

Efficiently Learning Significant Fourier Feature Pairs for Statistical Independence Testing

Yixin Ren¹, Yewei Xia^{1,4}, Hao Zhang³, Jihong Guan², Shuigeng Zhou¹

¹Fudan University, ²Tongji University, ³SIAT, ⁴MBZUAI



Statistical Independence Testing

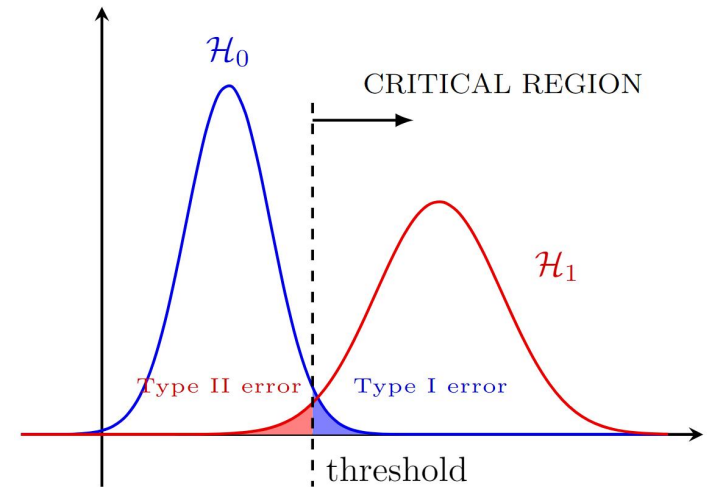
whether $\mathbb{P}_{XY} = \mathbb{P}_X \mathbb{P}_Y$?

Statistical Independence Testing

whether $\mathbb{P}_{XY} = \mathbb{P}_X \mathbb{P}_Y$?

Hypothesis testing

- Null hypothesis \mathcal{H}_0 : $\mathbb{P}_{XY} = \mathbb{P}_X \mathbb{P}_Y$.
- Alternative hypothesis \mathcal{H}_1 : $\mathbb{P}_{XY} \neq \mathbb{P}_X \mathbb{P}_Y$.



Given n i.i.d samples $Z := \{(x_i, y_i)\}_{i=1}^n$

Kernel-based Statistical Independence Testing

Definition (Hilbert-Schmidt Independence Criterion)

Let \mathcal{F} be an RKHS with kernel $k : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$ and let \mathcal{G} be a second RKHS on \mathcal{Y} with kernel $l : \mathcal{Y} \times \mathcal{Y} \mapsto \mathbb{R}$, the HSIC between X and Y , denoted as $\text{HSIC}(X, Y)$ is defined as

$$\mathbf{E}[k(X, X')l(Y, Y')] + \mathbf{E}[k(X, X')] \mathbf{E}[l(Y, Y')] - 2\mathbf{E}_{X', Y'}[\mathbf{E}_X k(X, X') \mathbf{E}_Y l(Y, Y')],$$

where (X', Y') is a independent copy of (X, Y) .

Kernel-based Statistical Independence Testing

Definition (Hilbert-Schmidt Independence Criterion)

Let \mathcal{F} be an RKHS with kernel $k : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$ and let \mathcal{G} be a second RKHS on \mathcal{Y} with kernel $l : \mathcal{Y} \times \mathcal{Y} \mapsto \mathbb{R}$, the HSIC between X and Y , denoted as $\text{HSIC}(X, Y)$ is defined as

$$\mathbf{E}[k(X, X')l(Y, Y')] + \mathbf{E}[k(X, X')] \mathbf{E}[l(Y, Y')] - 2\mathbf{E}_{X', Y'}[\mathbf{E}_X k(X, X') \mathbf{E}_Y l(Y, Y')],$$

where (X', Y') is a independent copy of (X, Y) .

- $\text{HSIC}(X, Y) = 0 \Leftrightarrow \mathbb{P}_{XY} = \mathbb{P}_X \mathbb{P}_Y$ with suitable kernels (e.g. Gaussian kernel).

Examples (Gaussian kernel with width σ)

The Gaussian kernel is defined as $k(x, x') := \exp\left(-\frac{\|x-x'\|^2}{2\sigma^2}\right)$, where σ is the width.

Kernel-based Statistical Independence Testing

Definition (Estimation of HSIC)

An estimator of $\text{HSIC}(X, Y)$ with sample Z is given by

$$\text{HSIC}_b(Z) := \frac{1}{n^2} \sum_{i,j} k_{ij} l_{ij} + \frac{1}{n^4} \sum_{i,j,q,r} k_{ij} l_{qr} - 2 \frac{1}{n^3} \sum_{i,j,q} k_{ij} l_{iq} = \frac{1}{n^2} \text{Tr}(\mathbf{KHLH}),$$

where $k_{ij} := k(x_i, x_j)$, $l_{ij} := l(y_i, y_j)$ are the entries of the $n \times n$ kernel matrix \mathbf{K} , \mathbf{L} , respectively and $\mathbf{H} = \mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}^T$ is the center matrix and $\mathbf{1}$ is a vector of ones.

$$\mathbf{K}_{n \times n} = \begin{bmatrix} k_{11} & k_{12} & \cdots & k_{1n} \\ k_{21} & k_{22} & \cdots & k_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ k_{n1} & k_{n2} & \cdots & k_{nn} \end{bmatrix}$$

Calculate $\text{HSIC}_b(Z)$ cost
 $\mathcal{O}(n^2)$ time and $\mathcal{O}(n^2)$ space.
n: sample size.

What's More?

$$\text{HSIC}_b(Z) := \frac{1}{n^2} \sum_{i,j} k_{ij} l_{ij} + \frac{1}{n^4} \sum_{i,j,q,r} k_{ij} l_{qr} - 2 \frac{1}{n^3} \sum_{i,j,q} k_{ij} l_{iq} = \frac{1}{n^2} \text{Tr}(\mathbf{KHLH}),$$

- **Faster:** The time/space computational complexity of current statistic are both quadratic computing time.
- **More flexible:** The kernel can not be adaptive.

What's More?

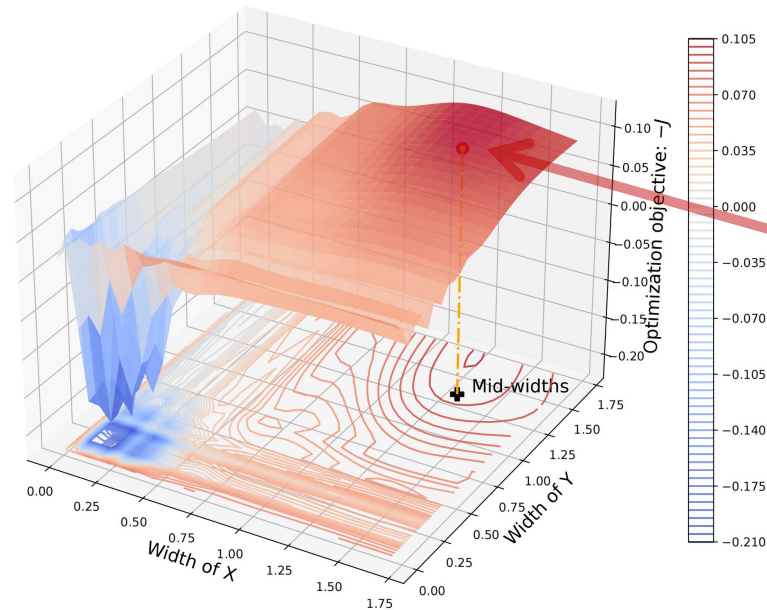
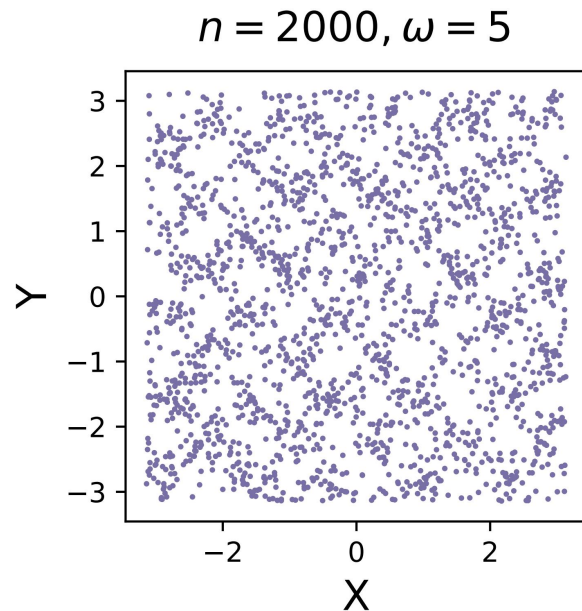
$$\text{HSIC}_b(Z) := \frac{1}{n^2} \sum_{i,j} k_{ij} l_{ij} + \frac{1}{n^4} \sum_{i,j,q,r} k_{ij} l_{qr} - 2 \frac{1}{n^3} \sum_{i,j,q} k_{ij} l_{iq} = \frac{1}{n^2} \text{Tr}(\mathbf{KHLH}),$$

- **Faster:** The time/space computational complexity of current statistic are both quadratic computing time.
- **More flexible:** The kernel can not be adaptive.

A Frequency Domain Perspective

$$\text{HSIC}(X, Y) = \int_{\mathbb{R}^{d_x} \times \mathbb{R}^{d_y}} |\phi_{\mathbb{P}_{XY}}(\omega) - \phi_{\mathbb{P}_X \mathbb{P}_Y}(\omega)|^2 (\mathcal{F}^{-1}\psi)(\omega) d\omega,$$

Example: $(X, Y) \sim p_{xy}(x, y) \propto 1 + \sin(\omega_0 x) \sin(\omega_0 y)$.



**Inflexible settings
result in
low test power**

We aim to obtain a **more flexible** $(\mathcal{F}^{-1}\psi)(\omega)$

A Frequency Domain Perspective

Enable more efficient calculation of statistic

- The statistic

$$\text{HSIC}(X, Y) = \int_{\mathbb{R}^{d_x} \times \mathbb{R}^{d_y}} |\phi_{\mathbb{P}_{XY}}(\omega) - \phi_{\mathbb{P}_X \mathbb{P}_Y}(\omega)|^2 (\mathcal{F}^{-1}\psi)(\omega) d\omega.$$

- Frequency samplings for integral approximation

$$\text{HSIC}_\omega(X, Y) := \frac{1}{D_x D_y} \sum_{i=1}^{D_x} \sum_{j=1}^{D_y} |\phi_{\mathbb{P}_{XY}}(\omega_{x;i}, \omega_{y;j}) - \phi_{\mathbb{P}_X \mathbb{P}_Y}(\omega_{x;i}, \omega_{y;j})|^2,$$

Sampling



where $\{\omega_{x;i}\}_{i=1}^{D_x}, \{\omega_{y;j}\}_{j=1}^{D_y}$ are sampled independently with the measure $\mathcal{F}^{-1}\psi_k, \mathcal{F}^{-1}\psi_l$, respectively. And $\mathcal{F}^{-1}\psi$ is a product measure, i.e., $\mathcal{F}^{-1}\psi = (\mathcal{F}^{-1}\psi_k) \otimes (\mathcal{F}^{-1}\psi_l)$.

This type of approximation also called random Fourier features (RFF).

Obtain Independence Criterion

Design $\mathcal{F}^{-1}\psi$ (take $\mathcal{F}^{-1}\psi_k(\omega)$ for example)

- Designing by kernels with adjustable parameters.

Kernel	$\psi_k(\Delta)$	$\mathcal{F}^{-1}\psi_k(\omega)$	$\mathcal{T}_{\theta_k}(x)$	$p_k(\omega)$
Gaussian	$e^{-\frac{\ \Delta\ _2^2}{2\sigma^2}}$	$(2\pi)^{-d_x/2}\sigma e^{-\sigma^2\ \omega\ _2^2/2}$	x/σ	$(2\pi)^{-d_x/2}e^{-\ \omega\ _2^2/2}$
Laplace	$e^{-\frac{\ \Delta\ _1}{\sigma}}$	$\sqrt{\frac{2}{\pi}} \prod_d \frac{\sigma}{\sigma^2 + \omega_d^2}$	x/σ	$\sqrt{\frac{2}{\pi}} \prod_d \frac{1}{1 + \omega_d^2}$
Mahalanobis	$e^{-\frac{1}{2}\Delta^T \Sigma^{-1} \Delta}$	$(2\pi)^{-d_x/2} \Sigma ^{-1/2} e^{-\omega^T \Sigma^{-1} \omega/2}$	$\Sigma^{1/2} x$	$(2\pi)^{-d_x/2} e^{-\ \omega\ _2^2/2}$

Table: Some popular kernels (parametered σ, Σ) with corresponding density functions.

Obtain Independence Criterion

Design $\mathcal{F}^{-1}\psi$ (take $\mathcal{F}^{-1}\psi_k(\omega)$ for example)

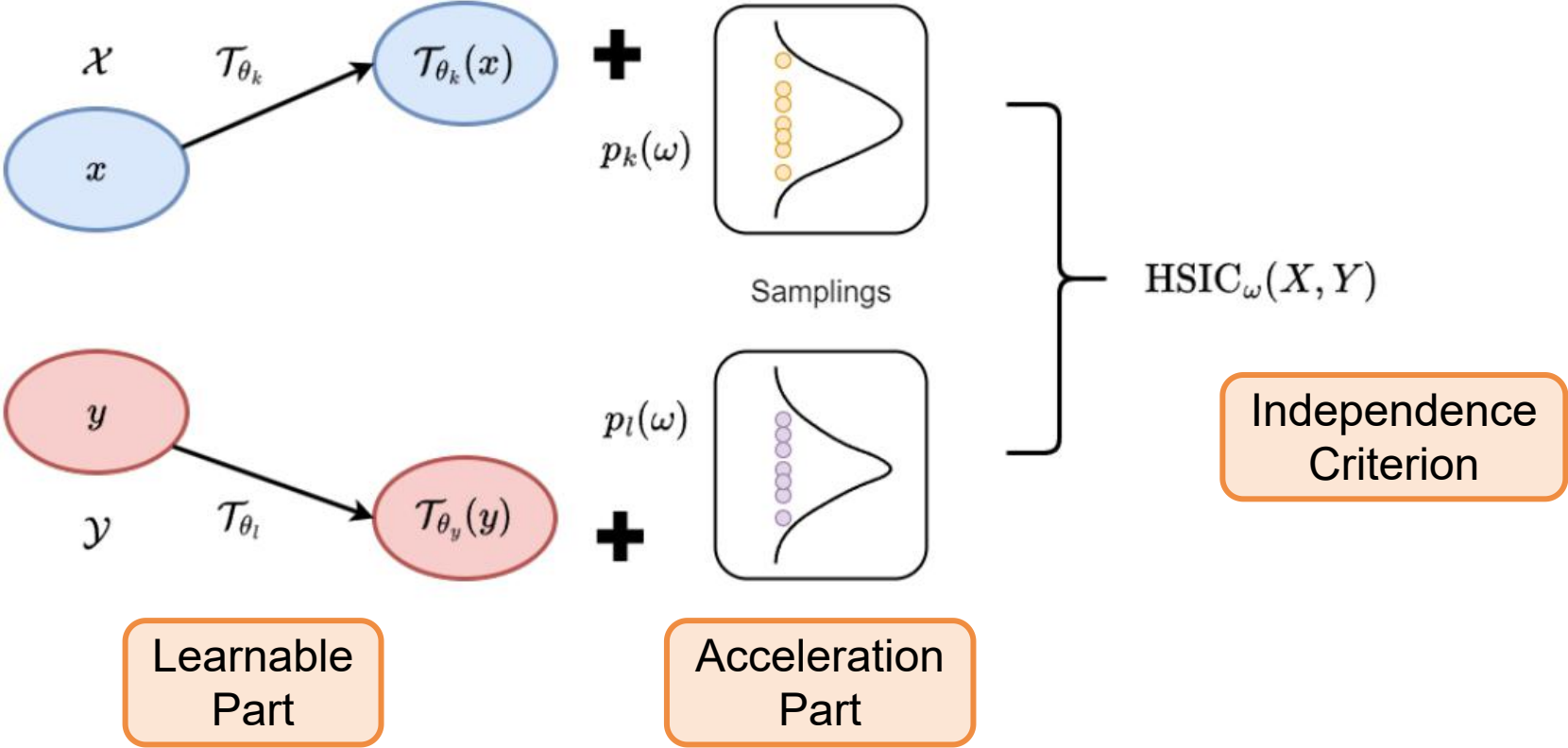
- Designing by kernels with adjustable parameters.

Kernel	$\psi_k(\Delta)$	$\mathcal{F}^{-1}\psi_k(\omega)$	$\mathcal{T}_{\theta_k}(x)$	$p_k(\omega)$
Gaussian	$e^{-\frac{\ \Delta\ _2^2}{2\sigma^2}}$	$(2\pi)^{-d_x/2} \sigma e^{-\sigma^2 \ \omega\ _2^2/2}$	x/σ	$(2\pi)^{-d_x/2} e^{-\ \omega\ _2^2/2}$
Laplace	$e^{-\frac{\ \Delta\ _1}{\sigma}}$	$\sqrt{\frac{2}{\pi}} \prod_d \frac{\sigma}{\sigma^2 + \omega_d^2}$	x/σ	$\sqrt{\frac{2}{\pi}} \prod_d \frac{1}{1 + \omega_d^2}$
Mahalanobis	$e^{-\frac{1}{2} \Delta^T \Sigma^{-1} \Delta}$	$(2\pi)^{-d_x/2} \Sigma ^{-1/2} e^{-\omega^T \Sigma^{-1} \omega/2}$	$\Sigma^{1/2} x$	$(2\pi)^{-d_x/2} e^{-\ \omega\ _2^2/2}$

Table: Some popular kernels (parametered σ, Σ) with corresponding density functions.

- Disentangling the sampled objects and the learnable parameters
 - ▶ Relocate the learnable component onto X using $\mathcal{T}_{\theta_k}(x)$.
 - ▶ Convert the probability measure $\mathcal{F}^{-1}\psi_k$ into a standard distribution $p_k(\omega)$.

Obtain Independence Criterion



Learnable RFF

We obtain the approximation of kernel

- Kernel with learnable mappings

$$\psi_k(\mathcal{T}_{\theta_k}x - \mathcal{T}_{\theta_k}x') = \mathcal{F}[\mathcal{F}^{-1}\psi_k(\omega)] = \int e^{-i\omega^T(\mathcal{T}_{\theta_k}x - \mathcal{T}_{\theta_k}x')} p_k(\omega) d\omega.$$

- Use the frequency sampling technique

$$\psi_k^{(\omega)}(\mathcal{T}_{\theta_k}x - \mathcal{T}_{\theta_k}x') := \frac{2}{D} \sum_{j=1}^{D/2} e^{-i\omega_{k;j}^T(\mathcal{T}_{\theta_k}x - \mathcal{T}_{\theta_k}x')} = \frac{2}{D} \sum_{j=1}^{D/2} \cos\left(\omega_{k;j}^T(\mathcal{T}_{\theta_k}x - \mathcal{T}_{\theta_k}x')\right),$$

where $\{\omega_{k;j}\}_{j=1}^{D/2}$ are sampled independently with distribution $p_k(\omega)$.

- The learnable RFF of k

$$\Lambda_k(x) := \sqrt{\frac{2}{D}} \left[\cos(\omega_1^T \mathcal{T}_{\theta_k}x), \sin(\omega_1^T \mathcal{T}_{\theta_k}x), \dots, \cos(\omega_{D/2}^T \mathcal{T}_{\theta_k}x), \sin(\omega_{D/2}^T \mathcal{T}_{\theta_k}x) \right],$$

- Hence $\psi_k^{(\omega)}(\mathcal{T}_{\theta_k}x - \mathcal{T}_{\theta_k}x') = \Lambda_k(x)\Lambda_k(x')^T$.

Independence Criterion

The statistic can be obtained as follows.

- The learnable RFF of samples in matrix form $\mathbf{\Lambda}_X := [\Lambda_k(x_1); \dots; \Lambda_k(x_n)]_{n \times D}$.
- Obtain $\mathbf{\Lambda}_Y$ by analogy. The same number of samplings are used for simplify.
- The statistic with sample Z

$$\text{HSIC}_\omega(Z) := \frac{1}{n^2} \text{Tr}(\mathbf{\Lambda}_X \mathbf{\Lambda}_X^T \mathbf{H} \mathbf{\Lambda}_Y \mathbf{\Lambda}_Y^T \mathbf{H}) = \frac{1}{n^2} \text{Tr}(\mathbf{\Lambda}_X^T \mathbf{H} \mathbf{\Lambda}_Y \mathbf{\Lambda}_Y^T \mathbf{H} \mathbf{\Lambda}_X) = \frac{1}{n^2} \|\mathbf{\Lambda}_{Xc}^T \mathbf{\Lambda}_{Yc}\|_F^2,$$

where $\mathbf{\Lambda}_{Xc} := \mathbf{H} \mathbf{\Lambda}_X$, $\mathbf{\Lambda}_{Yc} := \mathbf{H} \mathbf{\Lambda}_Y$.

- The time complexity is $\mathcal{O}(nD(d_x + d_y + D))$, i.e. the running time is linear with n .

Asymptotic Behavior

Proposition (Asymptotics)

Let $h_{ijqr}^{(\omega)} := \frac{1}{4!} \sum_{(t,u,v,w)}^{(i,j,q,r)} k_{tu}^{(\omega)} l_{tu}^{(\omega)} + k_{tu}^{(\omega)} l_{vw}^{(\omega)} - 2k_{uv}^{(\omega)} l_{tv}^{(\omega)}$. Then, Under the null hypothesis \mathcal{H}_0 , $HSIC_{\omega}(Z)$ coverages in distribution to

$$nHSIC_{\omega}(Z) \xrightarrow{d} \sum_{l=1}^{\infty} \lambda_l \chi_{1l}^2, \quad \lambda_l g_l(z_j) = \int_{z_i, z_q, z_r} h_{ijqr}^{(\omega)} g_l(z_i) dF_{z_i, z_q, z_r},$$

where $\chi_{11}^2, \chi_{12}^2, \dots$ are independent χ_1^2 variates and λ_l is the solution to eigenvalue problem. Also, under \mathcal{H}_1 , $HSIC_{\omega}(Z)$ converges in distribution as

$$n^{\frac{1}{2}} \left(HSIC_{\omega}(Z) - \mathbf{E}_Z HSIC_{\omega}(Z) \right) \xrightarrow{d} \mathcal{N}(0, \sigma_{\omega}^2), \quad \sigma_{\omega}^2 := 16 \left[\mathbf{E}_i \left(\mathbf{E}_{j,q,r} h_{ijqr}^{(\omega)} \right)^2 - \left(\mathbf{E}_Z h_{ijqr}^{(\omega)} \right)^2 \right]$$

with the simplified notation $\mathbf{E}_{j,q,r} := \mathbf{E}_{z_j, z_q, z_r}$ and $\mathbf{E}_Z := \mathbf{E}_{z_i, z_j, z_q, z_r}$.

Construct Optimization Objective

- According to the asymptotics, the power of the test

$$\mathbb{P}_{\mathcal{H}_1} (n\text{HSIC}_\omega(Z) > r_\omega) \rightarrow \Phi \left(\frac{n\mathbf{E}_Z \text{HSIC}_\omega(Z) - r_\omega}{\sqrt{n}\sigma_\omega} \right),$$

where Φ is the standard normal CDF.

- The optimization objective

$$J := \frac{(n\text{HSIC}_\omega(Z) - \hat{c}_\alpha)}{(\sqrt{n}\hat{\sigma}_\omega)},$$

where $\text{HSIC}_\omega(Z)$ is the criterion, \hat{c}_α is a estimate of threshold and the estimate of variance $\hat{\sigma}_\omega^2 := 16 \left[\frac{1}{n} \sum_i \left(\frac{1}{n^3} \sum_{j,q,r} h_{ijqr}^{(\omega)} \right)^2 - \text{HSIC}_\omega^2(Z) \right]$.

Obtain \widehat{c}_α with Gamma approximation.

- Determined completely by the first two moments of distribution under \mathcal{H}_0 .

$$n\text{HSIC}_\omega(Z) \sim \frac{x^{\gamma-1} e^{-x/\beta}}{\beta^\gamma \Gamma(\gamma)}, \text{ where } \gamma = \frac{(\mathbf{E}[\text{HSIC}_\omega(Z)])^2}{\text{var}[\text{HSIC}_\omega(Z)]}, \beta = \frac{n\mathbf{Var}[\text{HSIC}_\omega(Z)]}{\mathbf{E}[\text{HSIC}_\omega(Z)]}$$

- The $(1 - \alpha)$ -quantile of Gamma distribution

$$\int_0^{\widehat{c}_\alpha} \frac{x^{\gamma-1} e^{-x/\beta}}{\beta^\gamma \Gamma(\gamma)} dx = 1 - \alpha.$$

Obtain \widehat{c}_α with Gamma approximation.

- Determined completely by the first two moments of distribution under \mathcal{H}_0 .

$$n\text{HSIC}_\omega(Z) \sim \frac{x^{\gamma-1} e^{-x/\beta}}{\beta^\gamma \Gamma(\gamma)}, \text{ where } \gamma = \frac{(\mathbf{E}[\text{HSIC}_\omega(Z)])^2}{\text{var}[\text{HSIC}_\omega(Z)]}, \beta = \frac{n\mathbf{Var}[\text{HSIC}_\omega(Z)]}{\mathbf{E}[\text{HSIC}_\omega(Z)]}$$

- The $(1 - \alpha)$ -quantile of Gamma distribution

$$\int_0^{\widehat{c}_\alpha} \frac{x^{\gamma-1} e^{-x/\beta}}{\beta^\gamma \Gamma(\gamma)} dx = 1 - \alpha.$$

The linear-time estimations of two moments of distribution under \mathcal{H}_0 .

Theorem (Linear-Time Estimations)

Under \mathcal{H}_0 , the estimation of mean and variance with bias of $\mathcal{O}(n^{-1})$ to $\mathbf{E}_Z[n\text{HSIC}_\omega(Z)]$ and $\mathbf{Var}_Z[n\text{HSIC}_\omega(Z)]$, denote as \mathcal{E}_0 and \mathcal{V}_0 , respectively, are given by

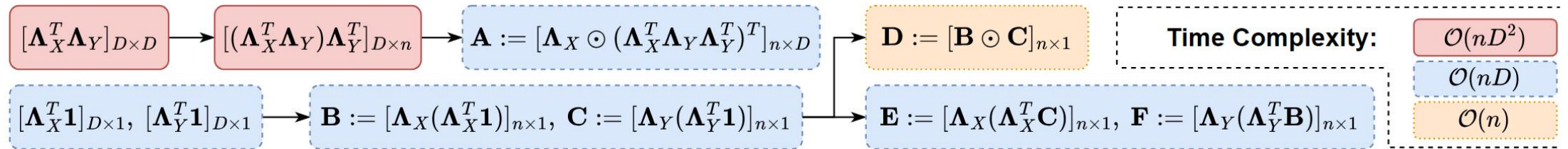
$$\mathcal{E}_0 := \frac{[\mathbf{1}^T \mathbf{\Lambda}_{X_c}^2 \mathbf{1}][\mathbf{1}^T \mathbf{\Lambda}_{Y_c}^2 \mathbf{1}]}{(n-1)^2}, \mathcal{V}_0 := \frac{2n(n-4)(n-5)}{(n-1)(n-2)(n-3)} \frac{[\mathbf{1}^T (\mathbf{\Lambda}_{X_c}^T \mathbf{\Lambda}_{X_c})^2 \mathbf{1}][\mathbf{1}^T (\mathbf{\Lambda}_{Y_c}^T \mathbf{\Lambda}_{Y_c})^2 \mathbf{1}]}{n^4},$$

where $()^2$ is the entrywise matrix power. Both \mathcal{E}_0 and \mathcal{V}_0 can be calculated in $\mathcal{O}(nD^2)$ time.

The estimate of variance $\hat{\sigma}_\omega^2$ also can be calculated in linear-time

- By the definition $\hat{\sigma}_\omega^2 := 16 \left[\frac{1}{n} \sum_i \left(\frac{1}{n^3} \sum_{j,q,r} h_{ijqr}^{(\omega)} \right)^2 - \text{HSIC}_\omega^2(Z) \right]$.
- Calculate $\sum_{j,q,r} h_{ijqr}^{(\omega)}$ in linear time

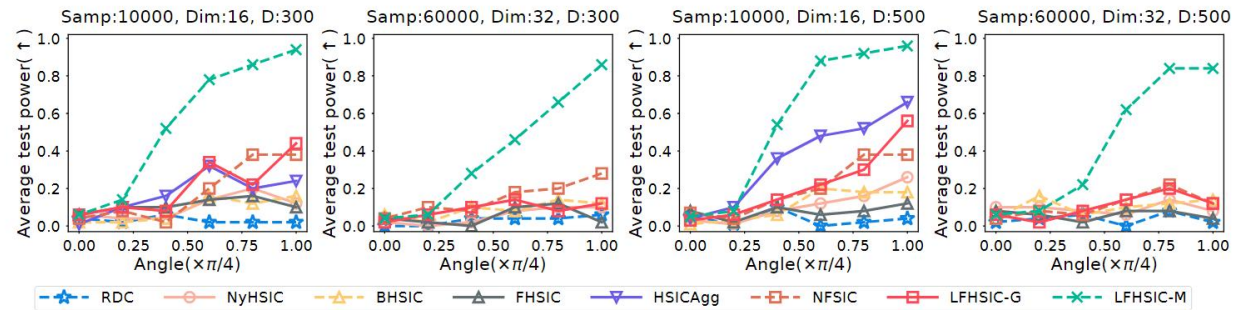
$$\sum_{j,q,r} h_{ijqr}^{(\omega)} = \frac{1}{2} \left[n \mathbf{1}^T \mathbf{A} \mathbf{1} + n^2 (\mathbf{A} \mathbf{1})_i + (\mathbf{1}^T \mathbf{C}) \mathbf{B}_i + (\mathbf{1}^T \mathbf{B}) \mathbf{C}_i - n \mathbf{E}_i - n \mathbf{F}_i - n \mathbf{D}_i - \mathbf{1}^T \mathbf{D} \right],$$



As a result, the optimization objective can be computed in linear time.

Experiments

Results on ISA Datasets (Large Scale)



Results on 3DShape (High-dimensional) and Real Datasets

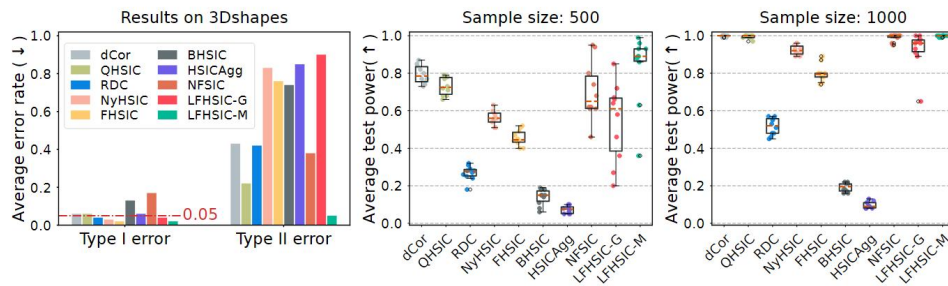
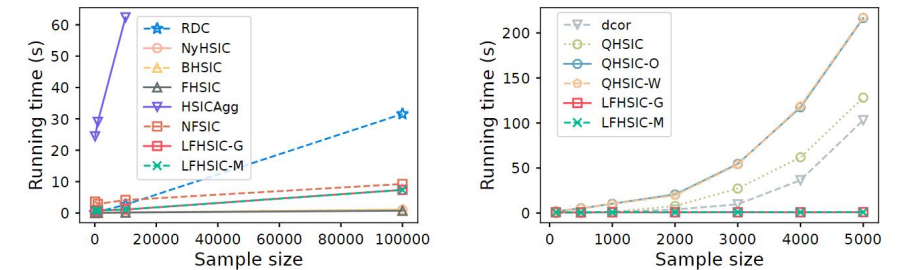


Figure: The results on two real data. Left: 3DShapes. Right two: MSD Dataset.

Running time comparing



Our test can handle 100,000 samples within 10 seconds.

- **Flexible:** handle large scale/high-dimensional settings well.
- **Fast:** linear-time time/space computation complexity.

Conclusions

The method achieves flexible independence testing in linear time (w.r.t sample size).

1. Flexible: The kernel can be *adaptive*.
2. Fast: *linear-time time/space* complexity.

Thank you for you attention!