-Paz et al. (2013), who computed the canonical correlation between finite sets of basis functions chosen as random Fourier features; in addition, they performed a copula transform on the inputs, with a total cost of $\mathcal{O}(n \log n)$. Finally, space partitioning approaches have also been proposed, based on statistics such as the KL divergence, however these apply only to univariate variables (Heller et al., 2016), or to multivariate variables of low dimension (Gretton & Györfi, 2010) (that said, these tests have other advantages of theoretical interest, notably distribution-independent test thresholds).

The approach we take is most closely related to HSIC on a finite set of features. Our simplest test statistic, the Finite Set Independence Criterion (FSIC), is an average of covariances of analytic functions (i.e., features) defined on each of $X$ and $Y$. A normalized version of the statistic (NFSIC) yields a distribution-independent asymptotic test threshold. We show that our test is consistent, despite a finite number of analytic features being used, via a generalization of arguments in Chwialkowski et al. (2015). As in recent work on two-sample testing by Jitkrittum et al. (2016), our test is *adaptive* in the sense that we choose our features on a held-out validation set to optimize a lower bound on the test power. The design of features for independence testing turns out to be quite different to the case of two-sample testing, however: the task is to find correlated feature *pairs* on the respective marginal domains, rather than attempting to find a single, high-dimensional feature representation on the *tensor product* of the marginals, as we would need to do if we were comparing distributions $P_{xy}$ and $Q_{xy}$. While the use of coupled feature pairs on the marginals entails a smaller feature space dimension, it introduces significant complications in the proof of the lower bound, compared with the two-sample case. We demonstrate the performance of our tests on several challenging artificial and real-world datasets, including detection of dependence between music and its year of appearance, and between videos and captions.

In these experiments, we outperform competing linear and $\mathcal{O}(n \log n)$ time tests.

## 2. Independence Criteria and Statistical Tests

We introduce two test statistics: first, the Finite Set Independence Criterion (FSIC), which builds on the principle that dependence can be measured in terms of the covariance between data features. Next, we propose a normalized version of this statistic (NFSIC), with a simpler asymptotic distribution when $P_{xy} = P_x P_y$. We show how to select features for the latter statistic to maximize a lower bound on the power of its corresponding statistical test.

### 2.1. The Finite Set Independence Criterion

We begin by recalling the Hilbert-Schmidt Independence Criterion (HSIC) as proposed in Gretton et al. (2005), since our unnormalized statistic is built along similar lines. Consider two random variables $X \in \mathcal{X} \subseteq \mathbb{R}^{d_x}$ and $Y \in \mathcal{Y} \subseteq \mathbb{R}^{d_y}$. Denote by $P_{xy}$ the joint distribution between $X$ and $Y$; $P_x$ and $P_y$ are the marginal distributions of $X$ and $Y$. Let $\otimes$ denote the tensor product, such that $(a \otimes b) c = a \langle b, c \rangle$. Assume that $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ and $l : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$ are positive definite kernels associated with reproducing kernel Hilbert spaces (RKHS) $\mathcal{H}_k$ and $\mathcal{H}_l$, respectively. Let $\| \cdot \|_{HS}$ be the norm on the space of $\mathcal{H}_l \to \mathcal{H}_k$ Hilbert-Schmidt operators. Then, HSIC between $X$ and $Y$ is defined as

$$
\begin{aligned}
\mathrm{HSIC}(X, Y) &= \left\| \mu_{xy} - \mu_x \otimes \mu_y \right\|_{\mathrm{HS}}^2 \\
&= \mathbb{E}_{(\mathbf{x},\mathbf{y}),(\mathbf{x}',\mathbf{y}')} \left[ k(\mathbf{x}, \mathbf{x}') l(\mathbf{y}, \mathbf{y}') \right] \\
&\quad + \mathbb{E}_{\mathbf{x}} \mathbb{E}_{\mathbf{x}'} [k(\mathbf{x}, \mathbf{x}')] \mathbb{E}_{\mathbf{y}} \mathbb{E}_{\mathbf{y}'} [l(\mathbf{y}, \mathbf{y}')] \\
&\quad - 2 \mathbb{E}_{(\mathbf{x},\mathbf{y})} \left[ \mathbb{E}_{\mathbf{x}'} [k(\mathbf{x}, \mathbf{x}')] \mathbb{E}_{\mathbf{y}'} [l(\mathbf{y}, \mathbf{y}')] \right], \quad (1)
\end{aligned}
$$

where $\mathbb{E}_{\mathbf{x}} := \mathbb{E}_{\mathbf{x} \sim P_x}$, $\mathbb{E}_{\mathbf{y}} := \mathbb{E}_{\mathbf{y} \sim P_y}$, $\mathbb{E}_{\mathbf{xy}} := \mathbb{E}_{(\mathbf{x},\mathbf{y}) \sim P_{xy}}$, and $\mathbf{x}'$ is an independent copy of $\mathbf{x}$. The mean embedding of $P_{xy}$ belongs to the space of Hilbert-Schmidt operators from $\mathcal{H}_l$ to $\mathcal{H}_k$, $\mu_{xy} := \int_{\mathcal{X} \times \mathcal{Y}} k(\mathbf{x}, \cdot) \otimes l(\mathbf{y}, \cdot) \, \mathrm{d}P_{xy}(\mathbf{x}, \mathbf{y}) \in \mathrm{HS}(\mathcal{H}_l, \mathcal{H}_k)$, and the marginal mean embeddings are $\mu_x := \int_{\mathcal{X}} k(\mathbf{x}, \cdot) \, \mathrm{d}P_x(\mathbf{x}) \in \mathcal{H}_k$ and $\mu_y := \int_{\mathcal{Y}} l(\mathbf{y}, \cdot) \, \mathrm{d}P_y(\mathbf{y}) \in \mathcal{H}_l$ (Smola et al., 2007). Gretton et al. (2005, Theorem 4) show that if the kernels $k$ and $l$ are universal (Steinwart & Christmann, 2008) on compact domains $\mathcal{X}$ and $\mathcal{Y}$, then $\mathrm{HSIC}(X, Y) = 0$ if and only if $X$ and $Y$ are independent. Given a joint sample $\mathsf{Z}_n = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n \sim P_{xy}$, an empirical estimator of HSIC can be computed in $\mathcal{O}(n^2)$ time by replacing the population expectations in (1) with their corresponding empirical expectations based on $\mathsf{Z}_n$.

We now propose our new linear-time dependence measure, the Finite Set Independence Criterion (FSIC). Let $\mathcal{X} \subseteq \mathbb{R}^{d_x}$ and $\mathcal{Y} \subseteq \mathbb{R}^{d_y}$ be open sets. Let $\mu_x \mu_y(\mathbf{x}, \mathbf{y}) := \mu_x(\mathbf{x}) \mu_y(\mathbf{y})$ The idea is to see $\mu_{xy}(\mathbf{v}, \mathbf{w}) = \mathbb{E}_{\mathbf{xy}}[k(\mathbf{x}, \mathbf{v}) l(\mathbf{y}, \mathbf{w})]$, $\mu_x(\mathbf{v}) = \mathbb{E}_{\mathbf{x}}[k(\mathbf{x}, \mathbf{v})]$ and $\mu_y(\mathbf{w}) = \mathbb{E}_{\mathbf{y}}[l(\mathbf{y}, \mathbf{w})]$ as smooth functions, and consider a new dis-

tance between $\mu_{xy}$ and $\mu_x\mu_y$ instead of a Hilbert-Schmidt distance as in HSIC (Gretton et al., 2005). The new measure is given by the average of squared differences between $\mu_{xy}$ and $\mu_x\mu_y$, evaluated at $J$ random test locations $V_J := \{(\mathbf{v}_i, \mathbf{w}_i)\}_{i=1}^J \subset \mathcal{X} \times \mathcal{Y}$.

$$\mathrm{FSIC}^2(X, Y) := \frac{1}{J} \sum_{i=1}^J [\mu_{xy}(\mathbf{v}_i, \mathbf{w}_i) - \mu_x(\mathbf{v}_i)\mu_y(\mathbf{w}_i)]^2$$
$$= \frac{1}{J} \sum_{i=1}^J u^2(\mathbf{v}_i, \mathbf{w}_i) = \frac{1}{J} \|\mathbf{u}\|_2^2,$$

where

$$u(\mathbf{v}, \mathbf{w}) := \mu_{xy}(\mathbf{v}, \mathbf{w}) - \mu_x(\mathbf{v})\mu_y(\mathbf{w})$$
$$= \mathbb{E}_{\mathbf{xy}}[k(\mathbf{x}, \mathbf{v})l(\mathbf{y}, \mathbf{w})] - \mathbb{E}_{\mathbf{x}}[k(\mathbf{x}, \mathbf{v})]\mathbb{E}_{\mathbf{y}}[l(\mathbf{y}, \mathbf{w})], \quad (2)$$
$$= \mathrm{cov}_{\mathbf{xy}}[k(\mathbf{x}, \mathbf{v}), l(\mathbf{y}, \mathbf{w})],$$

$\mathbf{u} := (u(\mathbf{v}_1, \mathbf{w}_1), \ldots, u(\mathbf{v}_J, \mathbf{w}_J))^\top$, and $\{(\mathbf{v}_i, \mathbf{w}_i)\}_{i=1}^J$ are realizations from an absolutely continuous distribution (wrt the Lebesgue measure).

Our first result in Proposition 2 states that $\mathrm{FSIC}(X, Y)$ almost surely defines a dependence measure for the random variables $X$ and $Y$, provided that the product kernel on the joint space $\mathcal{X} \times \mathcal{Y}$ is characteristic and analytic (see Definition 1).

**Definition 1** (Analytic kernels (Chwialkowski et al., 2015)). Let $\mathcal{X}$ be an open set in $\mathbb{R}^d$. A positive definite kernel $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is said to be analytic on its domain $\mathcal{X} \times \mathcal{X}$ if for all $\mathbf{v} \in \mathcal{X}$, $f(\mathbf{x}) := k(\mathbf{x}, \mathbf{v})$ is an analytic function on $\mathcal{X}$.

**Assumption A.** *The kernels* $k : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ *and* $l : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$ *are bounded by* $B_k$ *and* $B_l$ *respectively* $[\sup_{\mathbf{x}, \mathbf{x}' \in \mathcal{X}} k(\mathbf{x}, \mathbf{x}') \leq B_k, \sup_{\mathbf{y}, \mathbf{y}' \in \mathcal{Y}} l(\mathbf{y}, \mathbf{y}') \leq B_l]$, *and the product kernel* $g((\mathbf{x}, \mathbf{y}), (\mathbf{x}', \mathbf{y}')) := k(\mathbf{x}, \mathbf{x}')l(\mathbf{y}, \mathbf{y}')$ *is characteristic (Sriperumbudur et al., 2010, Definition 6), and analytic (Definition 1) on* $(\mathcal{X} \times \mathcal{Y}) \times (\mathcal{X} \times \mathcal{Y})$.

**Proposition 2** (FSIC is a dependence measure). *Assume that assumption A holds, and that the test locations* $V_J = \{(\mathbf{v}_i, \mathbf{w}_i)\}_{i=1}^J$ *are drawn from an absolutely continuous distribution* $\eta$. *Then,* $\eta$-*almost surely, it holds that* $\mathrm{FSIC}(X, Y) = \frac{1}{\sqrt{J}}\|\mathbf{u}\|_2 = 0$ *if and only if* $X$ *and* $Y$ *are independent.*

*Proof.* Since $g$ is characteristic, the mean embedding map $\Pi_g : P \mapsto \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim P}[g((\mathbf{x}, \mathbf{y}), \cdot)]$ is injective (Sriperumbudur et al., 2010, Section 3), where $P$ is a probability distribution on $\mathcal{X} \times \mathcal{Y}$. Since $g$ is analytic, by Lemma 10 (Appendix), $\mu_{xy}$ and $\mu_x\mu_y$ are analytic functions. Thus, Lemma 11 (Appendix, setting $\Lambda = \Pi_g$) guarantees that $\mathrm{FSIC}(X, Y) = 0 \iff P_{xy} = P_x P_y \iff X$ and $Y$ are independent almost surely. □

FSIC uses $\mu_{xy}$ as a proxy for $P_{xy}$, and $\mu_x\mu_y$ as a proxy for $P_x P_y$. Proposition 2 states that, to detect the dependence between $X$ and $Y$, it is sufficient to evaluate the difference of the population joint embedding $\mu_{xy}$ and the embedding of the product of the marginal distributions $\mu_x\mu_y$ at a finite number of locations (defined by $V_J$). The intuitive explanation of this property is as follows. If $P_{xy} = P_x P_y$, then $u(\mathbf{v}, \mathbf{w}) = 0$ everywhere, and $\mathrm{FSIC}(X, Y) = 0$ for any $V_J$. If $P_{xy} \neq P_x P_y$, then $u$ will not be a zero function, since the mean embedding map is injective (requires the product kernel to be characteristic). Using the same argument as in Chwialkowski et al. (2015), since $k$ and $l$ are analytic, $u$ is also analytic, and the set of roots $R_u := \{(\mathbf{v}, \mathbf{w}) \mid u(\mathbf{v}, \mathbf{w}) = 0\}$ has Lebesgue measure zero. Thus, it is sufficient to draw $(\mathbf{v}, \mathbf{w})$ from an absolutely continuous distribution to have $(\mathbf{v}, \mathbf{w}) \notin R_u$ $\eta$-almost surely, and hence $\mathrm{FSIC}(X, Y) > 0$. We note that a characteristic kernel which is not analytic may produce $u$ such that $R_u$ has a positive Lebesgue measure. In this case, there is a positive probability that $(\mathbf{v}, \mathbf{w}) \in R_u$, resulting in a potential failure to detect the dependence.

The next proposition shows that Gaussian kernels $k$ and $l$ yield a product kernel which is characteristic and analytic; in other words, this is an example when Assumption A holds.

**Proposition 3** (A product of Gaussian kernels is characteristic and analytic). *Let* $k(\mathbf{x}, \mathbf{x}') = \exp(-(\mathbf{x} - \mathbf{x}')^\top \mathbf{A}(\mathbf{x} - \mathbf{x}'))$ *and* $l(\mathbf{y}, \mathbf{y}') = \exp(-(\mathbf{y} - \mathbf{y}')^\top \mathbf{B}(\mathbf{y} - \mathbf{y}'))$ *be Gaussian kernels on* $\mathbb{R}^{d_x} \times \mathbb{R}^{d_x}$ *and* $\mathbb{R}^{d_y} \times \mathbb{R}^{d_y}$ *respectively, for positive definite matrices* $\mathbf{A}$ *and* $\mathbf{B}$. *Then,* $g((\mathbf{x}, \mathbf{y}), (\mathbf{x}', \mathbf{y}')) = k(\mathbf{x}, \mathbf{x}')l(\mathbf{y}, \mathbf{y}')$ *is characteristic and analytic on* $(\mathbb{R}^{d_x} \times \mathbb{R}^{d_y}) \times (\mathbb{R}^{d_x} \times \mathbb{R}^{d_y})$.

*Proof (sketch).* The main idea is to use the fact that a Gaussian kernel is analytic, and a product of Gaussian kernels is a Gaussian kernel on the pair of variables. See the full proof in Appendix D. □

**Plug-in Estimator** Assume that we observe a joint sample $\mathsf{Z}_n := \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n \overset{i.i.d.}{\sim} P_{xy}$. Unbiased estimators of $\mu_{xy}(\mathbf{v}, \mathbf{w})$ and $\mu_x\mu_y(\mathbf{v}, \mathbf{w})$ are $\hat{\mu}_{xy}(\mathbf{v}, \mathbf{w}) := \frac{1}{n} \sum_{i=1}^n k(\mathbf{x}_i, \mathbf{v})l(\mathbf{y}_i, \mathbf{w})$ and $\widehat{\mu_x\mu_y}(\mathbf{v}, \mathbf{w}) := \frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j \neq i} k(\mathbf{x}_i, \mathbf{v})l(\mathbf{y}_j, \mathbf{w})$, respectively. A straightforward empirical estimator of $\mathrm{FSIC}^2$ is then given by

$$\widehat{\mathrm{FSIC}^2}(\mathsf{Z}_n) = \frac{1}{J} \sum_{i=1}^J \hat{u}(\mathbf{v}_i, \mathbf{w}_i)^2,$$
$$\hat{u}(\mathbf{v}, \mathbf{w}) := \hat{\mu}_{xy}(\mathbf{v}, \mathbf{w}) - \widehat{\mu_x\mu_y}(\mathbf{v}, \mathbf{w}) \quad (3)$$
$$= \frac{2}{n(n-1)} \sum_{i<j} h_{(\mathbf{v}, \mathbf{w})}((\mathbf{x}_i, \mathbf{y}_i), (\mathbf{x}_j, \mathbf{y}_j)), \quad (4)$$

where $h_{(\mathbf{v}, \mathbf{w})}((\mathbf{x}, \mathbf{y}), (\mathbf{x}', \mathbf{y}')) := \frac{1}{2}(k(\mathbf{x}, \mathbf{v}) - k(\mathbf{x}', \mathbf{v}))(l(\mathbf{y}, \mathbf{w}) - l(\mathbf{y}', \mathbf{w}))$. For conciseness, we

define $\hat{\mathbf{u}} := (\hat{u}_1, \ldots, \hat{u}_J)^\top \in \mathbb{R}^J$ where $\hat{u}_i := \hat{u}(\mathbf{v}_i, \mathbf{w}_i)$ so that $\widehat{\text{FSIC}^2}(\mathsf{Z}_n) = \frac{1}{J}\hat{\mathbf{u}}^\top\hat{\mathbf{u}}$.

$\widehat{\text{FSIC}^2}$ can be efficiently computed in $\mathcal{O}((d_x + d_y)Jn)$ time which is linear in $n$ [see (3) which does not have nested double sums], assuming that the runtime complexity of evaluating $k(\mathbf{x}, \mathbf{v})$ is $\mathcal{O}(d_x)$ and that of $l(\mathbf{y}, \mathbf{w})$ is $\mathcal{O}(d_y)$.

Since FSIC satisfies $\text{FSIC}(X, Y) = 0 \iff X \perp Y$, in principle its empirical estimator can be used as a test statistic for an independence test proposing a null hypothesis $H_0$ : "$X$ and $Y$ are independent" against an alternative $H_1$ : "$X$ and $Y$ are dependent." The null distribution (i.e., distribution of the test statistic assuming that $H_0$ is true) is challenging to obtain, however, and depends on the unknown $P_{xy}$. This prompts us to consider a normalized version of FSIC whose asymptotic null distribution takes a more convenient form. We first derive the asymptotic distribution of $\hat{\mathbf{u}}$ in Proposition 4, which we use to derive the normalized test statistic in Theorem 5. As a shorthand, we write $\mathbf{z} := (\mathbf{x}, \mathbf{y})$, $\mathbf{t} := (\mathbf{v}, \mathbf{w})$, $\text{cov}_\mathbf{z}$ is covariance, $\mathbb{V}_\mathbf{z}$ stands for variance.

**Proposition 4** (Asymptotic distribution of $\hat{\mathbf{u}}$). *Define $\mathbf{u} := (u(\mathbf{t}_1), \ldots, u(\mathbf{t}_J))^\top$, $\tilde{k}(\mathbf{x}, \mathbf{v}) := k(\mathbf{x}, \mathbf{v}) - \mathbb{E}_{\mathbf{x}'}k(\mathbf{x}', \mathbf{v})$, and $\tilde{l}(\mathbf{y}, \mathbf{w}) := l(\mathbf{y}, \mathbf{w}) - \mathbb{E}_{\mathbf{y}'}l(\mathbf{y}', \mathbf{w})$. Let $\mathbf{\Sigma} = [\Sigma_{ij}] \in \mathbb{R}^{J \times J}$ be the positive semi-definite matrix with entries $\Sigma_{ij} = \text{cov}_\mathbf{z}(\hat{u}(\mathbf{t}_i), \hat{u}(\mathbf{t}_j)) = \mathbb{E}_{\mathbf{xy}}[\tilde{k}(\mathbf{x}, \mathbf{v}_i)\tilde{l}(\mathbf{y}, \mathbf{w}_i)\tilde{k}(\mathbf{x}, \mathbf{v}_j)\tilde{l}(\mathbf{y}, \mathbf{w}_j)] - u(\mathbf{t}_i)u(\mathbf{t}_j)$. Then, under both $H_0$ and $H_1$, for any fixed test locations $\{\mathbf{t}_1, \ldots, \mathbf{t}_J\}$ for which $\mathbf{\Sigma}$ is full rank, and $0 < \mathbb{V}_\mathbf{z}[h_{\mathbf{t}_j}(\mathbf{z})] < \infty$ for $j = 1, \ldots, J$, it holds that $\sqrt{n}(\hat{\mathbf{u}} - \mathbf{u}) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})$.*

*Proof.* For a fixed $\{\mathbf{t}_1, \ldots, \mathbf{t}_J\}$, $\hat{\mathbf{u}}$ is a one-sample second-order multivariate U-statistic with a U-statistic kernel $h_\mathbf{t}$. Thus, by Lehmann (1999, Theorem 6.1.6) and Kowalski & Tu (2008, Section 5.1, Theorem 1), it follows directly that $\sqrt{n}(\hat{\mathbf{u}} - \mathbf{u}) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})$ where we note that $\mathbb{E}_{\mathbf{xy}}[\tilde{k}(\mathbf{x}, \mathbf{v})\tilde{l}(\mathbf{y}, \mathbf{w})] = u(\mathbf{v}, \mathbf{w})$. $\square$

Recall from Proposition 2 that $\mathbf{u} = \mathbf{0}$ holds almost surely under $H_0$. The asymptotic normality described in Proposition 4 implies that $n\widehat{\text{FSIC}^2} = \frac{n}{J}\hat{\mathbf{u}}^\top\hat{\mathbf{u}}$ converges in distribution to a sum of $J$ dependent weighted $\chi^2$ random variables. The dependence comes from the fact that the coordinates $\hat{u}_1 \ldots, \hat{u}_J$ of $\hat{\mathbf{u}}$ all depend on the sample $\mathsf{Z}_n$. This null distribution is not analytically tractable, and requires a large number of simulations to compute the rejection threshold $T_\alpha$ for a given significance value $\alpha$.

## 2.2. Normalized FSIC and Adaptive Test

For the purpose of an independence test, we will consider a normalized variant of $\widehat{\text{FSIC}^2}$, which we call $\widehat{\text{NFSIC}^2}$, whose tractable asymptotic null distribution is $\chi^2(J)$, the chi-squared distribution with $J$ degrees of freedom. We then show that the independence test defined by $\widehat{\text{NFSIC}^2}$ is consistent. These results are given in Theorem 5.

**Theorem 5** (Independence test based on $\widehat{\text{NFSIC}^2}$ is consistent). *Let $\hat{\mathbf{\Sigma}}$ be a consistent estimate of $\mathbf{\Sigma}$ based on the joint sample $\mathsf{Z}_n$, where $\mathbf{\Sigma}$ is defined in Proposition 4. Assume that $V_J = \{(\mathbf{v}_i, \mathbf{w}_i)\}_{i=1}^J \sim \eta$ where $\eta$ is absolutely continuous wrt the Lebesgue measure. The $\widehat{\text{NFSIC}^2}$ statistic is defined as $\hat{\lambda}_n := n\hat{\mathbf{u}}^\top \left(\hat{\mathbf{\Sigma}} + \gamma_n\mathbf{I}\right)^{-1}\hat{\mathbf{u}}$ where $\gamma_n \geq 0$ is a regularization parameter. Assume that*

1. *Assumption A holds.*

2. *$\mathbf{\Sigma}$ is invertible $\eta$-almost surely.*

3. *$\lim_{n \to \infty} \gamma_n = 0$.*

*Then, for any $k, l$ and $V_J$ satisfying the assumptions,*

1. *Under $H_0$, $\hat{\lambda}_n \xrightarrow{d} \chi^2(J)$ as $n \to \infty$.*

2. *Under $H_1$, for any $r \in \mathbb{R}$, $\lim_{n \to \infty} \mathbb{P}\left(\hat{\lambda}_n \geq r\right) = 1$ $\eta$-almost surely. That is, the independence test based on $\widehat{\text{NFSIC}^2}$ is consistent.*

*Proof (sketch).* Under $H_0$, $n\hat{\mathbf{u}}^\top(\hat{\mathbf{\Sigma}} + \gamma_n\mathbf{I})^{-1}\hat{\mathbf{u}}$ asymptotically follows $\chi^2(J)$ because $\sqrt{n}\hat{\mathbf{u}}$ is asymptotically normally distributed (see Proposition 4). Claim 2 builds on the result in Proposition 2 stating that $\mathbf{u} \neq \mathbf{0}$ under $H_1$; it follows using the convergence of $\hat{\mathbf{u}}$ to $\mathbf{u}$. The full proof can be found in Appendix E. $\square$

Theorem 5 states that if $H_1$ holds, the statistic can be arbitrarily large as $n$ increases, allowing $H_0$ to be rejected for any fixed threshold. Asymptotically the test threshold $T_\alpha$ is given by the $(1 - \alpha)$-quantile of $\chi^2(J)$ and is independent of $n$. The assumption on the consistency of $\hat{\mathbf{\Sigma}}$ is required to obtain the asymptotic chi-squared distribution. The regularization parameter $\gamma_n$ is to ensure that $(\hat{\mathbf{\Sigma}} + \gamma_n\mathbf{I})^{-1}$ can be stably computed. In practice, $\gamma_n$ requires no tuning, and can be set to be a very small constant. We emphasize that $J$ need not increase with $n$ for test consistency.

The next proposition states that the computational complexity of the $\widehat{\text{NFSIC}^2}$ estimator is linear in both the input dimension and sample size, and that it can be expressed in terms of the $\mathbf{K} = [K_{ij}] = [k(\mathbf{v}_i, \mathbf{x}_j)] \in \mathbb{R}^{J \times n}, \mathbf{L} = [L_{ij}] = [l(\mathbf{w}_i, \mathbf{y}_j)] \in \mathbb{R}^{J \times n}$ matrices. In contrast to typical kernel methods, a large Gram matrix of size $n \times n$ is not needed to compute $\widehat{\text{NFSIC}^2}$.

**Proposition 6** (An empirical estimator of $\widehat{\text{NFSIC}^2}$). *Let $\mathbf{1}_n := (1, \ldots, 1)^\top \in \mathbb{R}^n$. Denote by $\circ$ the element-wise matrix product. Then,*

1. $\hat{\mathbf{u}} = \frac{(\mathbf{K} \circ \mathbf{L})\mathbf{1}_n}{n-1} - \frac{(\mathbf{K1}_n) \circ (\mathbf{L1}_n)}{n(n-1)}$.

2. *A consistent estimator for $\boldsymbol{\Sigma}$ is $\hat{\boldsymbol{\Sigma}} = \frac{\boldsymbol{\Gamma\Gamma}^\top}{n}$ where*

$$\boldsymbol{\Gamma} := (\mathbf{K} - n^{-1}\mathbf{K1}_n\mathbf{1}_n^\top) \circ (\mathbf{L} - n^{-1}\mathbf{L1}_n\mathbf{1}_n^\top) - \hat{\mathbf{u}}^b\mathbf{1}_n^\top,$$
$$\hat{\mathbf{u}}^b = n^{-1}(\mathbf{K} \circ \mathbf{L})\mathbf{1}_n - n^{-2}(\mathbf{K1}_n) \circ (\mathbf{L1}_n).$$

*Assume that the complexity of the kernel evaluation is linear in the input dimension. Then the test statistic $\hat{\lambda}_n = n\hat{\mathbf{u}}^\top (\hat{\boldsymbol{\Sigma}} + \gamma_n\mathbf{I})^{-1}\hat{\mathbf{u}}$ can be computed in $\mathcal{O}(J^3 + J^2n + (d_x + d_y)Jn)$ time.*

*Proof (sketch).* Claim 1 for $\hat{\mathbf{u}}$ is straightforward. The expression for $\hat{\boldsymbol{\Sigma}}$ in claim 2 follows directly from the asymptotic covariance expression in Proposition 4. The consistency of $\hat{\boldsymbol{\Sigma}}$ can be obtained by noting that the finite sample bound for $\mathbb{P}(\|\hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}\|_F > t)$ decreases as $n$ increases. This is implicitly shown in Appendix F.2.2 and its following sections. □

Although the dependency of the estimator on $J$ is cubic, we empirically observe that only a small value of $J$ is required (see Section 3). The number of test locations $J$ relates to the number of regions in $\mathcal{X} \times \mathcal{Y}$ of $p_{xy}$ and $p_xp_y$ that differ (see Figure 1).

Theorem 5 asserts the consistency of the test for any test locations $V_J$ drawn from an absolutely continuous distribution. In practice, $V_J$ can be further optimized to increase the test power for a fixed sample size. Our final theoretical result gives a lower bound on the test power of $\widehat{\text{NFSIC}}^2$ i.e., the probability of correctly rejecting $H_0$. We will use this lower bound as the objective function to determine $V_J$ and the kernel parameters. Let $\|\cdot\|_F$ be the Frobenius norm.

**Theorem 7** (A lower bound on the test power). *Let $\text{NFSIC}^2(X, Y) := \lambda_n := n\mathbf{u}^\top\boldsymbol{\Sigma}^{-1}\mathbf{u}$. Let $\mathcal{K}$ be a kernel class for $k$, $\mathcal{L}$ be a kernel class for $l$, and $\mathcal{V}$ be a collection with each element being a set of $J$ locations. Assume that*

1. *There exist finite $B_k$ and $B_l$ such that $\sup_{k\in\mathcal{K}}\sup_{\mathbf{x},\mathbf{x}'\in\mathcal{X}}|k(\mathbf{x},\mathbf{x}')| \leq B_k$ and $\sup_{l\in\mathcal{L}}\sup_{\mathbf{y},\mathbf{y}'\in\mathcal{Y}}|l(\mathbf{y},\mathbf{y}')| \leq B_l$.*

2. *$\tilde{c} := \sup_{k\in\mathcal{K}}\sup_{l\in\mathcal{L}}\sup_{V_J\in\mathcal{V}}\|\boldsymbol{\Sigma}^{-1}\|_F < \infty$.*

*Then, for any $k \in \mathcal{K}, l \in \mathcal{L}, V_J \in \mathcal{V}$, and $\lambda_n \geq r$, the test power satisfies $\mathbb{P}\left(\hat{\lambda}_n \geq r\right) \geq L(\lambda_n)$ where*

$$L(\lambda_n) = 1 - 62e^{-\xi_1\gamma_n^2(\lambda_n-r)^2/n} - 2e^{-\lfloor 0.5n \rfloor(\lambda_n-r)^2/[\xi_2 n^2]}$$
$$- 2e^{-[(\lambda_n-r)\gamma_n(n-1)/3 - \xi_3 n - c_3\gamma_n^2 n(n-1)]^2/[\xi_4 n^2(n-1)]},$$
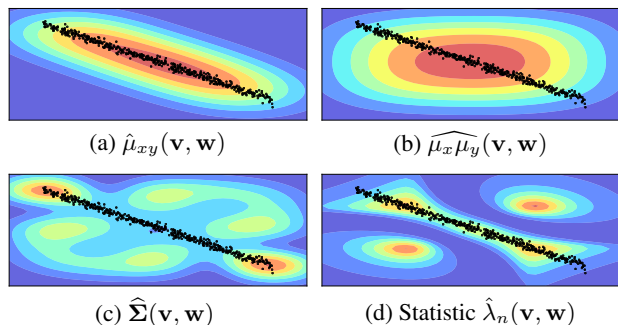
*$\lfloor\cdot\rfloor$ is the floor function, $\xi_1 := \frac{1}{3^2c_1^2J^2B^*}$, $B^*$ is a constant depending on only $B_k$ and $B_l$, $\xi_2 := 72c_2^2JB^2$, $B := B_kB_l$,*

*$\xi_3 := 8c_1B^2J$, $c_3 := 4B^2J\tilde{c}^2$, $\xi_4 := 2^8B^4J^2c_1^2$, $c_1 := 4B^2J\sqrt{J}\tilde{c}$, and $c_2 := 4B\sqrt{J}\tilde{c}$. Moreover, for sufficiently large fixed $n$, $L(\lambda_n)$ is increasing in $\lambda_n$.*

We provide the proof in Appendix F. To put Theorem 7 into perspective, assume that $\mathcal{K} = \left\{(\mathbf{x}, \mathbf{v}) \mapsto \exp\left(-\frac{\|\mathbf{x}-\mathbf{v}\|^2}{2\sigma_x^2}\right) \mid \sigma_x^2 \in [\sigma_{x,l}^2, \sigma_{x,u}^2]\right\} =: \mathcal{K}_g$ for some $0 < \sigma_{x,l}^2 < \sigma_{x,u}^2 < \infty$ and $\mathcal{L} = \left\{(\mathbf{y}, \mathbf{w}) \mapsto \exp\left(-\frac{\|\mathbf{y}-\mathbf{w}\|^2}{2\sigma_y^2}\right) \mid \sigma_y^2 \in [\sigma_{y,l}^2, \sigma_{y,u}^2]\right\} =: \mathcal{L}_g$ for some $0 < \sigma_{y,l}^2 < \sigma_{y,u}^2 < \infty$ are Gaussian kernel classes. Then, in Theorem 7, $B = B_k = B_l = 1$, and $B^* = 2$. The assumption $\tilde{c} < \infty$ is a technical condition to guarantee that the test power lower bound is finite for all $\theta$ defined by the feasible sets $\mathcal{K}, \mathcal{L}$, and $\mathcal{V}$. Let $\mathcal{V}_{\epsilon,r} := \{V_J \mid \|\mathbf{v}_i\|^2, \|\mathbf{w}_i\|^2 \leq r \text{ and } \|\mathbf{v}_i - \mathbf{v}_j\|_2^2 + \|\mathbf{w}_i - \mathbf{w}_j\|_2^2 \geq \epsilon, \text{ for all } i \neq j\}$. If we set $\mathcal{K} = \mathcal{K}_g, \mathcal{L} = \mathcal{L}_g$, and $\mathcal{V} = \mathcal{V}_{\epsilon,r}$ for some $\epsilon, r > 0$, then $\tilde{c} < \infty$ as $\mathcal{K}_g, \mathcal{L}_g$, and $\mathcal{V}_{\epsilon,r}$ are compact. In practice, these conditions do not necessarily create restrictions as they almost always hold implicitly. We show in Appendix C that the objective function used to choose $V_J$ will discourage any two locations to be in the same neighborhood.

**Parameter Tuning** Let $\theta$ be the collection of all tuning parameters of the test. If $k \in \mathcal{K}_g$ and $l \in \mathcal{L}_g$ (i.e., Gaussian kernels), then $\theta = \{\sigma_x^2, \sigma_y^2, V_J\}$. The test power lower bound $L(\lambda_n)$ in Theorem 7 is a function of $\lambda_n = n\mathbf{u}^\top\boldsymbol{\Sigma}^{-1}\mathbf{u}$ which is the population counterpart of the test statistic $\hat{\lambda}_n$. As in FSIC, it can be shown that $\lambda_n = 0$ if and only if $X$ are $Y$ are independent (from Proposition 2). According to Theorem 7, for a sufficiently large $n$, the test power lower bound is increasing in $\lambda_n$. One can therefore think of $\lambda_n$ (a function of $\theta$) as representing how easily the test rejects $H_0$ given a problem $P_{xy}$. The higher the $\lambda_n$, the greater the lower bound on the test power, and thus the more likely it is that the test will reject $H_0$ when it is false.

In light of this reasoning, we propose to set $\theta$ by maximizing the lower bound on the test power i.e., set $\theta$ to $\theta^* = \arg\max_\theta L(\lambda_n)$. Assume that $n$ is sufficiently large so that $\lambda_n \mapsto L(\lambda_n)$ is an increasing function. Then, $\arg\max_\theta L(\lambda_n) = \arg\max_\theta \lambda_n$. That this procedure is also valid under $H_0$ can be seen as follows. Under $H_0$, $\theta^* = \arg\max_\theta 0$ will be arbitrary. Since Theorem 7 guarantees that $\hat{\lambda}_n \xrightarrow{d} \chi^2(J)$ as $n \to \infty$ for any $\theta$, the asymptotic null distribution does not change by using $\theta^*$. In practice, $\lambda_n$ is a population quantity which is unknown. We propose dividing the sample $\mathsf{Z}_n$ into two disjoint sets: training and test sets. The training set is used to compute $\hat{\lambda}_n$ (an estimate of $\lambda_n$) to optimize for $\theta^*$, and the test set is used for the actual independence test with the optimized $\theta^*$. The splitting is to guarantee the independence of $\theta^*$ and the test sample to avoid overfitting.

(a) $\hat{\mu}_{xy}(\mathbf{v}, \mathbf{w})$      (b) $\widehat{\mu_x \mu_y}(\mathbf{v}, \mathbf{w})$

(c) $\hat{\boldsymbol{\Sigma}}(\mathbf{v}, \mathbf{w})$      (d) Statistic $\hat{\lambda}_n(\mathbf{v}, \mathbf{w})$

Figure 1: Illustration of $\widehat{\text{NFSIC}}^2$.

To better understand the behaviour of $\widehat{\text{NFSIC}}^2$, we visualize $\hat{\mu}_{xy}(\mathbf{v}, \mathbf{w})$, $\widehat{\mu_x \mu_y}(\mathbf{v}, \mathbf{w})$ and $\hat{\boldsymbol{\Sigma}}(\mathbf{v}, \mathbf{w})$ as a function of one test location $(\mathbf{v}, \mathbf{w})$ on a simple toy problem. In this problem, $Y = -X + Z$ where $Z \sim \mathcal{N}(0, 0.3^2)$ is an independent noise variable. As we consider only one location ($J = 1$), $\hat{\boldsymbol{\Sigma}}(\mathbf{v}, \mathbf{w})$ is a scalar. The statistic can be written as $\hat{\lambda}_n = n \frac{(\hat{\mu}_{xy}(\mathbf{v}, \mathbf{w}) - \widehat{\mu_x \mu_y}(\mathbf{v}, \mathbf{w}))^2}{\hat{\boldsymbol{\Sigma}}(\mathbf{v}, \mathbf{w})}$. These components are shown in Figure 1, where we use Gaussian kernels for both $X$ and $Y$, and the horizontal and vertical axes correspond to $\mathbf{v} \in \mathbb{R}$ and $\mathbf{w} \in \mathbb{R}$, respectively.

Intuitively, $\hat{u}(\mathbf{v}, \mathbf{w}) = \hat{\mu}_{xy}(\mathbf{v}, \mathbf{w}) - \widehat{\mu_x \mu_y}(\mathbf{v}, \mathbf{w})$ captures the difference of the joint distribution and the product of the marginals as a function of $(\mathbf{v}, \mathbf{w})$. Squaring $\hat{u}(\mathbf{v}, \mathbf{w})$ and dividing it by the variance shown in Figure 1c gives the statistic (also the parameter tuning objective) shown in Figure 1d. The latter figure illustrates that the parameter tuning objective function can be non-convex: non-convexity arises since there are multiple ways to detect the difference between the joint distribution and the product of the marginals. In this case, the lower left and upper right regions equally indicate the largest difference. A convex objective would not be able to capture this phenomenon.

## 3. Experiments

In this section, we empirically study the performance of the proposed method on both toy (Section 3.1) and real problems (Section 3.2). We are interested in challenging problems requiring a large number of samples, where a quadratic-time test might be computationally infeasible. Our goal is not to outperform a quadratic-time test with a linear-time test uniformly over *all* testing problems. We will find, however, that our test does outperform the quadratic-time test in some cases. Code is available at https://github.com/wittawatj/fsic-test.

We compare the proposed NFSIC with optimization (NFSIC-opt) to five multivariate nonparametric tests. The $\widehat{\text{NFSIC}}^2$ test without optimization (NFSIC-med) acts as a baseline, allowing the effect of parameter optimization to be clearly

seen. For pedagogical reason, we consider the original HSIC test of Gretton et al. (2005) denoted by QHSIC, which is a quadratic-time test. Nyström HSIC (NyHSIC) uses a Nyström approximation to the kernel matrices of $X$ and $Y$ when computing the HSIC statistic. FHSIC is another variant of HSIC in which a random Fourier feature approximation (Rahimi & Recht, 2008) to the kernel is used. NyHSIC and FHSIC are studied in Zhang et al. (2017) and can be computed in $\mathcal{O}(n)$, with quadratic dependency on the number of inducing points in NyHSIC, and quadratic dependency on the number of random features in FHSIC. Finally, the Randomized Dependence Coefficient (RDC) proposed in Lopez-Paz et al. (2013) is also considered. The RDC can be seen as the primal form (with random Fourier features) of the kernel canonical correlation analysis of Bach & Jordan (2002) on copula-transformed data. We consider RDC as a linear-time test even though preprocessing by an empirical copula transform costs $\mathcal{O}((d_x + d_y)n \log n)$.

We use Gaussian kernel classes $\mathcal{K}_g$ and $\mathcal{L}_g$ for both $X$ and $Y$ in all the methods. Except NFSIC-opt, all other tests use full sample to conduct the independence test, where the Gaussian widths $\sigma_x$ and $\sigma_y$ are set according to the widely used median heuristic i.e., $\sigma_x = \text{median}(\{\|\mathbf{x}_i - \mathbf{x}_j\|_2 \mid 1 \leq i < j \leq n\})$, and $\sigma_y$ is set in the same way using $\{\mathbf{y}_i\}_{i=1}^n$. The $J$ locations for NFSIC-med are randomly drawn from the standard multivariate normal distribution in each trial. For a sample of size $n$, NFSIC-opt uses half the sample for parameter tuning, and the other disjoint half for the test. We permute the sample 300 times in RDC[1] and HSIC to simulate from the null distribution and compute the test threshold. The null distributions for FHSIC and NyHSIC are given by a finite sum of weighted $\chi^2(1)$ random variables given in Eq. 8 of Zhang et al. (2017). Unless stated otherwise, we set the test threshold of the two NFSIC tests to be the $(1 - \alpha)$-quantile of $\chi^2(J)$. To provide a fair comparison, we set $J = 10$, use 10 inducing points in NyHSIC, and 10 random Fourier features in FHSIC and RDC.

**Optimization of NFSIC-opt** The parameters of NFSIC-opt are $\sigma_x, \sigma_y$, and $J$ locations of size $(d_x + d_y)J$. We treat all the parameters as a long vector in $\mathbb{R}^{2+(d_x+d_y)J}$ and use gradient ascent to optimize $\hat{\lambda}_{n/2}$. We observe that initializing $V_J$ by randomly picking $J$ points from the training sample yields good performance. The regularization parameter $\gamma_n$ in NFSIC is fixed to a small value, and is not optimized. It is worth emphasizing that the complexity of the optimization procedure is still linear-time.[2]

---

[1]We use a permutation test for RDC, following the authors' implementation (https://github.com/lopezpaz/randomized_dependence_coefficient, referred commit: b0ac6c0).

[2]Our claim on linear runtime (with respect to $n$) is for the gradient ascent procedure to find a local optimum for $\theta$. We do not
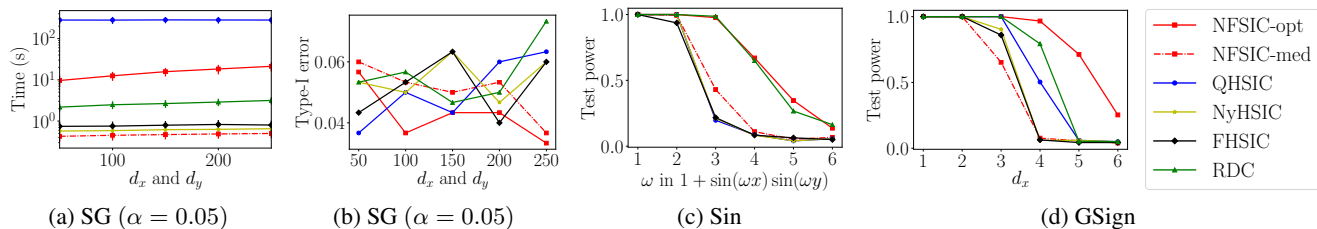
Figure 2: (a): Runtime. (b): Probability of rejecting $H_0$ as problem parameters vary. Fix $n = 4000$.

Since FSIC, NyHFSIC and RDC rely on a finite-dimensional kernel approximation, these tests are consistent only if both the number of features increases with $n$. By constrast, the proposed NFSIC requires only $n$ to go to infinity to achieve consistency i.e., $J$ can be fixed. We refer the reader to Appendix C for a brief investigation of the test power vs. increasing $J$. The test power does not necessarily monotonically increase with $J$.

### 3.1. Toy Problems

We consider three toy problems.

**1. Same Gaussian (SG).** The two variables are independently drawn from the standard multivariate normal distribution i.e., $X \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{d_x})$ and $Y \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{d_y})$ where $\mathbf{I}_d$ is the $d \times d$ identity matrix. This problem represents a case in which $H_0$ holds.

**2. Sinusoid (Sin).** Let $p_{xy}$ be the probability density of $P_{xy}$. In the Sinusoid problem, the dependency of $X$ and $Y$ is characterized by $(X, Y) \sim p_{xy}(x, y) \propto 1 + \sin(\omega x)\sin(\omega y)$, where the domains of $\mathcal{X}, \mathcal{Y} = (-\pi, \pi)$ and $\omega$ is the frequency of the sinusoid. As the frequency $\omega$ increases, the drawn sample becomes more similar to a sample drawn from $\text{Uniform}((-\pi, \pi)^2)$. That is, the higher $\omega$, the harder to detect the dependency between $X$ and $Y$. This problem was studied in Sejdinovic et al. (2013). Plots of the density for a few values of $\omega$ are shown in Figures 6 and 7 in the appendix. The main characteristic of interest in this problem is the local change in the density function.

**3. Gaussian Sign (GSign).** In this problem, $Y = |Z| \prod_{i=1}^{d_x} \text{sgn}(X_i)$, where $X \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{d_x})$, $\text{sgn}(\cdot)$ is the sign function, and $Z \sim \mathcal{N}(0, 1)$ serves as a source of noise. The full interaction of $X = (X_1, \ldots, X_{d_x})$ is what makes the problem challenging. That is, $Y$ is dependent on $X$, yet it is independent of any proper subset of $\{X_1, \ldots, X_d\}$. Thus, simultaneous consideration of all the coordinates of $X$ is required to successfully detect the dependency.

We fix $n = 4000$ and vary the problem parameters. Each problem is repeated for 300 trials, and the sample is redrawn each time. The significance level $\alpha$ is set to 0.05. The re-

claim a linear runtime to find a global optimum.

sults are shown in Figure 2. It can be seen that in the SG problem (Figure 2b) where $H_0$ holds, all the tests achieve roughly correct type-I errors at $\alpha = 0.05$. In particular, we point out that NFSIC-opt's rejection rate is well controlled as the sample used for testing and the sample used for parameter tuning are independent. The rejection rate would have been much higher had we done the optimization and testing on the same sample (i.e., overfitting). In the Sin problem, NFSIC-opt achieves high test power for all considered $\omega = 1, \ldots, 6$, highlighting its strength in detecting local changes in the joint density. The performance of NFSIC-med is significantly lower than that of NFSIC-opt. This phenomenon clearly emphasizes the importance of the optimization to place the locations at the relevant regions in $\mathcal{X} \times \mathcal{Y}$. RDC has a remarkably high performance in both Sin and GSign (Figure 2c, 2d) despite no parameter tuning. The ability to simultaneously consider interacting features of NFSIC-opt is indicated by its superior test power in GSign, especially at the challenging settings of $d_x = 5, 6$.

**NFSIC vs. QHSIC.** We observe that NFSIC-opt outperforms the quadratic-time QHSIC in these two problems. QHSIC is defined as the RKHS norm of the witness function $u$ (see (2)). Intuitively, one can think of the RKHS norm as taking into account all the locations $(\mathbf{v}, \mathbf{w})$. By contrast, the proposed NFSIC evaluates the witness function at $J$ locations. If the differences in $p_{xy}$ and $p_x p_y$ are local (e.g., Sin problem), or there are interacting features (e.g., GSign problem), then only small regions in the space of $(X, Y)$ are relevant in detecting the difference of $p_{xy}$ and $p_x p_y$. In these cases, pinpointing exact test locations by the optimization of NFSIC performs well. On the other hand, taking into account all possible test locations as done implicitly in QHSIC also integrates over regions where the difference between $p_{xy}$ and $p_x p_y$ is small, resulting in a weaker indication of dependence. Whether QHSIC is better than NFSIC depends heavily on the problem, and there is no one best answer. If the difference between $p_{xy}$ and $p_x p_y$ is large only in localized regions, then the proposed linear time statistic has an advantage. If the difference is spatially diffuse, then QHSIC has an advantage. No existing work has proposed a procedure to optimally tune kernel parameters for QHSIC; by contrast, NFSIC has a clearly defined objective for parameter tuning.
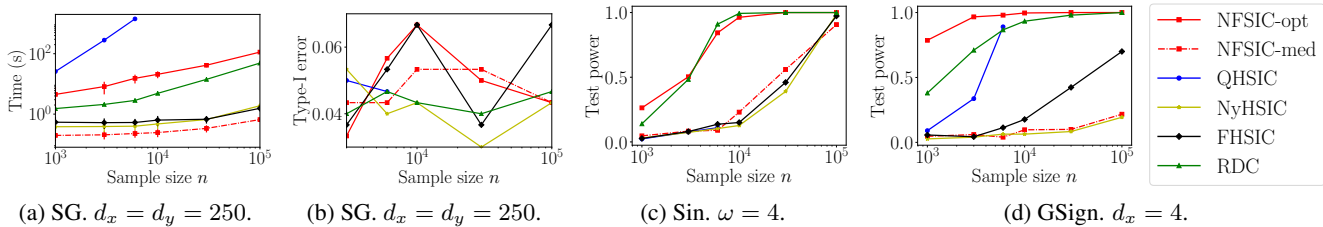
(a) SG. $d_x = d_y = 250$.    (b) SG. $d_x = d_y = 250$.    (c) Sin. $\omega = 4$.    (d) GSign. $d_x = 4$.

Figure 3: (a) Runtime. (b): Probability of rejecting $H_0$ as $n$ increases in the toy problems.

To investigate the sample efficiency of all the tests, we fix $d_x = d_y = 250$ in SG, $\omega = 4$ in Sin, $d_x = 4$ in GSign, and increase $n$. Figure 3 shows the results. The quadratic dependency on $n$ in QHSIC makes it infeasible both in terms of memory and runtime to consider $n$ larger than 6000 (Figure 3a). By contrast, although not the most time-efficient, NFSIC-opt has the highest sample-efficiency for GSign, and for Sin in the low-sample regime, significantly outperforming QHSIC. Despite the small additional overhead from the optimization, we are yet able to conduct an accurate test with $n = 10^5, d_x = d_y = 250$ in less than 100 seconds. We observe in Figure 3b that the two NFSIC variants have correct type-I errors across all sample sizes. We recall from Theorem 5 that the NFSIC test with random test locations will asymptotically reject $H_0$ if it is false. A demonstration of this property is given in Figure 3c, where the test power of NFSIC-med eventually reaches 1 with $n$ higher than $10^5$.

### 3.2. Real Problems

We now examine the performance of our proposed test on real problems.

**Million Song Data (MSD)** We consider a subset of the Million Song Data[3] (Bertin-Mahieux et al., 2011), in which each song $(X)$ out of 515,345 is represented by 90 features, of which 12 features are timbre average (over all segments) of the song, and 78 features are timbre covariance. Most of the songs are western commercial tracks from 1922 to 2011. The goal is to detect the dependency between each song and its year of release $(Y)$. We set $\alpha = 0.01$, and repeat for 300 trials where the full sample is randomly subsampled to $n$ points in each trial. Other settings are the same as in the toy problems. To make sure that the type-I error is correct, we use the permutation approach in the NFSIC tests to compute the threshold. Figure 4b shows the test powers as $n$ increases from 500 to 2000. To simulate the case where $H_0$ holds in the problem, we permute the sample to break the dependency of $X$ and $Y$. The results are shown in Figure 5 in the appendix.

Evidently, NFSIC-opt has the highest test power among all



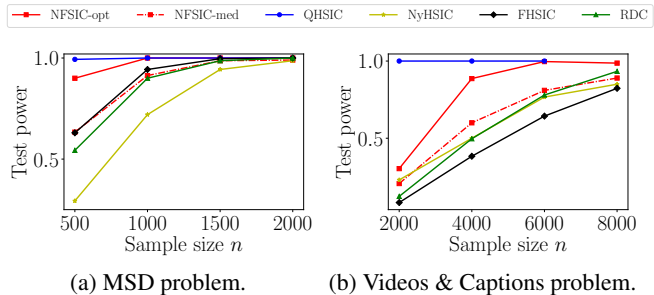(a) MSD problem.      (b) Videos & Captions problem.

Figure 4: Probability of rejecting $H_0$ as $n$ increases in the two real problems. $\alpha = 0.01$.

the linear-time tests for all the sample sizes. Its test power is second to only QHSIC. We recall that NFSIC-opt uses half of the sample for parameter tuning. Thus, at $n = 500$, the actual sample for testing is 250, which is relatively small. The fact that there is a vast power gain from 0.4 (NFSIC-med) to 0.8 (NFSIC-opt) at $n = 500$ suggests that the optimization procedure can perform well even at a lower sample sizes.

**Videos and Captions** Our last problem is based on the VideoStory46K[4] dataset (Habibian et al., 2014). The dataset contains 45,826 Youtube videos $(X)$ of an average length of roughly one minute, and their corresponding text captions $(Y)$ uploaded by the users. Each video is represented as a $d_x = 2000$ dimensional Fisher vector encoding of motion boundary histograms (MBH) descriptors of Wang & Schmid (2013). Each caption is represented as a bag of words with each feature being the frequency of one word. After filtering only words which occur in at least six video captions, we obtain $d_y = 1878$ words. We examine the test powers as $n$ increases from 2000 to 8000. The results are given in Figure 4. The problem is sufficiently challenging that all linear-time tests achieve a low power at $n = 2000$. QHSIC performs exceptionally well on this problem, achieving a maximum power throughout. NFSIC-opt has the highest sample efficiency among the linear-time tests, showing that the optimization procedure is also practical in a high dimensional setting.

---

[3]Million Song Data subset: `https://archive.ics.uci.edu/ml/datasets/YearPredictionMSD`.

[4]VideoStory46K dataset: `https://ivi.fnwi.uva.nl/isis/mediamill/datasets/videostory.php`.

# References

Anderson, Theodore W. *An Introduction to Multivariate Statistical Analysis*. Wiley, 2003.

Bach, Francis R. and Jordan, Michael I. Kernel independent component analysis. *Journal of Machine Learning Research*, 3:1–48, 2002.

Bertin-Mahieux, Thierry, Ellis, Daniel P.W., Whitman, Brian, and Lamere, Paul. The million song dataset. In *International Conference on Music Information Retrieval (ISMIR)*, 2011.

Chwialkowski, Kacper P., Ramdas, Aaditya, Sejdinovic, Dino, and Gretton, Arthur. Fast Two-Sample Testing with Analytic Representations of Probability Measures. In *Advances in Neural Information Processing Systems (NIPS)*, pp. 1981–1989. 2015.

Dauxois, Jacques and Nkiet, Guy Martial. Nonlinear canonical analysis and independence tests. *The Annals of Statistics*, 26(4):1254–1278, 1998.

Feuerverger, Andrey. A consistent test for bivariate dependence. *International Statistical Review*, 61(3):419–433, 1993.

Fukumizu, Kenji, Gretton, Arthur, Sun, Xiaohai, and Schölkopf, Bernhard. Kernel measures of conditional dependence. In *Advances in Neural Information Processing Systems (NIPS)*, pp. 489–496, 2008.

Gretton, Arthur and Györfi, László. Consistent nonparametric tests of independence. *Journal of Machine Learning Research*, 11:1391–1423, 2010.

Gretton, Arthur, Bousquet, Olivier, Smola, Alex, and Schölkopf, Bernhard. Measuring Statistical Dependence with Hilbert-Schmidt Norms. In *Algorithmic Learning Theory (ALT)*, pp. 63–77. 2005.

Gretton, Arthur, Fukumizu, Kenji, Teo, Choon H., Song, Le, Schölkopf, Bernhard, and Smola, Alex J. A Kernel Statistical Test of Independence. In *Advances in Neural Information Processing Systems (NIPS)*, pp. 585–592. 2008.

Habibian, Amirhossein, Mensink, Thomas, and Snoek, Cees GM. Videostory: A new multimedia embedding for few-example recognition and translation of events. In *ACM International Conference on Multimedia*, pp. 17–26, 2014.

Heller, Ruth, Heller, Yair, Kaufman, Shachar, Brill, Barak, and Gorfine, Malka. Consistent distribution-free $k$-sample and independence tests for univariate random variables. *Journal of Machine Learning Research*, 17(29):1–54, 2016.

Huo, Xiaoming and Székely, Gábor J. Fast computing for distance covariance. *Technometrics*, 58(4):435–447, 2016.

Jitkrittum, Wittawat, Szabó, Zoltán, Chwialkowski, Kacper, and Gretton, Arthur. Interpretable Distribution Features with Maximum Testing Power. In *Advances in Neural Information Processing Systems (NIPS)*, pp. 181–189. 2016.

Kowalski, Jeanne and Tu, Xin M. *Modern Applied U-Statistics*. John Wiley & Sons, 2008.

Lehmann, Eric L. *Elements of Large-Sample Theory*. Springer Science & Business Media, 1999.

Lopez-Paz, David, Hennig, Philipp, and Schölkopf, Bernhard. The Randomized Dependence Coefficient. In *Advances in Neural Information Processing Systems (NIPS)*, pp. 1–9. 2013.

Rahimi, Ali and Recht, Benjamin. Random features for large-scale kernel machines. In *Advances in Neural Information Processing Systems (NIPS)*, pp. 1177–1184. 2008.

Sejdinovic, Dino, Sriperumbudur, Bharath, Gretton, Arthur, and Fukumizu, Kenji. Equivalence of distance-based and RKHS-based statistics in hypothesis testing. *The Annals of Statistics*, 41(5):2263–2291, 2013.

Serfling, Robert J. *Approximation Theorems of Mathematical Statistics*. John Wiley & Sons, 2009.

Smola, Alex, Gretton, Arthur, Song, Le, and Schölkopf, Bernhard. A Hilbert space embedding for distributions. In *International Conference on Algorithmic Learning Theory (ALT)*, pp. 13–31, 2007.

Sriperumbudur, Bharath K., Gretton, Arthur, Fukumizu, Kenji, Schölkopf, Bernhard, and Lanckriet, Gert R. G. Hilbert Space Embeddings and Metrics on Probability Measures. *Journal of Machine Learning Research*, 11:1517–1561, 2010.

Steinwart, Ingo and Christmann, Andreas. *Support vector machines*. Springer Science & Business Media, 2008.

Székely, Gábor J. and Rizzo, Maria L. Brownian distance covariance. *The Annals of Applied Statistics*, 3(4):1236–1265, 2009.

Székely, Gábor J., Rizzo, Maria L., and Bakirov, Nail K. Measuring and testing dependence by correlation of distances. *The Annals of Statistics*, 35(6):2769–2794, 2007.

van der Vaart, Aad. *Asymptotic Statistics*. Cambridge University Press, 2000.

Wang, Heng and Schmid, Cordelia. Action recognition with improved trajectories. In *IEEE International Conference on Computer Vision (ICCV)*, pp. 3551–3558, 2013.

Zhang, Kun, Peters, Jonas, Janzing, Dominik, and Schölkopf, Bernhard. Kernel-based conditional independence test and application in causal discovery. In *Conference on Uncertainty in Artificial Intelligence (UAI)*, pp. 804–813, 2011.

Zhang, Qinyi, Filippi, Sarah, Gretton, Arthur, and Sejdinovic, Dino. Large-Scale Kernel Methods for Independence Testing. *Statistics and Computing*, pp. 1–18, 2017.