

Structured Data: Dependency, Testing

Zoltán Szabó

∈ Structured Data: Learning, Prediction, **Dependency, Testing**
M2 Data Science, University of Paris-Saclay

Paris & Palaiseau, France
Feb. 25, March 4, 18, 25, 2019

Contact information

- Email:

zoltan (dot) szabo (at) polytechnique (dot) edu

- Web:

<http://www.cmap.polytechnique.fr/~zoltan.szabo/>

Software (Python, Matlab)

- Dependency measures (KCCA, HSIC), divergences (MMD), etc.; several demos:

<https://bitbucket.org/szzoli/ite-in-python>
<https://bitbucket.org/szzoli/ite/>

- 2-sample, independence & goodness-of-fit tests (quadratic → linear-time methods):

<https://github.com/wittawatj/interpretable-test>
<https://github.com/wittawatj/fsic-test>
<https://github.com/wittawatj/kernel-gof>

Outline

- Motivation:
 - Objective functions: from dependency measures.
 - Testing.

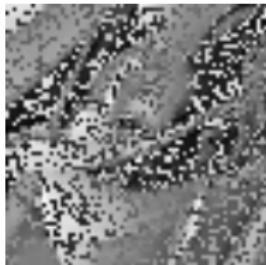
Outline

- Motivation:
 - Objective functions: from dependency measures.
 - Testing.
- Kernel, RKHS.
- Kernel canonical correlation analysis.
- Mean embedding:
 - Characteristic property,
 - Universality.
- Maximum mean discrepancy.
- Cross-covariance operator, HSIC.
- Hypothesis testing.

Dependency Measures as Objective Functions

Outlier-robust image registration [Kybic, 2004, Neemuchwala et al., 2007]

Given two images:

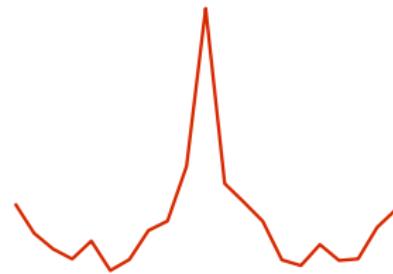
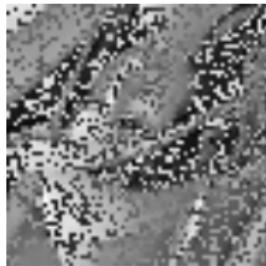


Goal: find the transformation which takes the right one to the left.

Outlier-robust image registration

[Kybic, 2004, Neemuchwala et al., 2007]

Given two images:



Goal: find the transformation which takes the right one to the left.

Outlier-robust image registration: equations

- Reference image: \mathbf{y}_{ref} ,
- test image: \mathbf{y}_{test} ,
- possible transformations: Θ .

Objective:

$$J(\theta) = \underbrace{I(\mathbf{y}_{\text{ref}}, \mathbf{y}_{\text{test}}(\theta))}_{\text{similarity}} \rightarrow \max_{\theta \in \Theta}$$

In the example: $I=KCCA$.

Independent Subspace Analysis [Cardoso, 1998]

Cocktail party problem:

- independent groups of people / music bands,
- observation = mixed sources.



ISA equations

Observation:

$$\mathbf{x}_t = \mathbf{A}\mathbf{s}_t, \quad \mathbf{s} = [\mathbf{s}^1; \dots; \mathbf{s}^M].$$

Goal: $\hat{\mathbf{s}}$ from $\{\mathbf{x}_1, \dots, \mathbf{x}_T\}$. Assumptions:

- independent groups: $I(\mathbf{s}^1, \dots, \mathbf{s}^M) = 0$,
- \mathbf{s}^m -s: non-Gaussian,
- \mathbf{A} : invertible.

Find \mathbf{W} which makes the estimated components independent:

$$\mathbf{y} = \mathbf{Wx} = \left[\mathbf{y}^1; \dots; \mathbf{y}^M \right],$$
$$J(\mathbf{W}) = I\left(\mathbf{y}^1, \dots, \mathbf{y}^M\right) \rightarrow \min_{\mathbf{W}}.$$

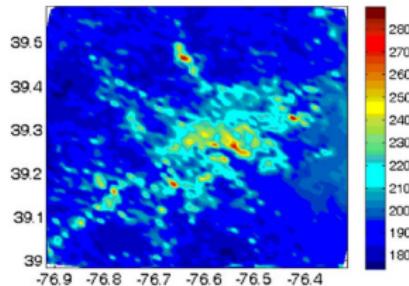
Distribution regression

[Póczos et al., 2013, Szabó et al., 2016]. Sustainability

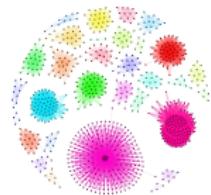
- **Goal:** aerosol prediction = air pollution → climate.



- Prediction using labelled bags:
 - bag := multi-spectral satellite measurements over an area,
 - label := local aerosol value.



Objects in the bags

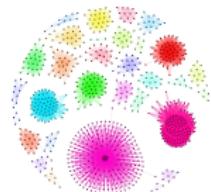


- Examples:
 - time-series modelling: user = set of **time-series**,
 - computer vision: image = collection of patch **vectors**,
 - NLP: corpus = bag of **documents**,
 - network analysis: group of people = bag of friendship **graphs**, ...

Objects in the bags



time series



- Examples:
 - time-series modelling: user = set of **time-series**,
 - computer vision: image = collection of patch **vectors**,
 - NLP: corpus = bag of **documents**,
 - network analysis: group of people = bag of friendship **graphs**, ...
- Wider context (statistics): point estimation tasks.

Regression on labelled bags

- Given:
 - labelled bags: $\hat{\mathbf{z}} = \{(\hat{P}_i, \mathbf{y}_i)\}_{i=1}^{\ell}$, \hat{P}_i : bag from P_i , $N := |\hat{P}_i|$.
 - test bag: \hat{P} .

Regression on labelled bags

- Given:
 - labelled bags: $\hat{\mathbf{z}} = \{(\hat{P}_i, y_i)\}_{i=1}^\ell$, \hat{P}_i : bag from P_i , $N := |\hat{P}_i|$.
 - test bag: \hat{P} .
- Estimator:

$$f_{\hat{\mathbf{z}}}^{\lambda} = \arg \min_{f \in \mathcal{H}} \frac{1}{\ell} \sum_{i=1}^{\ell} \left[f(\underbrace{\mu_{\hat{P}_i}}_{\text{feature of } \hat{P}_i}) - y_i \right]^2 + \lambda \|f\|_{\mathcal{H}}^2.$$

Regression on labelled bags

- Given:
 - labelled bags: $\hat{\mathbf{z}} = \{(\hat{P}_i, y_i)\}_{i=1}^\ell$, \hat{P}_i : bag from P_i , $N := |\hat{P}_i|$.
 - test bag: \hat{P} .
- Estimator:

$$f_{\hat{\mathbf{z}}}^\lambda = \arg \min_{f \in \mathcal{H}_K} \frac{1}{\ell} \sum_{i=1}^{\ell} \left[f(\mu_{\hat{P}_i}) - y_i \right]^2 + \lambda \|f\|_{\mathcal{H}}^2.$$

- Prediction:

$$\begin{aligned}\hat{y}(\hat{P}) &= \mathbf{g}^T (\mathbf{G} + \ell \lambda \mathbf{I})^{-1} \mathbf{y}, \\ \mathbf{g} &= [K(\mu_{\hat{P}}, \mu_{\hat{P}_i})], \mathbf{G} = [K(\mu_{\hat{P}_i}, \mu_{\hat{P}_j})], \mathbf{y} = [y_i].\end{aligned}$$

Regression on labelled bags

- Given:

- labelled bags: $\hat{\mathbf{z}} = \{(\hat{P}_i, y_i)\}_{i=1}^\ell$, \hat{P}_i : bag from P_i , $N := |\hat{P}_i|$.
- test bag: \hat{P} .

- Estimator:

$$f_{\hat{\mathbf{z}}}^{\lambda} = \arg \min_{f \in \mathcal{H}_K} \frac{1}{\ell} \sum_{i=1}^{\ell} \left[f(\mu_{\hat{P}_i}) - y_i \right]^2 + \lambda \|f\|_{\mathcal{H}}^2.$$

- Prediction:

$$\begin{aligned}\hat{y}(\hat{P}) &= \mathbf{g}^T (\mathbf{G} + \ell \lambda \mathbf{I})^{-1} \mathbf{y}, \\ \mathbf{g} &= [K(\mu_{\hat{P}}, \mu_{\hat{P}_i})], \mathbf{G} = [K(\mu_{\hat{P}_i}, \mu_{\hat{P}_j})], \mathbf{y} = [y_i].\end{aligned}$$

Challenge

Inner product of distributions: $K(\mu_{\hat{P}_i}, \mu_{\hat{P}_j}) = ?$

Feature selection

- **Goal:** find
 - the feature subset (# of rooms, criminal rate, local taxes)
 - most relevant for house price prediction (y).



Feature selection: equations

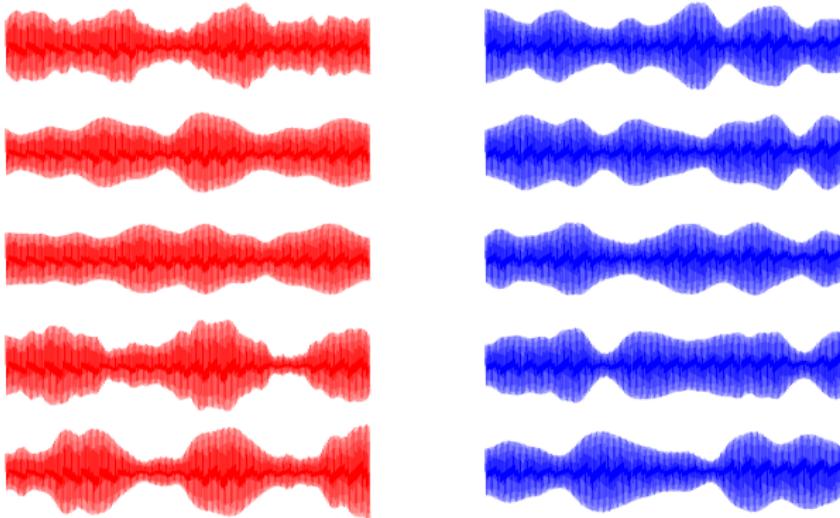
- Features: x^1, \dots, x^F . Subset: $S \subseteq \{1, \dots, F\}$.
- MaxRelevance - MinRedundancy principle [Peng et al., 2005]:

$$J(S) = \frac{1}{|S|} \sum_{i \in S} I(x^i, y) - \frac{1}{|S|^2} \sum_{i, j \in S} I(x^i, x^j) \rightarrow \max_{S \subseteq \{1, \dots, F\}} .$$

Testing

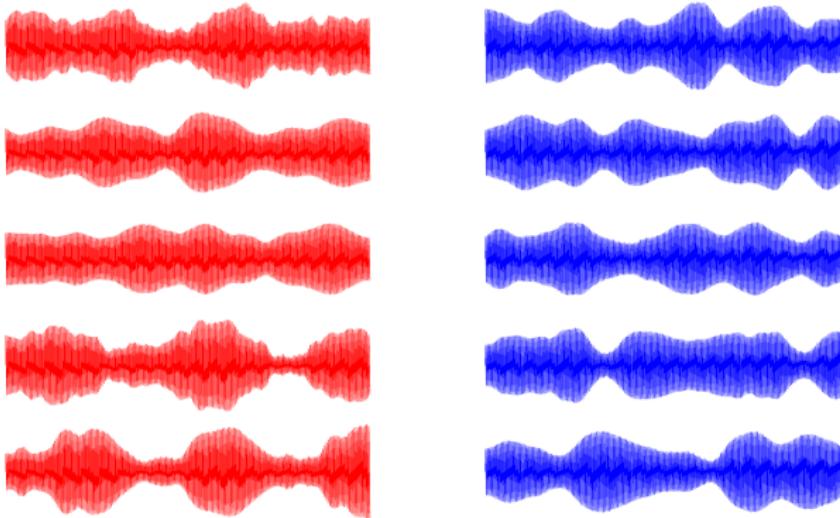
Motivation: detecting differences in AM signals

- Amplitude modulation:
 - simple technique to transmit voice over radio.
 - in the example: 2 songs.
- Fragments from song₁ ~ $\textcolor{red}{P}_x$, song₂ ~ $\textcolor{blue}{P}_y$.



Motivation: detecting differences in AM signals

- Amplitude modulation:
 - simple technique to transmit voice over radio.
 - in the example: 2 songs.
- Fragments from song₁ ~ \mathbb{P}_x , song₂ ~ \mathbb{P}_y .



Question: $\mathbb{P}_x = \mathbb{P}_y$?

Motivation: discrete domain - 2-sample testing

- How do we compare distributions?
- Given: 2 sets of text fragments (**fisheries, agriculture**).

x_1 : Now disturbing reports out of Newfoundland show that the fragile snow crab industry is in serious decline. First the west coast salmon, the east coast salmon and the cod, and now the snow crabs off Newfoundland.

x_2 : To my pleasant surprise he responded that he had personally visited those wharves and that he had already announced money to fix them. What wharves did the minister visit in my riding and how much additional funding is he going to provide for Delaps Cove, Hampton, Port Lorne, ...

...

y_1 : Honourable senators, I have a question for the Leader of the Government in the Senate with regard to the support funding to farmers that has been announced. Most farmers have not received any money yet.

y_2 : On the grain transportation system we have had the Estey report and the Kroeger report. We could go on and on. Recently programs have been announced over and over by the government such as money for the disaster in agriculture on the prairies and across Canada.

...

Motivation: discrete domain - 2-sample testing

- How do we compare distributions?
- Given: 2 sets of text fragments (**fisheries, agriculture**).

x_1 : Now disturbing reports out of Newfoundland show that the fragile snow crab industry is in serious decline. First the west coast salmon, the east coast salmon and the cod, and now the snow crabs off Newfoundland.

x_2 : To my pleasant surprise he responded that he had personally visited those wharves and that he had already announced money to fix them. What wharves did the minister visit in my riding and how much additional funding is he going to provide for Delaps Cove, Hampton, Port Lorne, ...

...

y_1 : Honourable senators, I have a question for the Leader of the Government in the Senate with regard to the support funding to farmers that has been announced. Most farmers have not received any money yet.

y_2 : On the grain transportation system we have had the Estey report and the Kroeger report. We could go on and on. Recently programs have been announced over and over by the government such as money for the disaster in agriculture on the prairies and across Canada.

...

Do $\{x_i\}$ and $\{y_j\}$ come from the same distribution, i.e. $\mathbb{P}_x = \mathbb{P}_y$?

Motivation: discrete domain - independence testing

- How do we detect dependency? (paired samples)

x_1 : Honourable senators, I have a question for the Leader of the Government in the Senate with regard to the support funding to farmers that has been announced. Most farmers have not received any money yet.

x_2 : No doubt there is great pressure on provincial and municipal governments in relation to the issue of child care, but the reality is that there have been no cuts to child care funding from the federal government to the provinces. In fact, we have increased federal investments for early childhood development.

...

y_1 : Honorables sénateurs, ma question s'adresse au leader du gouvernement au Sénat et concerne l'aide financière qu'on a annoncée pour les agriculteurs. La plupart des agriculteurs n'ont encore rien reçu de cet argent.

y_2 : Il est évident que les ordres de gouvernements provinciaux et municipaux subissent de fortes pressions en ce qui concerne les services de garde, mais le gouvernement n'a pas réduit le financement qu'il verse aux provinces pour les services de garde. Au contraire, nous avons augmenté le financement fédéral pour le développement des jeunes enfants.

...

Motivation: discrete domain - independence testing

- How do we detect dependency? (paired samples)

x_1 : Honourable senators, I have a question for the Leader of the Government in the Senate with regard to the support funding to farmers that has been announced. Most farmers have not received any money yet.

x_2 : No doubt there is great pressure on provincial and municipal governments in relation to the issue of child care, but the reality is that there have been no cuts to child care funding from the federal government to the provinces. In fact, we have increased federal investments for early childhood development.

...

y_1 : Honorables sénateurs, ma question s'adresse au leader du gouvernement au Sénat et concerne l'aide financière qu'on a annoncée pour les agriculteurs. La plupart des agriculteurs n'ont encore rien reçu de cet argent.

y_2 : Il est évident que les ordres de gouvernements provinciaux et municipaux subissent de fortes pressions en ce qui concerne les services de garde, mais le gouvernement n'a pas réduit le financement qu'il verse aux provinces pour les services de garde. Au contraire, nous avons augmenté le financement fédéral pour le développement des jeunes enfants.

...

Are the French paragraphs translations of the English ones, or have nothing to do with it, i.e. $\mathbb{P}_{XY} = \mathbb{P}_X \otimes \mathbb{P}_Y$?

We will use **kernels** to tackle these problems

They exist essentially **on any data type**

We will use **kernels** to tackle these problems

They exist essentially **on any data type**

trees [Collins and Duffy, 2001, Kashima and Koyanagi, 2002], time series [Cuturi, 2011], strings [Lodhi et al., 2002], mixture models, hidden Markov models or linear dynamical systems [Jebara et al., 2004], sets [Haussler, 1999, Gärtner et al., 2002], fuzzy domains [Guevara et al., 2017], distributions [Hein and Bousquet, 2005, Martins et al., 2009, Muandet et al., 2011], groups [Cuturi et al., 2005] with specific constructions on permutations [Jiao and Vert, 2016], graphs [Vishwanathan et al., 2010, Kondor and Pan, 2016], ...



Kernel Canonical Correlation Analysis (KCCA)

Independence measures

- Given: random variable $(x, y) \in \mathcal{X} \times \mathcal{Y}$, $(x, y) \sim \mathbb{P}_{xy}$.
- Goal:** measure the dependence of x and y .

Independence measures

- Given: random variable $(x, y) \in \mathcal{X} \times \mathcal{Y}$, $(x, y) \sim \mathbb{P}_{xy}$.
- Goal:** measure the dependence of x and y .
- Desiderata** for a $Q(\mathbb{P}_{xy})$ independence measure [Rényi, 1959]:
 - $Q(\mathbb{P}_{xy})$ is well-defined,
 - $Q(\mathbb{P}_{xy}) \in [0, 1]$,
 - $Q(\mathbb{P}_{xy}) = 0$ iff. $x \perp y$.
 - $Q(\mathbb{P}_{xy}) = 1$ iff. $y = f(x)$ or $x = g(y)$.

Independence measures

- He showed:

$$Q(\mathbb{P}_{xy}) = \sup_{f,g: \text{ measurable}} \text{corr}(f(x), g(y)),$$

satisfies 1-4.

- Too ambitious:
 - computationally intractable.
 - many measurable functions.

- $C_b(\mathcal{X}) = \{f : \mathcal{X} \text{ metric} \rightarrow \mathbb{R}, \text{ bounded continuous}\}$ would also work.
- Still too large!

- $C_b(\mathcal{X}) = \{f : \mathcal{X} \text{ metric} \rightarrow \mathbb{R}, \text{ bounded continuous}\}$ would also work.
- Still too large!
- Idea:
 - certain RKHS-s are dense in $C_b(\mathcal{X})$.
 - computationally tractable.

KCCA: definition

- Given: $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$.
- Associated:
 - feature maps $\varphi(x) = k(\cdot, x)$, $\psi(y) = \ell(\cdot, y)$,
 - RKHS-s \mathcal{H}_k , \mathcal{H}_ℓ .

KCCA: definition

- Given: $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$.
- Associated:
 - feature maps $\varphi(x) = k(\cdot, x)$, $\psi(y) = \ell(\cdot, y)$,
 - RKHS-s \mathcal{H}_k , \mathcal{H}_ℓ .
- KCCA measure of $(x, y) \in \mathcal{X} \times \mathcal{Y}$

$$\rho_{\text{KCCA}}(x, y; \mathcal{H}_k, \mathcal{H}_\ell) = \sup_{f \in \mathcal{H}_k, g \in \mathcal{H}_\ell} \text{corr}(f(x), g(y)),$$

$$\text{corr}(f(x), g(y)) = \frac{\text{cov}_{xy}(f(x), g(y))}{\sqrt{\text{var}_x f(x) \text{var}_y g(y)}}.$$

- Optimization domain: $\mathcal{H}_k \times \mathcal{H}_\ell \ni (f, g)$.
- By **reproducing property**: we will get a **finite-D task**.
- k, ℓ linear: traditional CCA.
- In **practice**: we have $\{(x_n, y_n)\}_{n=1}^N$ **samples** from (x, y) .

- Optimization domain: $\mathcal{H}_k \times \mathcal{H}_\ell \ni (f, g)$.
- By **reproducing property**: we will get a **finite-D task**.
- k, ℓ linear: traditional CCA.
- In **practice**: we have $\{(x_n, y_n)\}_{n=1}^N$ **samples** from (x, y) .

Recall the reproducing property

$$f(x) = \langle f, k(\cdot, x) \rangle_{\mathcal{H}_k} \quad \forall f \in \mathcal{H}_k, x \in \mathcal{X}.$$

KCCA: empirical estimate

$$\widehat{\text{cov}}_{xy}(f(x), g(y)) = \frac{1}{N} \sum_{n=1}^N \left[\underbrace{f(x_n) - \frac{1}{N} \sum_{i=1}^N f(x_i)}_{\langle f, \varphi(x_n) - \frac{1}{N} \sum_{i=1}^N \varphi(x_i) \rangle_{\mathcal{H}_k}} \right] \left[\underbrace{g(y_n) - \frac{1}{N} \sum_{i=1}^N g(y_i)}_{\langle g, \psi(y_n) - \frac{1}{N} \sum_{i=1}^N \psi(y_i) \rangle_{\mathcal{H}_\ell}} \right]$$

KCCA: empirical estimate

$$\widehat{\text{cov}}_{xy}(f(x), g(y)) = \frac{1}{N} \sum_{n=1}^N \left[\underbrace{f(x_n) - \frac{1}{N} \sum_{i=1}^N f(x_i)}_{\langle f, \varphi(x_n) - \frac{1}{N} \sum_{i=1}^N \varphi(x_i) \rangle_{\mathcal{H}_k}} \right] \left[\underbrace{g(y_n) - \frac{1}{N} \sum_{i=1}^N g(y_i)}_{\langle g, \psi(y_n) - \frac{1}{N} \sum_{i=1}^N \psi(y_i) \rangle_{\mathcal{H}_\ell}} \right]$$
$$= \frac{1}{N} \sum_{n=1}^N \langle \mathbf{f}, \tilde{\varphi}(x_n) \rangle_{\mathcal{H}_k} \langle \mathbf{g}, \tilde{\psi}(y_n) \rangle_{\mathcal{H}_\ell},$$

KCCA: empirical estimate

$$\begin{aligned}\widehat{\text{cov}}_{xy}(f(x), g(y)) &= \frac{1}{N} \sum_{n=1}^N \left[\underbrace{f(x_n) - \frac{1}{N} \sum_{i=1}^N f(x_i)}_{\langle f, \varphi(x_n) - \frac{1}{N} \sum_{i=1}^N \varphi(x_i) \rangle_{\mathcal{H}_k}} \right] \left[\underbrace{g(y_n) - \frac{1}{N} \sum_{i=1}^N g(y_i)}_{\langle g, \psi(y_n) - \frac{1}{N} \sum_{i=1}^N \psi(y_i) \rangle_{\mathcal{H}_\ell}} \right] \\ &= \frac{1}{N} \sum_{n=1}^N \langle \color{blue}f\color{black}, \tilde{\varphi}(x_n) \rangle_{\mathcal{H}_k} \langle \color{red}g\color{black}, \tilde{\psi}(y_n) \rangle_{\mathcal{H}_\ell},\end{aligned}$$

Similarly:

$$\widehat{\text{var}}_x f(x) = \frac{1}{N} \sum_{n=1}^N \left[f(x_n) - \frac{1}{N} \sum_{i=1}^N f(x_i) \right]^2$$

KCCA: empirical estimate

$$\begin{aligned}\widehat{\text{cov}}_{xy}(f(x), g(y)) &= \frac{1}{N} \sum_{n=1}^N \left[\underbrace{f(x_n) - \frac{1}{N} \sum_{i=1}^N f(x_i)}_{\langle f, \varphi(x_n) - \frac{1}{N} \sum_{i=1}^N \varphi(x_i) \rangle_{\mathcal{H}_k}} \right] \left[\underbrace{g(y_n) - \frac{1}{N} \sum_{i=1}^N g(y_i)}_{\langle g, \psi(y_n) - \frac{1}{N} \sum_{i=1}^N \psi(y_i) \rangle_{\mathcal{H}_\ell}} \right] \\ &= \frac{1}{N} \sum_{n=1}^N \langle f, \tilde{\varphi}(x_n) \rangle_{\mathcal{H}_k} \langle g, \tilde{\psi}(y_n) \rangle_{\mathcal{H}_\ell},\end{aligned}$$

Similarly:

$$\widehat{\text{var}}_x f(x) = \frac{1}{N} \sum_{n=1}^N \left[f(x_n) - \frac{1}{N} \sum_{i=1}^N f(x_i) \right]^2 = \frac{1}{N} \sum_{n=1}^N \langle f, \tilde{\varphi}(x_n) \rangle_{\mathcal{H}_k}^2,$$

KCCA: empirical estimate

$$\begin{aligned}\widehat{\text{cov}}_{xy}(f(x), g(y)) &= \frac{1}{N} \sum_{n=1}^N \left[\underbrace{f(x_n) - \frac{1}{N} \sum_{i=1}^N f(x_i)}_{\langle f, \varphi(x_n) - \frac{1}{N} \sum_{i=1}^N \varphi(x_i) \rangle_{\mathcal{H}_k}} \right] \left[\underbrace{g(y_n) - \frac{1}{N} \sum_{i=1}^N g(y_i)}_{\langle g, \psi(y_n) - \frac{1}{N} \sum_{i=1}^N \psi(y_i) \rangle_{\mathcal{H}_\ell}} \right] \\ &= \frac{1}{N} \sum_{n=1}^N \langle f, \tilde{\varphi}(x_n) \rangle_{\mathcal{H}_k} \langle g, \tilde{\psi}(y_n) \rangle_{\mathcal{H}_\ell},\end{aligned}$$

Similarly:

$$\begin{aligned}\widehat{\text{var}}_x f(x) &= \frac{1}{N} \sum_{n=1}^N \left[f(x_n) - \frac{1}{N} \sum_{i=1}^N f(x_i) \right]^2 = \frac{1}{N} \sum_{n=1}^N \langle f, \tilde{\varphi}(x_n) \rangle_{\mathcal{H}_k}^2, \\ \widehat{\text{var}}_y g(y) &= \frac{1}{N} \sum_{n=1}^N \langle g, \tilde{\psi}(y_n) \rangle_{\mathcal{H}_\ell}^2.\end{aligned}$$

KCCA: empirical estimate

- f : appears only as $\langle f, \tilde{\varphi}(x_n) \rangle_{\mathcal{H}_k}$ [similarly: g in $\langle g, \tilde{\psi}(y_n) \rangle_{\mathcal{H}_\ell}$]. \Rightarrow

KCCA: empirical estimate

- f : appears only as $\langle f, \tilde{\varphi}(x_n) \rangle_{\mathcal{H}_k}$ [similarly: g in $\langle g, \tilde{\psi}(y_n) \rangle_{\mathcal{H}_\ell}$]. \Rightarrow
- \forall component of f \perp

$$span \left(\{ \tilde{\varphi}(x_n) \}_{n=1}^N \right) = \left\{ \sum_{n=1}^N c_n \tilde{\varphi}(x_n), \mathbf{c} = [c_n] \in \mathbb{R}^N \right\}$$

has no affect in the objective.

KCCA: empirical estimate

- f : appears only as $\langle f, \tilde{\varphi}(x_n) \rangle_{\mathcal{H}_k}$ [similarly: g in $\langle g, \tilde{\psi}(y_n) \rangle_{\mathcal{H}_\ell}$]. \Rightarrow
- \forall component of $f \perp$

$$\text{span} \left(\{ \tilde{\varphi}(x_n) \}_{n=1}^N \right) = \left\{ \sum_{n=1}^N c_n \tilde{\varphi}(x_n), \mathbf{c} = [c_n] \in \mathbb{R}^N \right\}$$

has no affect in the objective.

Key idea

Enough to consider $f = \sum_{i=1}^N c_i \tilde{\varphi}(x_i)$, $g = \sum_{i=1}^N d_i \tilde{\psi}(y_i)$.

KCCA: empirical estimate

Using that $\mathbf{f} = \sum_{i=1}^N \mathbf{c}_i \tilde{\varphi}(x_i)$, $\mathbf{g} = \sum_{i=1}^N \mathbf{d}_i \tilde{\psi}(y_i)$:

$$\langle f, \tilde{\varphi}(x_n) \rangle_{\mathcal{H}_k} = \sum_{i=1}^N c_i \langle \tilde{\varphi}(x_i), \tilde{\varphi}(x_n) \rangle_{\mathcal{H}_k}$$

KCCA: empirical estimate

Using that $f = \sum_{i=1}^N c_i \tilde{\varphi}(x_i)$, $g = \sum_{i=1}^N d_i \tilde{\psi}(y_i)$:

$$\langle f, \tilde{\varphi}(x_n) \rangle_{\mathcal{H}_k} = \sum_{i=1}^N c_i \langle \tilde{\varphi}(x_i), \tilde{\varphi}(x_n) \rangle_{\mathcal{H}_k} = \sum_{i=1}^N c_i \tilde{k}(x_i, x_n)$$

KCCA: empirical estimate

Using that $\mathbf{f} = \sum_{i=1}^N \mathbf{c}_i \tilde{\varphi}(x_i)$, $\mathbf{g} = \sum_{i=1}^N \mathbf{d}_i \tilde{\psi}(y_i)$:

$$\langle f, \tilde{\varphi}(x_n) \rangle_{\mathcal{H}_k} = \sum_{i=1}^N c_i \langle \tilde{\varphi}(x_i), \tilde{\varphi}(x_n) \rangle_{\mathcal{H}_k} = \sum_{i=1}^N c_i \tilde{k}(x_i, x_n) = (\mathbf{c}^T \tilde{\mathbf{G}}_x)_n,$$

KCCA: empirical estimate

Using that $f = \sum_{i=1}^N c_i \tilde{\varphi}(x_i)$, $g = \sum_{i=1}^N d_i \tilde{\psi}(y_i)$:

$$\begin{aligned}\langle f, \tilde{\varphi}(x_n) \rangle_{\mathcal{H}_k} &= \sum_{i=1}^N c_i \langle \tilde{\varphi}(x_i), \tilde{\varphi}(x_n) \rangle_{\mathcal{H}_k} = \sum_{i=1}^N c_i \tilde{k}(x_i, x_n) = (\mathbf{c}^T \tilde{\mathbf{G}}_x)_n, \\ \langle g, \tilde{\psi}(y_n) \rangle_{\mathcal{H}_\ell} &= (\mathbf{d}^T \tilde{\mathbf{G}}_y)_n,\end{aligned}$$

with the centered kernels $(\tilde{k}, \tilde{\ell})$ and Gram matrices $(\tilde{\mathbf{G}}_x, \tilde{\mathbf{G}}_y)$.

Until now

All the objective terms can be expressed by \mathbf{c} , \mathbf{d} , $\tilde{\mathbf{G}}_x$, $\tilde{\mathbf{G}}_y$.

KCCA: empirical estimate

$$\widehat{\text{cov}}_{xy}(f(x), g(y)) = \frac{1}{N} \sum_{n=1}^N \langle f, \tilde{\varphi}(x_n) \rangle_{\mathcal{H}_k} \langle g, \tilde{\psi}(y_n) \rangle_{\mathcal{H}_\ell},$$

$$\widehat{\text{var}}_x f(x) = \frac{1}{N} \sum_{n=1}^N \langle f, \tilde{\varphi}(x_n) \rangle_{\mathcal{H}_k}^2, \quad \widehat{\text{var}}_y g(y) = \frac{1}{N} \sum_{n=1}^N \langle g, \tilde{\psi}(y_n) \rangle_{\mathcal{H}_\ell}^2,$$

and we have

$$\langle f, \tilde{\varphi}(x_n) \rangle_{\mathcal{H}_k} = (\mathbf{c}^T \tilde{\mathbf{G}}_x)_n, \quad \langle g, \tilde{\psi}(y_n) \rangle_{\mathcal{H}_\ell} = (\mathbf{d}^T \tilde{\mathbf{G}}_y)_n.$$

KCCA: empirical estimate

$$\widehat{\text{cov}}_{xy}(f(x), g(y)) = \frac{1}{N} \sum_{n=1}^N \langle f, \tilde{\varphi}(x_n) \rangle_{\mathcal{H}_k} \langle g, \tilde{\psi}(y_n) \rangle_{\mathcal{H}_\ell},$$

$$\widehat{\text{var}}_x f(x) = \frac{1}{N} \sum_{n=1}^N \langle f, \tilde{\varphi}(x_n) \rangle_{\mathcal{H}_k}^2, \quad \widehat{\text{var}}_y g(y) = \frac{1}{N} \sum_{n=1}^N \langle g, \tilde{\psi}(y_n) \rangle_{\mathcal{H}_\ell}^2,$$

and we have

$$\langle f, \tilde{\varphi}(x_n) \rangle_{\mathcal{H}_k} = (\mathbf{c}^T \tilde{\mathbf{G}}_x)_n, \quad \langle g, \tilde{\psi}(y_n) \rangle_{\mathcal{H}_\ell} = (\mathbf{d}^T \tilde{\mathbf{G}}_y)_n.$$

Thus,

$$\widehat{\text{cov}}_{xy}(f(x), g(y)) = \frac{1}{N} \mathbf{c}^T \tilde{\mathbf{G}}_x \tilde{\mathbf{G}}_y \mathbf{d},$$

$$\widehat{\text{var}}_x f(x) = \frac{1}{N} \mathbf{c}^T (\tilde{\mathbf{G}}_x)^2 \mathbf{c}, \quad \widehat{\text{var}}_y g(y) = \frac{1}{N} \mathbf{d}^T (\tilde{\mathbf{G}}_y)^2 \mathbf{d}.$$

KCCA: finite-D form

Empirical estimate of KCCA:

$$\widehat{\rho_{\text{KCCA}}}^{\text{temp}}(x, y; \mathcal{H}_k, \mathcal{H}_\ell) = \sup_{\mathbf{c} \in \mathbb{R}^N, \mathbf{d} \in \mathbb{R}^N} \frac{\mathbf{c}^T \tilde{\mathbf{G}}_x \tilde{\mathbf{G}}_y \mathbf{d}}{\sqrt{\mathbf{c}^T (\tilde{\mathbf{G}}_x)^2 \mathbf{c}} \sqrt{\mathbf{d}^T (\tilde{\mathbf{G}}_y)^2 \mathbf{d}}}.$$

KCCA: finite-D form

Empirical estimate of KCCA:

$$\widehat{\rho_{\text{KCCA}}}^{\text{temp}}(x, y; \mathcal{H}_k, \mathcal{H}_\ell) = \sup_{\mathbf{c} \in \mathbb{R}^N, \mathbf{d} \in \mathbb{R}^N} \frac{\mathbf{c}^T \tilde{\mathbf{G}}_x \tilde{\mathbf{G}}_y \mathbf{d}}{\sqrt{\mathbf{c}^T (\tilde{\mathbf{G}}_x)^2 \mathbf{c}} \sqrt{\mathbf{d}^T (\tilde{\mathbf{G}}_y)^2 \mathbf{d}}}.$$

In practice ($\kappa > 0$):

$$\begin{aligned} \widehat{\rho_{\text{KCCA}}}(x, y) &:= \widehat{\rho_{\text{KCCA}}}(x, y; \mathcal{H}_k, \mathcal{H}_\ell, \kappa) \\ &= \sup_{\mathbf{c} \in \mathbb{R}^N, \mathbf{d} \in \mathbb{R}^N} \frac{\mathbf{c}^T \tilde{\mathbf{G}}_x \tilde{\mathbf{G}}_y \mathbf{d}}{\sqrt{\mathbf{c}^T (\tilde{\mathbf{G}}_x + \kappa \mathbf{I}_N)^2 \mathbf{c}} \sqrt{\mathbf{d}^T (\tilde{\mathbf{G}}_y + \kappa \mathbf{I}_N)^2 \mathbf{d}}}. \end{aligned}$$

KCCA: finite-D form

Empirical estimate of KCCA:

$$\widehat{\rho_{\text{KCCA}}}^{\text{temp}}(x, y; \mathcal{H}_k, \mathcal{H}_\ell) = \sup_{\mathbf{c} \in \mathbb{R}^N, \mathbf{d} \in \mathbb{R}^N} \frac{\mathbf{c}^T \tilde{\mathbf{G}}_x \tilde{\mathbf{G}}_y \mathbf{d}}{\sqrt{\mathbf{c}^T (\tilde{\mathbf{G}}_x)^2 \mathbf{c}} \sqrt{\mathbf{d}^T (\tilde{\mathbf{G}}_y)^2 \mathbf{d}}}.$$

In practice ($\kappa > 0$):

$$\begin{aligned}\widehat{\rho_{\text{KCCA}}}(x, y) &:= \widehat{\rho_{\text{KCCA}}}(x, y; \mathcal{H}_k, \mathcal{H}_\ell, \kappa) \\ &= \sup_{\mathbf{c} \in \mathbb{R}^N, \mathbf{d} \in \mathbb{R}^N} \frac{\mathbf{c}^T \tilde{\mathbf{G}}_x \tilde{\mathbf{G}}_y \mathbf{d}}{\sqrt{\mathbf{c}^T (\tilde{\mathbf{G}}_x + \kappa \mathbf{I}_N)^2 \mathbf{c}} \sqrt{\mathbf{d}^T (\tilde{\mathbf{G}}_y + \kappa \mathbf{I}_N)^2 \mathbf{d}}}.\end{aligned}$$

Question

How do we solve it?

KCCA: solution

Stationary points of $\widehat{\rho_{\text{KCCA}}}(x, y)$:

$$\mathbf{0} = \frac{\partial \widehat{\rho_{\text{KCCA}}}(x, y)}{\partial \mathbf{c}}, \quad \mathbf{0} = \frac{\partial \widehat{\rho_{\text{KCCA}}}(x, y)}{\partial \mathbf{d}},$$

which simplifies to

$$\tilde{\mathbf{G}}_x \tilde{\mathbf{G}}_y \mathbf{d} = \frac{(\mathbf{c}^T \tilde{\mathbf{G}}_x \tilde{\mathbf{G}}_y \mathbf{d})(\tilde{\mathbf{G}}_x + \kappa \mathbf{I}_N)^2 \mathbf{c}}{\mathbf{c}^T (\tilde{\mathbf{G}}_x + \kappa \mathbf{I}_N)^2 \mathbf{c}}, \quad \tilde{\mathbf{G}}_y \tilde{\mathbf{G}}_x \mathbf{c} = \frac{(\mathbf{d}^T \tilde{\mathbf{G}}_y \tilde{\mathbf{G}}_x \mathbf{c})(\tilde{\mathbf{G}}_y + \kappa \mathbf{I}_N)^2 \mathbf{d}}{\mathbf{d}^T (\tilde{\mathbf{G}}_y + \kappa \mathbf{I}_N)^2 \mathbf{d}}.$$

KCCA: solution

Stationary points of $\widehat{\rho_{\text{KCCA}}}(x, y)$:

$$\mathbf{0} = \frac{\partial \widehat{\rho_{\text{KCCA}}}(x, y)}{\partial \mathbf{c}}, \quad \mathbf{0} = \frac{\partial \widehat{\rho_{\text{KCCA}}}(x, y)}{\partial \mathbf{d}},$$

which simplifies to

$$\tilde{\mathbf{G}}_x \tilde{\mathbf{G}}_y \mathbf{d} = \frac{(\mathbf{c}^T \tilde{\mathbf{G}}_x \tilde{\mathbf{G}}_y \mathbf{d})(\tilde{\mathbf{G}}_x + \kappa \mathbf{I}_N)^2 \mathbf{c}}{\mathbf{c}^T (\tilde{\mathbf{G}}_x + \kappa \mathbf{I}_N)^2 \mathbf{c}}, \quad \tilde{\mathbf{G}}_y \tilde{\mathbf{G}}_x \mathbf{c} = \frac{(\mathbf{d}^T \tilde{\mathbf{G}}_y \tilde{\mathbf{G}}_x \mathbf{c})(\tilde{\mathbf{G}}_y + \kappa \mathbf{I}_N)^2 \mathbf{d}}{\mathbf{d}^T (\tilde{\mathbf{G}}_y + \kappa \mathbf{I}_N)^2 \mathbf{d}}.$$

Normalization:

- (\mathbf{c}, \mathbf{d}) : solution $\Rightarrow (a\mathbf{c}, b\mathbf{d})$: solution $a, b \in \mathbb{R}, \neq 0$.
- denominators := 1.

KCCA: final task

Find the maximal eigenvalue, $\lambda := \mathbf{c}^T \tilde{\mathbf{G}}_x \tilde{\mathbf{G}}_y \mathbf{d}$, of the generalized eigenvalue problem:

$$\begin{bmatrix} \mathbf{0} & \tilde{\mathbf{G}}_x \tilde{\mathbf{G}}_y \\ \tilde{\mathbf{G}}_y \tilde{\mathbf{G}}_x & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{c} \\ \mathbf{d} \end{bmatrix} = \mathbf{c}^T \tilde{\mathbf{G}}_x \tilde{\mathbf{G}}_y \mathbf{d} \begin{bmatrix} (\tilde{\mathbf{G}}_x + \kappa \mathbf{I}_N)^2 & \mathbf{0} \\ \mathbf{0} & (\tilde{\mathbf{G}}_y + \kappa \mathbf{I}_N)^2 \end{bmatrix} \begin{bmatrix} \mathbf{c} \\ \mathbf{d} \end{bmatrix}$$
$$\mathbf{A}\mathbf{z} = \lambda \mathbf{B}\mathbf{z}.$$

KCCA as an independence measure

If $x \perp y$, then $\rho_{\text{KCCA}}(x, y; \mathcal{H}_k, \mathcal{H}_\ell, \kappa) = 0$. Opposite direction:

- For 'rich' $\mathcal{H}_k, \mathcal{H}_\ell$
[Bach and Jordan, 2002, Gretton et al., 2005b].

KCCA as an independence measure

If $x \perp y$, then $\rho_{\text{KCCA}}(x, y; \mathcal{H}_k, \mathcal{H}_\ell, \kappa) = 0$. Opposite direction:

- For 'rich' $\mathcal{H}_k, \mathcal{H}_\ell$
[Bach and Jordan, 2002, Gretton et al., 2005b].
- Enough: universal kernel on a compact metric domain ([later](#)).

KCCA as an independence measure

If $x \perp y$, then $\rho_{\text{KCCA}}(x, y; \mathcal{H}_k, \mathcal{H}_\ell, \kappa) = 0$. Opposite direction:

- For 'rich' $\mathcal{H}_k, \mathcal{H}_\ell$
[Bach and Jordan, 2002, Gretton et al., 2005b].
- Enough: universal kernel on a compact metric domain ([later](#)).
- Example ($\gamma > 0$):
 - Gaussian: $k(x, x') = e^{-\gamma \|x-x'\|_2^2}$.
 - Laplacian kernel: $k(x, x') = e^{-\gamma \|x-x'\|_2}$.

KCCA: regularization

In fact, we estimated

$$\rho_{\text{KCCA}}(x, y; \mathcal{H}_k, \mathcal{H}_\ell, \kappa) = \sup_{f \in \mathcal{H}_k, g \in \mathcal{H}_\ell} \text{corr}(f(x), g(y); \kappa),$$

$$\text{corr}(f(x), g(y); \kappa) = \frac{\text{cov}_{xy}(f(x), g(y))}{\sqrt{\text{var}_x f(x) + \kappa \|f\|_{\mathcal{H}_k}^2} \sqrt{\text{var}_y g(y) + \kappa \|g\|_{\mathcal{H}_\ell}^2}}.$$

KCCA: regularization

In fact, we estimated

$$\rho_{\text{KCCA}}(x, y; \mathcal{H}_k, \mathcal{H}_\ell, \kappa) = \sup_{f \in \mathcal{H}_k, g \in \mathcal{H}_\ell} \text{corr}(f(x), g(y); \kappa),$$

$$\text{corr}(f(x), g(y); \kappa) = \frac{\text{cov}_{xy}(f(x), g(y))}{\sqrt{\text{var}_x f(x) + \kappa \|f\|_{\mathcal{H}_k}^2} \sqrt{\text{var}_y g(y) + \kappa \|g\|_{\mathcal{H}_\ell}^2}}.$$

- **Regularization is important:** With $\kappa = 0, \lambda \in \{0, \pm 1\} \Rightarrow$

$$\rho_{\text{KCCA}}(x, y; \mathcal{H}_k, \mathcal{H}_\ell, \kappa) = 1$$

would be data-independently [Gretton et al., 2005b],
[Bach and Jordan, 2002].

KCCA: regularization

In fact, we estimated

$$\rho_{\text{KCCA}}(x, y; \mathcal{H}_k, \mathcal{H}_\ell, \kappa) = \sup_{f \in \mathcal{H}_k, g \in \mathcal{H}_\ell} \text{corr}(f(x), g(y); \kappa),$$

$$\text{corr}(f(x), g(y); \kappa) = \frac{\text{cov}_{xy}(f(x), g(y))}{\sqrt{\text{var}_x f(x) + \kappa \|f\|_{\mathcal{H}_k}^2} \sqrt{\text{var}_y g(y) + \kappa \|g\|_{\mathcal{H}_\ell}^2}}.$$

- For consistent KCCA estimate:
 - $\kappa_N \rightarrow 0$ [Leurgans et al., 1993] (spline-RKHS),
[Fukumizu et al., 2007] (general RKHS).
 - analysis: covariance operators (later).

KCCA: symmetry, other form

For

$$\begin{bmatrix} \mathbf{0} & \tilde{\mathbf{G}}_x \tilde{\mathbf{G}}_y \\ \tilde{\mathbf{G}}_y \tilde{\mathbf{G}}_x & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{c} \\ \mathbf{d} \end{bmatrix} = \mathbf{c}^T \tilde{\mathbf{G}}_x \tilde{\mathbf{G}}_y \mathbf{d} \begin{bmatrix} (\tilde{\mathbf{G}}_x + \kappa \mathbf{I}_N)^2 & \mathbf{0} \\ \mathbf{0} & (\tilde{\mathbf{G}}_y + \kappa \mathbf{I}_N)^2 \end{bmatrix} \begin{bmatrix} \mathbf{c} \\ \mathbf{d} \end{bmatrix}$$

$([\mathbf{c}, \mathbf{d}], \lambda)$ solution \Rightarrow $([-\mathbf{c}; \mathbf{d}], -\lambda)$: solution. Thus, eigenvalues:

$$\{\lambda_1, -\lambda_1, \dots, \lambda_N, -\lambda_N\}.$$

KCCA: symmetry, other form

For

$$\begin{bmatrix} \mathbf{0} & \tilde{\mathbf{G}}_x \tilde{\mathbf{G}}_y \\ \tilde{\mathbf{G}}_y \tilde{\mathbf{G}}_x & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{c} \\ \mathbf{d} \end{bmatrix} = \mathbf{c}^T \tilde{\mathbf{G}}_x \tilde{\mathbf{G}}_y \mathbf{d} \begin{bmatrix} (\tilde{\mathbf{G}}_x + \kappa \mathbf{I}_N)^2 & \mathbf{0} \\ \mathbf{0} & (\tilde{\mathbf{G}}_y + \kappa \mathbf{I}_N)^2 \end{bmatrix} \begin{bmatrix} \mathbf{c} \\ \mathbf{d} \end{bmatrix}$$

$([\mathbf{c}, \mathbf{d}], \lambda)$ solution \Rightarrow $([-\mathbf{c}; \mathbf{d}], -\lambda)$: solution. Thus, eigenvalues:

$$\{\lambda_1, -\lambda_1, \dots, \lambda_N, -\lambda_N\}.$$

Adding the r.h.s. to both sides:

$$\begin{bmatrix} (\tilde{\mathbf{G}}_x + \kappa \mathbf{I}_N)^2 & \tilde{\mathbf{G}}_x \tilde{\mathbf{G}}_y \\ \tilde{\mathbf{G}}_y \tilde{\mathbf{G}}_x & (\tilde{\mathbf{G}}_y + \kappa \mathbf{I}_N)^2 \end{bmatrix} \begin{bmatrix} \mathbf{c} \\ \mathbf{d} \end{bmatrix} = (1 + \lambda) \begin{bmatrix} (\tilde{\mathbf{G}}_x + \kappa \mathbf{I}_N)^2 & \mathbf{0} \\ \mathbf{0} & (\tilde{\mathbf{G}}_x + \kappa \mathbf{I}_N)^2 \end{bmatrix} \begin{bmatrix} \mathbf{c} \\ \mathbf{d} \end{bmatrix}$$

with eigenvalues $\{1 + \lambda_1, 1 - \lambda_1, \dots, 1 + \lambda_N, 1 - \lambda_N\}$.

KCCA: M -variables

2-variables $[(x, y)]$:

$$\begin{bmatrix} (\tilde{\mathbf{G}}_x + \kappa \mathbf{I}_N)^2 & \tilde{\mathbf{G}}_x \tilde{\mathbf{G}}_y \\ \tilde{\mathbf{G}}_y \tilde{\mathbf{G}}_x & (\tilde{\mathbf{G}}_y + \kappa \mathbf{I}_N)^2 \end{bmatrix} \begin{bmatrix} \mathbf{c} \\ \mathbf{d} \end{bmatrix} = (1 + \lambda) \begin{bmatrix} (\tilde{\mathbf{G}}_x + \kappa \mathbf{I}_N)^2 & \mathbf{0} \\ \mathbf{0} & (\tilde{\mathbf{G}}_x + \kappa \mathbf{I}_N)^2 \end{bmatrix} \begin{bmatrix} \mathbf{c} \\ \mathbf{d} \end{bmatrix}$$

KCCA: M -variables

2-variables $[(x, y)]$:

$$\begin{bmatrix} (\tilde{\mathbf{G}}_x + \kappa \mathbf{I}_N)^2 & \tilde{\mathbf{G}}_x \tilde{\mathbf{G}}_y \\ \tilde{\mathbf{G}}_y \tilde{\mathbf{G}}_x & (\tilde{\mathbf{G}}_y + \kappa \mathbf{I}_N)^2 \end{bmatrix} \begin{bmatrix} \mathbf{c} \\ \mathbf{d} \end{bmatrix} = (1 + \lambda) \begin{bmatrix} (\tilde{\mathbf{G}}_x + \kappa \mathbf{I}_N)^2 & \mathbf{0} \\ \mathbf{0} & (\tilde{\mathbf{G}}_x + \kappa \mathbf{I}_N)^2 \end{bmatrix} \begin{bmatrix} \mathbf{c} \\ \mathbf{d} \end{bmatrix}$$

For M -variables (pairwise dependence):

$$\begin{bmatrix} (\tilde{\mathbf{G}}_1 + \kappa \mathbf{I}_N)^2 & \tilde{\mathbf{G}}_1 \tilde{\mathbf{G}}_2 & \dots & \tilde{\mathbf{G}}_1 \tilde{\mathbf{G}}_M \\ \tilde{\mathbf{G}}_2 \tilde{\mathbf{G}}_1 & (\tilde{\mathbf{G}}_2 + \kappa \mathbf{I}_N)^2 & \dots & \tilde{\mathbf{G}}_2 \tilde{\mathbf{G}}_M \\ \vdots & \vdots & & \vdots \\ \tilde{\mathbf{G}}_M \tilde{\mathbf{G}}_1 & \tilde{\mathbf{G}}_M \tilde{\mathbf{G}}_2 & \dots & (\tilde{\mathbf{G}}_M + \kappa \mathbf{I}_N)^2 \end{bmatrix} \begin{bmatrix} \mathbf{c}_1 \\ \mathbf{c}_2 \\ \vdots \\ \mathbf{c}_M \end{bmatrix} =$$
$$\gamma \begin{bmatrix} (\tilde{\mathbf{G}}_1 + \kappa \mathbf{I}_N)^2 & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & (\tilde{\mathbf{G}}_2 + \kappa \mathbf{I}_N)^2 & \dots & \mathbf{0} \\ \vdots & \vdots & & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & (\tilde{\mathbf{G}}_M + \kappa \mathbf{I}_N)^2 \end{bmatrix} \begin{bmatrix} \mathbf{c}_1 \\ \mathbf{c}_2 \\ \vdots \\ \mathbf{c}_M \end{bmatrix}.$$

Centered Gram matrix

In short

$$\tilde{\mathbf{G}}_x = \mathbf{H}\mathbf{G}_x\mathbf{H} \text{ with } \mathbf{H} = \mathbf{I}_N - \frac{\mathbf{E}_N}{N}; \quad \mathbf{H}, \mathbf{E}_N \in \mathbb{R}^{N \times N}.$$

$$(\tilde{\mathbf{G}}_x)_{ij} = \tilde{k}(x_i, x_j) = \langle \tilde{\varphi}(x_i), \tilde{\varphi}(x_j) \rangle_{\mathcal{H}_k}$$

Centered Gram matrix

In short

$$\tilde{\mathbf{G}}_x = \mathbf{H}\mathbf{G}_x\mathbf{H} \text{ with } \mathbf{H} = \mathbf{I}_N - \frac{\mathbf{E}_N}{N}; \quad \mathbf{H}, \mathbf{E}_N \in \mathbb{R}^{N \times N}.$$

$$\begin{aligned}(\tilde{\mathbf{G}}_x)_{ij} &= \tilde{k}(x_i, x_j) = \langle \tilde{\varphi}(x_i), \tilde{\varphi}(x_j) \rangle_{\mathcal{H}_k} \\&= \left\langle \varphi(x_i) - \frac{1}{N} \sum_{n=1}^N \varphi(x_n), \varphi(x_j) - \frac{1}{N} \sum_{m=1}^N \varphi(x_m) \right\rangle_{\mathcal{H}_k}\end{aligned}$$

Centered Gram matrix

In short

$$\tilde{\mathbf{G}}_x = \mathbf{H}\mathbf{G}_x\mathbf{H} \text{ with } \mathbf{H} = \mathbf{I}_N - \frac{\mathbf{E}_N}{N}; \mathbf{H}, \mathbf{E}_N \in \mathbb{R}^{N \times N}.$$

$$\begin{aligned}(\tilde{\mathbf{G}}_x)_{ij} &= \tilde{k}(x_i, x_j) = \langle \tilde{\varphi}(x_i), \tilde{\varphi}(x_j) \rangle_{\mathcal{H}_k} \\&= \langle \varphi(x_i) - \frac{1}{N} \sum_{n=1}^N \varphi(x_n), \varphi(x_j) - \frac{1}{N} \sum_{m=1}^N \varphi(x_m) \rangle_{\mathcal{H}_k} \\&= (\mathbf{G}_x)_{ij} - \frac{1}{N} \sum_{m=1}^N (\mathbf{G}_x)_{im} - \frac{1}{N} \sum_{n=1}^N (\mathbf{G}_x)_{ni} - \frac{1}{N^2} \sum_{n,m=1}^N (\mathbf{G}_x)_{nm}\end{aligned}$$

Centered Gram matrix

In short

$$\tilde{\mathbf{G}}_x = \mathbf{H}\mathbf{G}_x\mathbf{H} \text{ with } \mathbf{H} = \mathbf{I}_N - \frac{\mathbf{E}_N}{N}; \mathbf{H}, \mathbf{E}_N \in \mathbb{R}^{N \times N}.$$

$$\begin{aligned}(\tilde{\mathbf{G}}_x)_{ij} &= \tilde{k}(x_i, x_j) = \langle \tilde{\varphi}(x_i), \tilde{\varphi}(x_j) \rangle_{\mathcal{H}_k} \\&= \langle \varphi(x_i) - \frac{1}{N} \sum_{n=1}^N \varphi(x_n), \varphi(x_j) - \frac{1}{N} \sum_{m=1}^N \varphi(x_m) \rangle_{\mathcal{H}_k} \\&= (\mathbf{G}_x)_{ij} - \frac{1}{N} \sum_{m=1}^N (\mathbf{G}_x)_{im} - \frac{1}{N} \sum_{n=1}^N (\mathbf{G}_x)_{ni} - \frac{1}{N^2} \sum_{n,m=1}^N (\mathbf{G}_x)_{nm} \\&= \left(\mathbf{G}_x - \mathbf{G}_x \frac{\mathbf{E}_N}{N} - \frac{\mathbf{E}_N}{N} \mathbf{G}_x - \frac{\mathbf{E}_N}{N} \mathbf{G}_x \frac{\mathbf{E}_N}{N} \right)_{ij},\end{aligned}$$

Centered Gram matrix

In short

$$\tilde{\mathbf{G}}_x = \mathbf{H}\mathbf{G}_x\mathbf{H} \text{ with } \mathbf{H} = \mathbf{I}_N - \frac{\mathbf{E}_N}{N}; \mathbf{H}, \mathbf{E}_N \in \mathbb{R}^{N \times N}.$$

$$\begin{aligned}(\tilde{\mathbf{G}}_x)_{ij} &= \tilde{k}(x_i, x_j) = \langle \tilde{\varphi}(x_i), \tilde{\varphi}(x_j) \rangle_{\mathcal{H}_k} \\&= \langle \varphi(x_i) - \frac{1}{N} \sum_{n=1}^N \varphi(x_n), \varphi(x_j) - \frac{1}{N} \sum_{m=1}^N \varphi(x_m) \rangle_{\mathcal{H}_k} \\&= (\mathbf{G}_x)_{ij} - \frac{1}{N} \sum_{m=1}^N (\mathbf{G}_x)_{im} - \frac{1}{N} \sum_{n=1}^N (\mathbf{G}_x)_{ni} - \frac{1}{N^2} \sum_{n,m=1}^N (\mathbf{G}_x)_{nm} \\&= \left(\mathbf{G}_x - \mathbf{G}_x \frac{\mathbf{E}_N}{N} - \frac{\mathbf{E}_N}{N} \mathbf{G}_x - \frac{\mathbf{E}_N}{N} \mathbf{G}_x \frac{\mathbf{E}_N}{N} \right)_{ij}, \\&= (\mathbf{H}\mathbf{G}_x\mathbf{H})_{ij},\end{aligned}$$

\mathbf{H} : symmetric ($\mathbf{H} = \mathbf{H}^T$), idempotent ($\mathbf{H}^2 = \mathbf{H}$).

KCCA: finished.

Mean embedding

Mean embedding: pioneers

- Nonparametric probability distribution representation.
- Late 70s-; survey in [Berlinet and Thomas-Agnan, 2004].

Mean embedding: pioneers

- Nonparametric probability distribution representation.
- Late 70s-; survey in [Berlinet and Thomas-Agnan, 2004].
- **Pioneers in ML:** Bharath Sriperumbudur, Arthur Gretton, Kenji Fukumizu, Alex Smola, Bernhard Schölkopf, Le Song.

Mean embedding: further pointers

- [Names+:](#) Ingo Steinwart, Francis Bach, Dino Sejdinovic, Wittawat Jitkrittum, Krikamol Maundet, Kacper P. Chwialkowski, Ilya Tolstikhin, Carl Johann Simon-Gabriel, David Lopez-Paz, Dougal Sutherland, Aaditya Ramdas, Karsten Borgwardt, Me;)

Mean embedding: further pointers

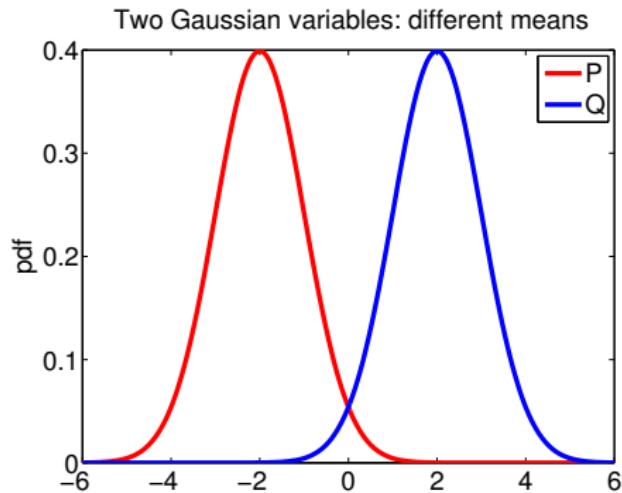
- [Names+:](#) Ingo Steinwart, Francis Bach, Dino Sejdinovic, Wittawat Jitkrittum, Krikamol Maundet, Kacper P. Chwialkowski, Ilya Tolstikhin, Carl Johann Simon-Gabriel, David Lopez-Paz, Dougal Sutherland, Aaditya Ramdas, Karsten Borgwardt, Me;)
- [Wiki:](#) https://en.wikipedia.org/wiki/Kernel_embedding_of_distributions.

Mean embedding: further pointers

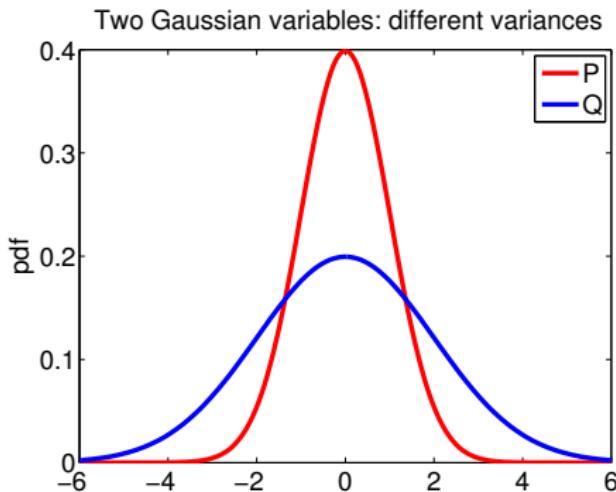
- **Names+:** Ingo Steinwart, Francis Bach, Dino Sejdinovic, Wittawat Jitkrittum, Krikamol Muandet, Kacper P. Chwialkowski, Ilya Tolstikhin, Carl Johann Simon-Gabriel, David Lopez-Paz, Dougal Sutherland, Aaditya Ramdas, Karsten Borgwardt, Me;)
- **Wiki:** https://en.wikipedia.org/wiki/Kernel_embedding_of_distributions.
- **Recent review:** [Muandet et al., 2017].

Towards representations of distributions: EX

- Given: 2 Gaussians with different means.
- Solution: *t*-test.

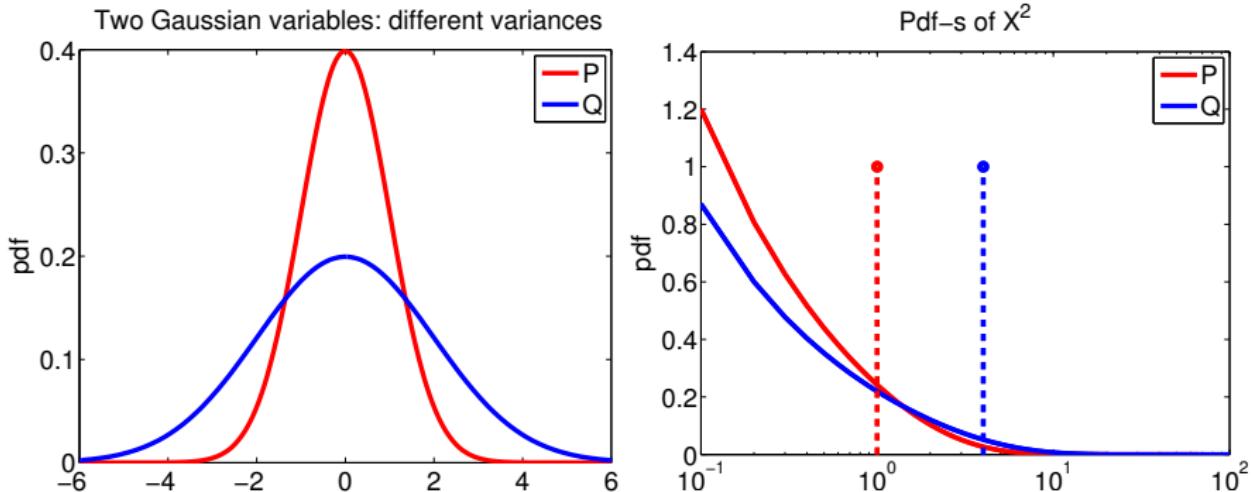


- Setup: 2 Gaussians; same means, different variances.
- Idea: look at the 2nd-order features of RVs.



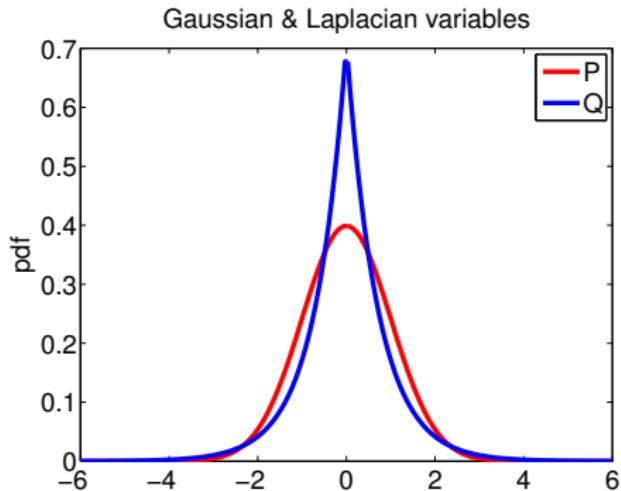
Towards representations of distributions: $\mathbb{E}X^2$

- Setup: 2 Gaussians; same means, different variances.
- Idea: look at the 2nd-order features of RVs.
- $\varphi_x = x^2 \Rightarrow$ difference in $\mathbb{E}X^2$.



Towards representations of distributions: further moments

- Setup: a Gaussian and a Laplacian distribution.
- Challenge: their means *and* variances are the same.
- Idea: look at higher-order features.



Let us consider feature representations!

From kernel trick to mean trick

- Recall:

- $\varphi(x) \in \mathcal{H}_k$: feature of $x \in \mathcal{X}$.
- Kernel: $k(x, x') = \langle \varphi(x), \varphi(x') \rangle_{\mathcal{H}_k}$.

From kernel trick to mean trick

- Recall:
 - $\varphi(x) \in \mathcal{H}_k$: feature of $x \in \mathcal{X}$.
 - Kernel: $k(x, x') = \langle \varphi(x), \varphi(x') \rangle_{\mathcal{H}_k}$.
- Mean embedding:
 - Feature of \mathbb{P} :

$$\mu_{\mathbb{P}} = \mathbb{E}_{x \sim \mathbb{P}}[\varphi(x)] \in \mathcal{H}_k.$$

- Inner product: $\langle \mu_{\mathbb{P}}, \mu_{\mathbb{Q}} \rangle_{\mathcal{H}_k} = \mathbb{E}_{x \sim \mathbb{P}, x' \sim \mathbb{Q}} k(x, x')$.

From kernel trick to mean trick

- Recall:
 - $\varphi(x) \in \mathcal{H}_k$: feature of $x \in \mathcal{X}$.
 - Kernel: $k(x, x') = \langle \varphi(x), \varphi(x') \rangle_{\mathcal{H}_k}$.
 - Mean embedding:
 - Feature of \mathbb{P} :
- $$\mu_{\mathbb{P}} = \mathbb{E}_{x \sim \mathbb{P}}[\varphi(x)] \in \mathcal{H}_k.$$
- Inner product: $\langle \mu_{\mathbb{P}}, \mu_{\mathbb{Q}} \rangle_{\mathcal{H}_k} = \mathbb{E}_{x \sim \mathbb{P}, x' \sim \mathbb{Q}} k(x, x')$.
 - $\mu_{\mathbb{P}}$: well-defined for all distributions (bounded k).

From kernel trick to mean trick

- Recall:
 - $\varphi(x) \in \mathcal{H}_k$: feature of $x \in \mathcal{X}$.
 - Kernel: $k(x, x') = \langle \varphi(x), \varphi(x') \rangle_{\mathcal{H}_k}$.
- Mean embedding:

- Feature of \mathbb{P} :

$$\mu_{\mathbb{P}} = \mathbb{E}_{x \sim \mathbb{P}}[\varphi(x)] \in \mathcal{H}_k.$$

- Inner product: $\langle \mu_{\mathbb{P}}, \mu_{\mathbb{Q}} \rangle_{\mathcal{H}_k} = \mathbb{E}_{x \sim \mathbb{P}, x' \sim \mathbb{Q}} k(x, x')$.
- $\mu_{\mathbb{P}}$: well-defined for all distributions (bounded k).

Commonly used construction

$\mu_{\mathbb{P}} = \mathbb{E}_{x \sim \mathbb{P}}[\varphi(x)]$. Indeed...

Distribution Representation via Functions

- Cumulative density function:

$$\mathbb{P} \mapsto F(z) = \mathbb{P}(x < z)$$

Distribution Representation via Functions

- Cumulative density function:

$$\mathbb{P} \mapsto F(z) = \mathbb{P}(x < z) = \mathbb{E}_{x \sim \mathbb{P}} \chi_{(-\infty, z)}(x).$$

Distribution Representation via Functions

- Cumulative density function:

$$\mathbb{P} \mapsto F(z) = \mathbb{P}(x < z) = \mathbb{E}_{x \sim \mathbb{P}} \chi_{(-\infty, z)}(x).$$

- Characteristic function:

$$\mathbb{P} \mapsto \phi_{\mathbb{P}}(z) = \int e^{i \langle z, x \rangle} d\mathbb{P}(x).$$

Distribution Representation via Functions

- Cumulative density function:

$$\mathbb{P} \mapsto F(z) = \mathbb{P}(x < z) = \mathbb{E}_{x \sim \mathbb{P}} \chi_{(-\infty, z)}(x).$$

- Characteristic function:

$$\mathbb{P} \mapsto \phi_{\mathbb{P}}(z) = \int e^{i\langle z, x \rangle} d\mathbb{P}(x).$$

- Moment generating function:

$$\mathbb{P} \mapsto M_{\mathbb{P}}(z) = \int e^{\langle z, x \rangle} d\mathbb{P}(x).$$

Distribution Representation via Functions

- Cumulative density function:

$$\mathbb{P} \mapsto F(z) = \mathbb{P}(x < z) = \mathbb{E}_{x \sim \mathbb{P}} \chi_{(-\infty, z)}(x).$$

- Characteristic function:

$$\mathbb{P} \mapsto \phi_{\mathbb{P}}(z) = \int e^{i \langle z, x \rangle} d\mathbb{P}(x).$$

- Moment generating function:

$$\mathbb{P} \mapsto M_{\mathbb{P}}(z) = \int e^{\langle z, x \rangle} d\mathbb{P}(x).$$

Pattern

$$\mathbb{P} \mapsto \mu_{\mathbb{P}} = \int_{\mathcal{X}} \varphi(x) d\mathbb{P}(x)$$

Distribution Representation via Functions

- Cumulative density function:

$$\mathbb{P} \mapsto F(z) = \mathbb{P}(x < z) = \mathbb{E}_{x \sim \mathbb{P}} \chi_{(-\infty, z)}(x).$$

- Characteristic function:

$$\mathbb{P} \mapsto \phi_{\mathbb{P}}(z) = \int e^{i \langle z, x \rangle} d\mathbb{P}(x).$$

- Moment generating function:

$$\mathbb{P} \mapsto M_{\mathbb{P}}(z) = \int e^{\langle z, x \rangle} d\mathbb{P}(x).$$

Pattern

$\mathbb{P} \mapsto \mu_{\mathbb{P}} = \int_{\mathcal{X}} \varphi(x) d\mathbb{P}(x)$, in our case: $\varphi(x) = k(\cdot, x)$.

Bochner integral: quick summary [Diestel and Uhl, 1977, Dinculeanu, 2000, Steinwart and Christmann, 2008]

- Given:
 - $(\mathcal{X}, \mathcal{A}, \mu)$: σ -finite measure space,
 - $f : (\mathcal{X}, \mathcal{A}) \rightarrow B$ (anach space)-valued measurable function.

Bochner integral: quick summary [Diestel and Uhl, 1977, Dinculeanu, 2000, Steinwart and Christmann, 2008]

- Given:
 - $(\mathcal{X}, \mathcal{A}, \mu)$: σ -finite measure space,
 - $f : (\mathcal{X}, \mathcal{A}) \rightarrow B$ (anach space)-valued measurable function.
- For $f = \sum_{i=1}^n c_i \chi_{A_i}$ ($A_i \in \mathcal{A}, c_i \in B$) **measurable step functions**

$$\int_{\mathcal{X}} f d\mu := \sum_{i=1}^n c_i \mu(A_i) \in B.$$

Bochner integral: quick summary [Diestel and Uhl, 1977, Dinculeanu, 2000, Steinwart and Christmann, 2008]

- Given:
 - $(\mathcal{X}, \mathcal{A}, \mu)$: σ -finite measure space,
 - $f : (\mathcal{X}, \mathcal{A}) \rightarrow B$ (anach space)-valued measurable function.
- For $f = \sum_{i=1}^n c_i \chi_{A_i}$ ($A_i \in \mathcal{A}, c_i \in B$) **measurable step functions**

$$\int_{\mathcal{X}} f d\mu := \sum_{i=1}^n c_i \mu(A_i) \in B.$$

- f **measurable function** is Bochner μ -integrable if
 - $\exists (f_n)$ measurable step functions: $\lim_{n \rightarrow \infty} \int_{\mathcal{X}} \|f - f_n\|_B d\mu = 0$.
 - In this case $\lim_{n \rightarrow \infty} \int_{\mathcal{X}} f_n d\mu$ exists, $=: \int_{\mathcal{X}} f d\mu$.

Bochner integral: properties

- $f : \mathcal{X} \rightarrow B$ is Bochner integrable $\Leftrightarrow \int_{\mathcal{X}} \|f\|_B \, d\mu < \infty$.

Bochner integral: properties

- $f : \mathcal{X} \rightarrow B$ is Bochner integrable $\Leftrightarrow \int_{\mathcal{X}} \|f\|_B d\mu < \infty$.
- In this case $\|\int_{\mathcal{X}} f d\mu\|_B \leq \int_{\mathcal{X}} \|f\|_B d\mu$. ('Jensen inequality')

Bochner integral: properties

- $f : \mathcal{X} \rightarrow B$ is Bochner integrable $\Leftrightarrow \int_{\mathcal{X}} \|f\|_B \, d\mu < \infty$.
- In this case $\left\| \int_{\mathcal{X}} f \, d\mu \right\|_B \leq \int_{\mathcal{X}} \|f\|_B \, d\mu$. ('Jensen inequality')
- If
 - $S : B \rightarrow B_2$: bounded linear operator,
 - $f : X \rightarrow B$: Bochner integrable, then

$S \circ f : X \rightarrow B_2$ is Bochner integrable and

$$S \left(\int_{\mathcal{X}} f \, d\mu \right) = \int_{\mathcal{X}} Sf \, d\mu.$$

Bochner integral: properties

- $f : \mathcal{X} \rightarrow B$ is Bochner integrable $\Leftrightarrow \int_{\mathcal{X}} \|f\|_B d\mu < \infty$.
- In this case $\|\int_{\mathcal{X}} f d\mu\|_B \leq \int_{\mathcal{X}} \|f\|_B d\mu$. ('Jensen inequality')
- If
 - $S : B \rightarrow B_2$: bounded linear operator,
 - $f : X \rightarrow B$: Bochner integrable, then

$S \circ f : X \rightarrow B_2$ is Bochner integrable and

$$S \left(\int_{\mathcal{X}} f d\mu \right) = \int_{\mathcal{X}} Sf d\mu.$$

In short

$|\int f d\mu| \leq \int |f| d\mu$ and $c \int f d\mu = \int c f d\mu$ generalize nicely.

Mean embedding: \exists , $\mathbb{E}_{\mathbb{P}}$ -reproducing property

Given:

- $(\mathcal{X}, \mathcal{A})$ measurable space,
- $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ kernel.

Theorem

$\mu_{\mathbb{P}} := \int_{\mathcal{X}} k(\cdot, x) d\mathbb{P}(x)$ exists, $\mu_{\mathbb{P}} \in \mathcal{H}_k$, and

$$\mathbb{P}f := \mathbb{E}_{x \sim \mathbb{P}} f(x) = \langle f, \mu_{\mathbb{P}} \rangle_{\mathcal{H}_k} \quad \forall f \in \mathcal{H}_k$$

under mild conditions:

- $\mathbb{E}_{x \sim \mathbb{P}} \sqrt{k(x, x)} < \infty$, and
- $y \mapsto k(y, x)$ is measurable for any $x \in \mathcal{X}$.

Existence of $\mu_{\mathbb{P}}$: proof

- $\exists \int_{\mathcal{X}} k(\cdot, x) d\mathbb{P}(x) \ (\& \in \mathcal{H}_k) \Leftrightarrow$

$$\infty > \int_{\mathcal{X}} \|k(\cdot, x)\|_{\mathcal{H}_k} d\mathbb{P}(x) = \mathbb{E}_{x \sim \mathbb{P}} \sqrt{k(x, x)}.$$

Existence of $\mu_{\mathbb{P}}$: proof

- $\exists \int_{\mathcal{X}} k(\cdot, x) d\mathbb{P}(x)$ ($\& \in \mathcal{H}_k$) \Leftrightarrow

$$\infty > \int_{\mathcal{X}} \|k(\cdot, x)\|_{\mathcal{H}_k} d\mathbb{P}(x) = \mathbb{E}_{x \sim \mathbb{P}} \sqrt{k(x, x)}.$$

- $\mathbb{E}_{x \sim \mathbb{P}} f(x) = \mathbb{E}_{x \sim \mathbb{P}} \langle f, k(\cdot, x) \rangle_{\mathcal{H}_k} = \langle f, \mathbb{E}_{x \sim \mathbb{P}} k(\cdot, x) \rangle_{\mathcal{H}_k} = \langle f, \mu_{\mathbb{P}} \rangle_{\mathcal{H}_k}$ by
 - reproducing property of k ,
 - $g \in \mathcal{H}_k \mapsto \langle f, g \rangle \in \mathbb{R}$: bounded linear ($S \leftrightarrow \int$).

Existence of $\mu_{\mathbb{P}}$: proof

- $\exists \int_{\mathcal{X}} k(\cdot, x) d\mathbb{P}(x)$ ($\& \in \mathcal{H}_k$) \Leftrightarrow

$$\infty > \int_{\mathcal{X}} \|k(\cdot, x)\|_{\mathcal{H}_k} d\mathbb{P}(x) = \mathbb{E}_{x \sim \mathbb{P}} \sqrt{k(x, x)}.$$

- $\mathbb{E}_{x \sim \mathbb{P}} f(x) = \mathbb{E}_{x \sim \mathbb{P}} \langle f, k(\cdot, x) \rangle_{\mathcal{H}_k} = \langle f, \mathbb{E}_{x \sim \mathbb{P}} k(\cdot, x) \rangle_{\mathcal{H}_k} = \langle f, \mu_{\mathbb{P}} \rangle_{\mathcal{H}_k}$ by
 - reproducing property of k ,
 - $g \in \mathcal{H}_k \mapsto \langle f, g \rangle \in \mathbb{R}$: bounded linear ($S \leftrightarrow \int$).
- Measurability of $x \in \mathcal{X} \mapsto k(\cdot, x) \in \mathcal{H}_k$: $\Leftrightarrow y \mapsto k(y, x)$ is measurable $\forall x$ [Berlinet and Thomas-Agnan, 2004].

Mean embedding: specific cases

For

- $k(x, x') = e^{\langle x, x' \rangle}$: $\mu_{\mathbb{P}}$ = moment generating function of \mathbb{P} .
- $k(x, y) = e^{i\langle x, y \rangle}$: $\mu_{\mathbb{P}}$ = characteristic function of \mathbb{P} .
 - Only formally: $k(x, y) = k(y, x)^*$ fails.
- $\mathbb{P} = \delta_x$, $\mu_{\mathbb{P}} = k(\cdot, x)$.

Condition:

- $y \mapsto k(y, x)$ is measurable $\forall x$: super-mild.
- $\mathbb{E}_{x \sim \mathbb{P}} \sqrt{k(x, x)} < \infty$: holds for **bounded kernels**, i.e. when

$$\sup_{x, x' \in \mathcal{X}} k(x, x') \leq B_k < \infty.$$

Mean embedding: empirical estimate

- $\mu_{\mathbb{P}}$: typically **analytically not available**.
- Empirical estimate: from $\{x_i\}_{i=1}^n \stackrel{i.i.d.}{\sim} \mathbb{P}$

$$\widehat{\mu}_{\mathbb{P}} = \frac{1}{n} \sum_{i=1}^n k(\cdot, x_i) = \mu_{\mathbb{P}_n} \in \mathcal{H}_k,$$

where $\mathbb{P}_n = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$ is the empirical measure.

Empirical mean embedding: finite-sample guarantees

Theorem ([Altun and Smola, 2006])

For a *k bounded* kernel $[\sup_{x,y \in \mathcal{X}} k(x,y) \leq B_k]$, with probability $\geq 1 - \delta$

$$\|\mu_{\mathbb{P}} - \mu_{\mathbb{P}_n}\|_{\mathcal{H}_k} \leq \frac{\left[1 + \sqrt{\log\left(\frac{1}{\delta}\right)}\right] \sqrt{2B_k}}{\sqrt{n}}.$$

Finite-sample guarantee: proof idea

- $g(x_1, \dots, x_n) = \|\mu_{\mathbb{P}} - \mu_{\mathbb{P}_n}\|_{\mathcal{H}_k}$: bounded difference property \Rightarrow
- McDiarmid inequality: concentration around $\mathbb{E}g$.
- $\mathbb{E}g \leq$ expected kernel values (B_k appears).

Finite-sample guarantee: note

Alternative of

$$\mathbb{P} \left(\|\mu_{\mathbb{P}} - \mu_{\mathbb{P}_n}\|_{\mathcal{H}_k} \leqslant \frac{\left[1 + \sqrt{\log \left(\frac{1}{\delta} \right)} \right] \sqrt{2B_k}}{\sqrt{n}} \right) \geqslant 1 - \delta.$$

Directly by the Bernstein inequality [Caponnetto and De Vito, 2007]:

$$\mathbb{P} \left(\|\mu_{\mathbb{P}} - \mu_{\mathbb{P}_n}\|_{\mathcal{H}_k} \leqslant 2\sqrt{B_k} \left[\frac{2}{n} + \frac{1}{\sqrt{n}} \log \left(\frac{2}{\delta} \right) \right] \right) \geqslant 1 - \delta$$

would give a bit **worse** dependence.

- Mean embeddings define a semi-metric (MMD):

$$d_k(\mathbb{P}, \mathbb{Q}) := \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k}.$$

- Mean embeddings define a semi-metric (MMD):

$$d_k(\mathbb{P}, \mathbb{Q}) := \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k}.$$

- d_k is metric $\Leftrightarrow \mathbb{P} \mapsto \mu_{\mathbb{P}}$ is injective.
- Characteristic kernel [Fukumizu et al., 2004, Fukumizu et al., 2008]:
 - characteristic function analogy.
 - L -order polynomial kernel: encodes moments $\leq L$. (not)

Mean embedding: universality (k)

Let $C(\mathcal{X}) = \{f : \mathcal{X} \rightarrow \mathbb{R} \text{ continuous}\}$.

Definition

Assume:

- \mathcal{X} : compact metric space.
- k : continuous kernel on \mathcal{X} .

k is called *(c)-universal* [Steinwart, 2001] if \mathcal{H}_k is dense in $(C(\mathcal{X}), \|\cdot\|_\infty)$.

Let $C(\mathcal{X}) = \{f : \mathcal{X} \rightarrow \mathbb{R} \text{ continuous}\}$.

Definition

Assume:

- \mathcal{X} : compact metric space.
- k : continuous kernel on \mathcal{X} .

k is called *(c)-universal* [Steinwart, 2001] if \mathcal{H}_k is dense in $(C(\mathcal{X}), \|\cdot\|_\infty)$.

\mathcal{X} assumption \Rightarrow

$C(\mathcal{X}) = C_b(\mathcal{X}) = \{f : \mathcal{X} \rightarrow \mathbb{R} \text{ continuous bounded}\}$

$\mathcal{H}_k \subset C(\mathcal{X})$? Non-compact spaces?

Notes:

- k : continuous, \mathcal{X} : compact $\Rightarrow k$: bounded.
- k : continuous, bounded $\Rightarrow \mathcal{H}_k \subset C(\mathcal{X})$
[Steinwart and Christmann, 2008].

$\mathcal{H}_k \subset C(\mathcal{X})$? Non-compact spaces?

Notes:

- Extensions of c-universality to non-compact spaces:
 - c_0 -universality, cc-universality,
... [Carmeli et al., 2010, Sriperumbudur et al., 2010a, Simon-Gabriel and Schölkopf, 2018].

≥ 3 different proof options:

- ① [Micchelli et al., 2006]: k is c-universal $\Leftrightarrow \mu$ is injective on $\mathcal{M}_b(\mathcal{X})$, the set of finite signed Borel measures on \mathcal{X} .

≥ 3 different proof options:

- ① [Micchelli et al., 2006]: k is c-universal $\Leftrightarrow \mu$ is injective on $\mathcal{M}_b(\mathcal{X})$, the set of finite signed Borel measures on \mathcal{X} .
- ② Direct reasoning [Gretton et al., 2012].

≥ 3 different proof options:

- ① [Micchelli et al., 2006]: k is c-universal $\Leftrightarrow \mu$ is injective on $\mathcal{M}_b(\mathcal{X})$, the set of finite signed Borel measures on \mathcal{X} .
- ② Direct reasoning [Gretton et al., 2012].
- ③ Denseness of $\mathcal{H}_k + \mathbb{R}$ in $L^2(\mathbb{P})$
[Fukumizu et al., 2008, Fukumizu et al., 2009a].

k : universal $\Rightarrow k$: characteristic

≥ 3 different proof options:

- ① [Micchelli et al., 2006]: k is c-universal $\Leftrightarrow \mu$ is injective on $\mathcal{M}_b(\mathcal{X})$, the set of finite signed Borel measures on \mathcal{X} .
- ② Direct reasoning [Gretton et al., 2012].
- ③ Denseness of $\mathcal{H}_k + \mathbb{R}$ in $L^2(\mathbb{P})$
[Fukumizu et al., 2008, Fukumizu et al., 2009a].

Let us construct some *examples* first! (then prove 1-2)

Properties of universal kernels

[Steinwart, 2001, Steinwart and Christmann, 2008]

If k is universal, then

- $k(x, x) > 0$ for all $x \in \mathcal{X}$.

Properties of universal kernels

[Steinwart, 2001, Steinwart and Christmann, 2008]

If k is universal, then

- $k(x, x) > 0$ for all $x \in \mathcal{X}$.
- Every restriction of k to an $\mathcal{X}' \subseteq \mathcal{X}$ compact set is universal.

Properties of universal kernels

[Steinwart, 2001, Steinwart and Christmann, 2008]

If k is universal, then

- $k(x, x) > 0$ for all $x \in \mathcal{X}$.
- Every restriction of k to an $\mathcal{X}' \subseteq \mathcal{X}$ compact set is universal.
- $\varphi(x) = k(\cdot, x)$ is injective, i.e.

$$\rho_k(x, y) = \|\varphi(x) - \varphi(y)\|_{\mathcal{H}_k}$$

is a metric.

Properties of universal kernels

[Steinwart, 2001, Steinwart and Christmann, 2008]

If k is universal, then

- $k(x, x) > 0$ for all $x \in \mathcal{X}$.
- Every restriction of k to an $\mathcal{X}' \subseteq \mathcal{X}$ compact set is universal.
- $\varphi(x) = k(\cdot, x)$ is injective, i.e.

$$\rho_k(x, y) = \|\varphi(x) - \varphi(y)\|_{\mathcal{H}_k}$$

is a metric.

- The normalized kernel (recall: corr)

$$\tilde{k}(x, y) := \frac{k(x, y)}{\sqrt{k(x, x)k(y, y)}}$$

is universal.

Universal Taylor kernels

[Steinwart, 2001, Steinwart and Christmann, 2008]

- For an $C^\infty \ni f : (-r, r) \rightarrow \mathbb{R}$

$$f(t) = \sum_{n=0}^{\infty} \textcolor{blue}{a_n} t^n \quad t \in (-r, r), \quad r \in (0, \infty].$$

Universal Taylor kernels

[Steinwart, 2001, Steinwart and Christmann, 2008]

- For an $C^\infty \ni f : (-r, r) \rightarrow \mathbb{R}$

$$f(t) = \sum_{n=0}^{\infty} a_n t^n \quad t \in (-r, r), \quad r \in (0, \infty].$$

- If $a_n > 0 \ \forall n$, then

$$k(x, y) = f(\langle x, y \rangle)$$

is **universal** on $\mathcal{X} := \{x \in \mathbb{R}^d : \|x\|_2 \leq \sqrt{r}\}$.

Universal kernels on compact subsets of \mathbb{R}^d , $\alpha > 0$

- $k(x, y) = e^{\alpha \langle x, y \rangle}$: previous result with $f(t) = e^{\alpha t} \Rightarrow a_n = \frac{\alpha^n}{n!}$.

Universal kernels on compact subsets of \mathbb{R}^d , $\alpha > 0$

- $k(x, y) = e^{\alpha \langle x, y \rangle}$: previous result with $f(t) = e^{\alpha t} \Rightarrow a_n = \frac{\alpha^n}{n!}$.
- $k(x, y) = e^{-\alpha \|x - y\|_2^2}$: exp. kernel & normalization.

Universal kernels on compact subsets of \mathbb{R}^d , $\alpha > 0$

- $k(x, y) = (1 - \langle x, y \rangle)^{-\alpha}$ binomial kernel
 - on \mathcal{X} compact $\subset \{x \in \mathbb{R}^d : \|x\|_2 < 1\}$.
 - $f(t) = (1 - t)^{-\alpha} = \sum_{n=0}^{\infty} \underbrace{\binom{-\alpha}{n}}_{>0} (-1)^n t^n \quad (|t| < 1),$

where $\binom{b}{n} = \sum_{i=1}^n \frac{b-i+1}{i}$.

Universal \Rightarrow characteristic: proof-1

Injectivity on finite signed measures (proof):

- k : universal $\Rightarrow \mathcal{H}_k$ is dense in $C(\mathcal{X})$.

Universal \Rightarrow characteristic: proof-1

Injectivity on finite signed measures (proof):

- k : universal $\Rightarrow \mathcal{H}_k$ is dense in $C(\mathcal{X})$.
- Hahn-Banach theorem [Rudin, 1991]: Let H be a subspace of a normed space C . H is dense in C iff.

$$\{0\} = H^\perp := \{F \in C' : \forall f \in H, F(f) = 0\}.$$

Universal \Rightarrow characteristic: proof-1

Injectivity on finite signed measures (proof):

- k : universal $\Rightarrow \mathcal{H}_k$ is dense in $C(\mathcal{X})$.
- Hahn-Banach theorem [Rudin, 1991]: Let H be a subspace of a normed space C . H is dense in C iff.

$$\{0\} = H^\perp := \{F \in C' : \forall f \in H, F(f) = 0\}.$$

- Denseness \Leftrightarrow

$$\{0\} = \mathcal{H}_k^\perp = \left\{ \mathbb{F} \in \underbrace{C(\mathcal{X})'}_{=\mathcal{M}_b(\mathcal{X})} : \forall f \in \mathcal{H}_k, 0 = T_{\mathbb{F}}(f) = \underbrace{\int_{\mathcal{X}} f d\mathbb{F}}_{\langle f, \mu_{\mathbb{F}} \rangle_{\mathcal{H}_k}} \right\}$$

Universal \Rightarrow characteristic: proof-1

Injectivity on finite signed measures (proof):

- k : universal $\Rightarrow \mathcal{H}_k$ is dense in $C(\mathcal{X})$.
- Hahn-Banach theorem [Rudin, 1991]: Let H be a subspace of a normed space C . H is dense in C iff.

$$\{0\} = H^\perp := \{F \in C' : \forall f \in H, F(f) = 0\}.$$

- Denseness \Leftrightarrow

$$\begin{aligned}\{0\} &= \mathcal{H}_k^\perp = \left\{ \mathbb{F} \in \underbrace{C(\mathcal{X})'}_{=\mathcal{M}_b(\mathcal{X})} : \forall f \in \mathcal{H}_k, 0 = T_{\mathbb{F}}(f) = \underbrace{\int_{\mathcal{X}} f d\mathbb{F}}_{\langle f, \mu_{\mathbb{F}} \rangle_{\mathcal{H}_k}} \right\} \\ &= \{ \mathbb{F} \in \mathcal{M}_b(\mathcal{X}) : \mu_{\mathbb{F}} = 0 \}.\end{aligned}$$

Universal \Rightarrow characteristic: proof-2

Direct reasoning: We have already 'mentioned' [Dudley, 2004]:

- Let \mathcal{X} : metric space, $\mathbb{P}, \mathbb{Q} \in \mathcal{M}_1^+(\mathcal{X})$.
- Then $\mathbb{P} = \mathbb{Q}$ (Borel probability measures) \Leftrightarrow

$$\mathbb{P}f = \mathbb{Q}f \left(:= \int_{\mathcal{X}} f(x) d\mathbb{Q}(x) \right) \quad \forall f \in C_b(\mathcal{X}).$$

We have a characterization of $\mathbb{P} = \mathbb{Q}$ in terms of expectations.

Universal \Rightarrow characteristic: proof-2

- Goal: $\mu_{\mathbb{P}} = \mu_{\mathbb{Q}} \Rightarrow \mathbb{P} = \mathbb{Q}$ [$\Leftrightarrow \mathbb{P}f = \mathbb{Q}f, \forall f \in C_b(\mathcal{X})$].

Universal \Rightarrow characteristic: proof-2

- Goal: $\mu_{\mathbb{P}} = \mu_{\mathbb{Q}} \Rightarrow \mathbb{P} = \mathbb{Q}$ [$\Leftrightarrow \mathbb{P}f = \mathbb{Q}f, \forall f \in C_b(\mathcal{X})$].
- We want: for any $f \in C_b(\mathcal{X})$ and $\epsilon > 0$, $|\mathbb{P}f - \mathbb{Q}f| \stackrel{?}{\leqslant} \epsilon$.

Universal \Rightarrow characteristic: proof-2

- Goal: $\mu_{\mathbb{P}} = \mu_{\mathbb{Q}} \Rightarrow \mathbb{P} = \mathbb{Q}$ [$\Leftrightarrow \mathbb{P}f = \mathbb{Q}f, \forall f \in C_b(\mathcal{X})$].
- We want: for any $f \in C_b(\mathcal{X})$ and $\epsilon > 0$, $|\mathbb{P}f - \mathbb{Q}f| \stackrel{?}{\leqslant} \epsilon$.
- Universality of $k \Rightarrow \mathcal{H}_k$ is **dense** in $C_b(\mathcal{X})$.

Universal \Rightarrow characteristic: proof-2

- Goal: $\mu_{\mathbb{P}} = \mu_{\mathbb{Q}} \Rightarrow \mathbb{P} = \mathbb{Q}$ [$\Leftrightarrow \mathbb{P}f = \mathbb{Q}f, \forall f \in C_b(\mathcal{X})$].
- We want: for any $f \in C_b(\mathcal{X})$ and $\epsilon > 0$, $|\mathbb{P}f - \mathbb{Q}f| \stackrel{?}{\leq} \epsilon$.
- Universality of $k \Rightarrow \mathcal{H}_k$ is **dense** in $C_b(\mathcal{X})$.
- $\mathcal{H}_k \ni g := \epsilon$ -approximation of f ,

$$|\mathbb{P}f - \mathbb{Q}f| \leq \underbrace{|\mathbb{P}f - \mathbb{P}g|}_{\leq \mathbb{P}|f-g| \leq \epsilon} + \underbrace{|\mathbb{P}g - \mathbb{Q}g|}_{\stackrel{?}{\leq} \epsilon} + \underbrace{|\mathbb{Q}g - \mathbb{Q}f|}_{\leq \epsilon},$$

Universal \Rightarrow characteristic: proof-2

- Goal: $\mu_{\mathbb{P}} = \mu_{\mathbb{Q}} \Rightarrow \mathbb{P} = \mathbb{Q}$ [$\Leftrightarrow \mathbb{P}f = \mathbb{Q}f, \forall f \in C_b(\mathcal{X})$].
- We want: for any $f \in C_b(\mathcal{X})$ and $\epsilon > 0$, $|\mathbb{P}f - \mathbb{Q}f| \stackrel{?}{\leq} \epsilon$.
- Universality of $k \Rightarrow \mathcal{H}_k$ is **dense** in $C_b(\mathcal{X})$.
- $\mathcal{H}_k \ni g := \epsilon\text{-approximation of } f$,

$$|\mathbb{P}f - \mathbb{Q}f| \leq \underbrace{|\mathbb{P}f - \mathbb{P}g|}_{\leq \mathbb{P}|f-g| \leq \epsilon} + |\mathbb{P}g - \mathbb{Q}g| + \underbrace{|\mathbb{Q}g - \mathbb{Q}f|}_{\leq \epsilon},$$

$$|\mathbb{P}g - \mathbb{Q}g| = \left| \underbrace{\langle g, \mu_{\mathbb{P}} \rangle_{\mathcal{H}_k} - \langle g, \mu_{\mathbb{Q}} \rangle_{\mathcal{H}_k}}_{\langle g, \underbrace{\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}}_{=0} \rangle_{\mathcal{H}_k}} \right| = 0. \text{ Thus } |\mathbb{P}f - \mathbb{Q}f| \leq 2\epsilon.$$

Universality: finished. Now: characteristic
property.

[Gretton et al., 2007]:

$$d_k^2(\mathbb{P}, \mathbb{Q}) = \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k}^2 = \left\| \int_{\mathcal{X}} k(\cdot, x) d\mathbb{P}(x) - \int_{\mathcal{X}} k(\cdot, y) d\mathbb{Q}(y) \right\|_{\mathcal{H}_k}^2$$

$d_k(\mathbb{P}, \mathbb{Q})$ (=MMD) in terms of kernel evaluations

[Gretton et al., 2007]:

$$\begin{aligned} d_k^2(\mathbb{P}, \mathbb{Q}) &= \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k}^2 = \left\| \int_{\mathcal{X}} k(\cdot, x) d\mathbb{P}(x) - \int_{\mathcal{X}} k(\cdot, y) d\mathbb{Q}(y) \right\|_{\mathcal{H}_k}^2 \\ &= \langle \mathbf{a} - \mathbf{b}, \mathbf{a} - \mathbf{b} \rangle_{\mathcal{H}_k} \end{aligned}$$

$d_k(\mathbb{P}, \mathbb{Q})$ (=MMD) in terms of kernel evaluations

[Gretton et al., 2007]:

$$\begin{aligned} d_k^2(\mathbb{P}, \mathbb{Q}) &= \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k}^2 = \left\| \int_{\mathcal{X}} k(\cdot, x) d\mathbb{P}(x) - \int_{\mathcal{X}} k(\cdot, y) d\mathbb{Q}(y) \right\|_{\mathcal{H}_k}^2 \\ &= \langle \mathbf{a} - \mathbf{b}, \mathbf{a} - \mathbf{b} \rangle_{\mathcal{H}_k} \\ &= \langle \mu_{\mathbb{P}}, \mu_{\mathbb{P}} \rangle_{\mathcal{H}_k} + \langle \mu_{\mathbb{Q}}, \mu_{\mathbb{Q}} \rangle_{\mathcal{H}_k} - 2 \langle \mu_{\mathbb{P}}, \mu_{\mathbb{Q}} \rangle_{\mathcal{H}_k}, \end{aligned}$$

$d_k(\mathbb{P}, \mathbb{Q})$ (=MMD) in terms of kernel evaluations

[Gretton et al., 2007]:

$$\begin{aligned} d_k^2(\mathbb{P}, \mathbb{Q}) &= \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k}^2 = \left\| \int_{\mathcal{X}} k(\cdot, x) d\mathbb{P}(x) - \int_{\mathcal{X}} k(\cdot, y) d\mathbb{Q}(y) \right\|_{\mathcal{H}_k}^2 \\ &= \langle \mathbf{a} - \mathbf{b}, \mathbf{a} - \mathbf{b} \rangle_{\mathcal{H}_k} \\ &= \langle \mu_{\mathbb{P}}, \mu_{\mathbb{P}} \rangle_{\mathcal{H}_k} + \langle \mu_{\mathbb{Q}}, \mu_{\mathbb{Q}} \rangle_{\mathcal{H}_k} - 2 \langle \mu_{\mathbb{P}}, \mu_{\mathbb{Q}} \rangle_{\mathcal{H}_k}, \\ &= \int_{\mathcal{X}} \int_{\mathcal{X}} k(x, x') d\mathbb{P}(x) d\mathbb{P}(x') + \int_{\mathcal{X}} \int_{\mathcal{X}} k(y, y') d\mathbb{Q}(y) d\mathbb{Q}(y') \\ &\quad - 2 \int_{\mathcal{X}} \int_{\mathcal{X}} k(x, y) d\mathbb{P}(x) d\mathbb{Q}(y) \end{aligned}$$

$d_k(\mathbb{P}, \mathbb{Q})$ (=MMD) in terms of kernel evaluations

[Gretton et al., 2007]:

$$\begin{aligned} d_k^2(\mathbb{P}, \mathbb{Q}) &= \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k}^2 = \left\| \int_{\mathcal{X}} k(\cdot, x) d\mathbb{P}(x) - \int_{\mathcal{X}} k(\cdot, y) d\mathbb{Q}(y) \right\|_{\mathcal{H}_k}^2 \\ &= \langle \mathbf{a} - \mathbf{b}, \mathbf{a} - \mathbf{b} \rangle_{\mathcal{H}_k} \\ &= \langle \mu_{\mathbb{P}}, \mu_{\mathbb{P}} \rangle_{\mathcal{H}_k} + \langle \mu_{\mathbb{Q}}, \mu_{\mathbb{Q}} \rangle_{\mathcal{H}_k} - 2 \langle \mu_{\mathbb{P}}, \mu_{\mathbb{Q}} \rangle_{\mathcal{H}_k}, \\ &= \int_{\mathcal{X}} \int_{\mathcal{X}} k(x, x') d\mathbb{P}(x) d\mathbb{P}(x') + \int_{\mathcal{X}} \int_{\mathcal{X}} k(y, y') d\mathbb{Q}(y) d\mathbb{Q}(y') \\ &\quad - 2 \int_{\mathcal{X}} \int_{\mathcal{X}} k(x, y) d\mathbb{P}(x) d\mathbb{Q}(y) \\ &= \int_{\mathcal{X}} \int_{\mathcal{X}} k(x, y) d(\mathbb{P} - \mathbb{Q})(x) d(\mathbb{P} - \mathbb{Q})(y). \end{aligned}$$

⇒ Polynomial kernels are *not* characteristic

[Sriperumbudur et al., 2010b]:

- $k(x, y) = \langle x, y \rangle$: linear kernel ($L = 1$).

$$d_k^2(\mathbb{P}, \mathbb{Q}) = \|\textcolor{blue}{m}_{\mathbb{P}} - m_{\mathbb{Q}}\|_2^2, \quad \textcolor{blue}{m}_{\mathbb{P}} = \int_{\mathcal{X}} x d\mathbb{P}(x).$$

⇒ Polynomial kernels are *not* characteristic

[Sriperumbudur et al., 2010b]:

- $k(x, y) = \langle x, y \rangle$: linear kernel ($L = 1$).

$$d_k^2(\mathbb{P}, \mathbb{Q}) = \| \mathbf{m}_{\mathbb{P}} - \mathbf{m}_{\mathbb{Q}} \|_2^2, \quad \mathbf{m}_{\mathbb{P}} = \int_{\mathcal{X}} x d\mathbb{P}(x).$$

- $k(x, y) = (\langle x, y \rangle + 1)^2$ ($L = 2$):

$$d_k^2(\mathbb{P}, \mathbb{Q}) = 2 \| \mathbf{m}_{\mathbb{P}} - \mathbf{m}_{\mathbb{Q}} \|_2^2 + \| \Sigma_{\mathbb{P}} - \Sigma_{\mathbb{Q}} + \mathbf{m}_{\mathbb{P}} \mathbf{m}_{\mathbb{P}}^T - \mathbf{m}_{\mathbb{Q}} \mathbf{m}_{\mathbb{Q}}^T \|_F^2,$$

where $\|\cdot\|_F$: Frobenious norm; $\Sigma_{\mathbb{P}}$: cov. matrix w.r.t. \mathbb{P} .

Characteristic property

Well-understood for

- ➊ Continuous bounded shift-invariant kernels on \mathbb{R}^d :

$$k(x, y) = k_0(\textcolor{blue}{x} - \textcolor{blue}{y}), \quad k_0 \in C_b(\mathbb{R}^d).$$

Characteristic property

Well-understood for

- ① Continuous bounded shift-invariant kernels on \mathbb{R}^d :

$$k(x, y) = k_0(\textcolor{blue}{x} - \textcolor{blue}{y}), \quad k_0 \in C_b(\mathbb{R}^d).$$

- ② Continuous bounded radial kernels on \mathbb{R}^d :

$$k(x, y) = k_0(\|\textcolor{green}{x} - \textcolor{green}{y}\|_2), \quad k_0 \in C_b(\mathbb{R}^d),$$

$$k_0(z) = \int_{[0, \infty)} e^{-tz^2} d\nu(t)$$

$\nu \in \mathcal{M}_b^+[0, \infty)$, i.e. it is a finite measure on $[0, \infty)$.

Bochner's theorem

We focus on continuous bounded shift-invariant kernels:

Theorem (Bochner's theorem [Wendland, 2005], $k \leftrightarrow \Lambda$)

$$k_0(z) = \int_{\mathbb{R}^d} e^{-i\langle z, \omega \rangle} d\Lambda(\omega),$$

where Λ is a finite Borel measure (w.l.o.g. probability).

MMD in terms of characteristic functions

Using Bochner's theorem:

$$d_k^2(\mathbb{P}, \mathbb{Q}) = \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} k(x, y) d(\mathbb{P} - \mathbb{Q})(x) d(\mathbb{P} - \mathbb{Q})(y)$$

MMD in terms of characteristic functions

Using Bochner's theorem:

$$\begin{aligned} d_k^2(\mathbb{P}, \mathbb{Q}) &= \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} k(x, y) d(\mathbb{P} - \mathbb{Q})(x) d(\mathbb{P} - \mathbb{Q})(y) \\ &= \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} e^{-i\langle x-y, \omega \rangle} d\Lambda(\omega) d(\mathbb{P} - \mathbb{Q})(x) d(\mathbb{P} - \mathbb{Q})(y) \end{aligned}$$

MMD in terms of characteristic functions

Using Bochner's theorem:

$$\begin{aligned} d_k^2(\mathbb{P}, \mathbb{Q}) &= \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} k(x, y) d(\mathbb{P} - \mathbb{Q})(x) d(\mathbb{P} - \mathbb{Q})(y) \\ &= \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} e^{-i\langle x-y, \omega \rangle} d\Lambda(\omega) d(\mathbb{P} - \mathbb{Q})(x) d(\mathbb{P} - \mathbb{Q})(y) \\ &= \int_{\mathbb{R}^d} \underbrace{\left[\int_{\mathbb{R}^d} e^{-i\langle x, \omega \rangle} d(\mathbb{P} - \mathbb{Q})(x) \right]}_{\overline{\phi_{\mathbb{P}}(\omega) - \phi_{\mathbb{Q}}(\omega)}} \underbrace{\left[\int_{\mathbb{R}^d} e^{i\langle y, \omega \rangle} d(\mathbb{P} - \mathbb{Q})(y) \right]}_{\phi_{\mathbb{P}}(\omega) - \phi_{\mathbb{Q}}(\omega)} d\Lambda(\omega) \end{aligned}$$

MMD in terms of characteristic functions

Using Bochner's theorem:

$$\begin{aligned} d_k^2(\mathbb{P}, \mathbb{Q}) &= \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} k(x, y) d(\mathbb{P} - \mathbb{Q})(x) d(\mathbb{P} - \mathbb{Q})(y) \\ &= \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} e^{-i\langle x-y, \omega \rangle} d\Lambda(\omega) d(\mathbb{P} - \mathbb{Q})(x) d(\mathbb{P} - \mathbb{Q})(y) \\ &= \int_{\mathbb{R}^d} \underbrace{\left[\int_{\mathbb{R}^d} e^{-i\langle x, \omega \rangle} d(\mathbb{P} - \mathbb{Q})(x) \right]}_{\overline{\phi_{\mathbb{P}}(\omega) - \phi_{\mathbb{Q}}(\omega)}} \underbrace{\left[\int_{\mathbb{R}^d} e^{i\langle y, \omega \rangle} d(\mathbb{P} - \mathbb{Q})(y) \right]}_{\phi_{\mathbb{P}}(\omega) - \phi_{\mathbb{Q}}(\omega)} d\Lambda(\omega) \\ &= \int_{\mathbb{R}^d} |\phi_{\mathbb{P}}(\omega) - \phi_{\mathbb{Q}}(\omega)|^2 d\Lambda(\omega) \end{aligned}$$

MMD in terms of characteristic functions

Using Bochner's theorem:

$$\begin{aligned} d_k^2(\mathbb{P}, \mathbb{Q}) &= \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} k(x, y) d(\mathbb{P} - \mathbb{Q})(x) d(\mathbb{P} - \mathbb{Q})(y) \\ &= \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} e^{-i\langle x-y, \omega \rangle} d\Lambda(\omega) d(\mathbb{P} - \mathbb{Q})(x) d(\mathbb{P} - \mathbb{Q})(y) \\ &= \int_{\mathbb{R}^d} \underbrace{\left[\int_{\mathbb{R}^d} e^{-i\langle x, \omega \rangle} d(\mathbb{P} - \mathbb{Q})(x) \right]}_{\overline{\phi_{\mathbb{P}}(\omega) - \phi_{\mathbb{Q}}(\omega)}} \underbrace{\left[\int_{\mathbb{R}^d} e^{i\langle y, \omega \rangle} d(\mathbb{P} - \mathbb{Q})(y) \right]}_{\phi_{\mathbb{P}}(\omega) - \phi_{\mathbb{Q}}(\omega)} d\Lambda(\omega) \\ &= \int_{\mathbb{R}^d} |\phi_{\mathbb{P}}(\omega) - \phi_{\mathbb{Q}}(\omega)|^2 d\Lambda(\omega) = \|\phi_{\mathbb{P}} - \phi_{\mathbb{Q}}\|_{L^2(\Lambda)}^2. \end{aligned}$$

Theorem ([Sriperumbudur et al., 2010b])

They are characteristic iff. $\text{supp}(\Lambda) = \mathbb{R}^d$.

Theorem ([Sriperumbudur et al., 2010b])

They are characteristic iff. $\text{supp}(\Lambda) = \mathbb{R}^d$.

$\text{supp}(\Lambda) := \{x \in \mathcal{X}: \text{for any open set } U \text{ such that } x \in U, \Lambda(U) > 0\}$.

Theorem ([Sriperumbudur et al., 2010b])

They are characteristic iff. $\text{supp}(\Lambda) = \mathbb{R}^d$.

$\text{supp}(\Lambda) := \{x \in \mathcal{X}: \text{for any open set } U \text{ such that } x \in U, \Lambda(U) > 0\}$.

- **Example:** Gaussian, Laplacian, Matérn kernel, B-spline kernel.

Theorem ([Sriperumbudur et al., 2010b])

They are characteristic iff. $\text{supp}(\Lambda) = \mathbb{R}^d$.

$\text{supp}(\Lambda) := \{x \in \mathcal{X}: \text{for any open set } U \text{ such that } x \in U, \Lambda(U) > 0\}$.

- **Example:** Gaussian, Laplacian, Matérn kernel, B-spline kernel.
- Similar characterization \exists on '**Bochner domains**' (LCA groups [Berg et al., 1984], orthogonal matrices, \mathbb{R}_+^d)
[Fukumizu et al., 2009b].

Matérn kernel

$$k(x, y) = k_0(x - y) = \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu} \|x - y\|_2}{\sigma} \right)^{\nu} K_{\nu} \left(\frac{\sqrt{2\nu} \|x - y\|_2}{\sigma} \right),$$

where K_ν : modified Bessel function of the second kind of order ν

Matérn kernel

$$k(x, y) = k_0(x - y) = \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu} \|x - y\|_2}{\sigma} \right)^{\nu} K_{\nu} \left(\frac{\sqrt{2\nu} \|x - y\|_2}{\sigma} \right),$$
$$\hat{k}_0(\omega) = \frac{2^{d+\nu} \pi^{\frac{d}{2}} \Gamma(\nu + d/2) \nu^{\nu}}{\Gamma(\nu) \sigma^{2\nu}} \left(\frac{2\nu}{\sigma^2} + 4\pi^2 \|\omega\|_2^2 \right)^{-(\nu+d/2)} > 0 \quad \forall \omega \in \mathbb{R}^d,$$

where K_{ν} : modified Bessel function of the second kind of order ν , Γ : Gamma function.

Matérn kernel

$$k(x, y) = k_0(x - y) = \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu} \|x - y\|_2}{\sigma} \right)^{\nu} K_{\nu} \left(\frac{\sqrt{2\nu} \|x - y\|_2}{\sigma} \right),$$
$$\hat{k}_0(\omega) = \frac{2^{d+\nu} \pi^{\frac{d}{2}} \Gamma(\nu + d/2) \nu^{\nu}}{\Gamma(\nu) \sigma^{2\nu}} \left(\frac{2\nu}{\sigma^2} + 4\pi^2 \|\omega\|_2^2 \right)^{-(\nu+d/2)} > 0 \quad \forall \omega \in \mathbb{R}^d,$$

where K_{ν} : modified Bessel function of the second kind of order ν , Γ : Gamma function.

- For $\nu = \frac{1}{2}$: one gets $k(x, y) = e^{-\frac{\|x-y\|_2}{\sigma}}$.

Matérn kernel

$$k(x, y) = k_0(x - y) = \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu} \|x - y\|_2}{\sigma} \right)^{\nu} K_{\nu} \left(\frac{\sqrt{2\nu} \|x - y\|_2}{\sigma} \right),$$
$$\hat{k}_0(\omega) = \frac{2^{d+\nu} \pi^{\frac{d}{2}} \Gamma(\nu + d/2) \nu^{\nu}}{\Gamma(\nu) \sigma^{2\nu}} \left(\frac{2\nu}{\sigma^2} + 4\pi^2 \|\omega\|_2^2 \right)^{-(\nu+d/2)} > 0 \quad \forall \omega \in \mathbb{R}^d,$$

where K_{ν} : modified Bessel function of the second kind of order ν , Γ : Gamma function.

- For $\nu = \frac{1}{2}$: one gets $k(x, y) = e^{-\frac{\|x-y\|_2}{\sigma}}$.
- Gaussian kernel: $\nu \rightarrow \infty$.

Shift-invariant kernels on \mathbb{R} [Sriperumbudur et al., 2010b]

For Poisson kernel: $\sigma \in (0, 1)$.

kernel name k_0	$\hat{k}_0(\omega)$	$supp(\hat{k}_0)$
Gaussian	$e^{-\frac{x^2}{2\sigma^2}}$	\mathbb{R}
Laplacian	$e^{-\sigma x }$	\mathbb{R}
B_{2n+1} -spline	$*^{2n+2} \chi_{[-\frac{1}{2}, \frac{1}{2}]}(x) \frac{4^{n+1}}{\sqrt{2\pi}} \frac{\sin^{2n+2}(\frac{\omega}{2})}{\omega^{2n+2}}$	\mathbb{R}
Sinc	$\frac{\sin(\sigma x)}{x}$	$[-\sigma, \sigma]$
Poisson	$\frac{1-\sigma^2}{\sigma^2 - 2\sigma \cos(x) + 1}$	\mathbb{Z}
Dirichlet	$\frac{\sin(\frac{(2n+1)x}{2})}{\sin(\frac{x}{2})}$	$\{0, \pm 1, \pm 2, \dots, \pm n\}$
Fejér	$\frac{1}{n+1} \frac{\sin^2(\frac{n+1}{2}x)}{\sin^2(\frac{x}{2})}$	$\{0, \pm 1, \pm 2, \dots, \pm n\}$
Cosine	$\cos(\sigma x)$	$\{-\sigma, \sigma\}$

Shift-invariant kernels on \mathbb{R} [Sriperumbudur et al., 2010b]

For Poisson kernel: $\sigma \in (0, 1)$.

kernel name k_0	$\hat{k}_0(\omega)$	$supp(\hat{k}_0)$
Gaussian	$e^{-\frac{x^2}{2\sigma^2}}$	\mathbb{R}
Laplacian	$e^{-\sigma x }$	\mathbb{R}
B_{2n+1} -spline	$*^{2n+2} \chi_{[-\frac{1}{2}, \frac{1}{2}]}(x) \frac{4^{n+1}}{\sqrt{2\pi}} \frac{\sin^{2n+2}(\frac{\omega}{2})}{\omega^{2n+2}}$	\mathbb{R}
Sinc	$\frac{\sin(\sigma x)}{x}$	$[-\sigma, \sigma]$
Poisson	$\frac{1-\sigma^2}{\sigma^2 - 2\sigma \cos(x) + 1}$	\mathbb{Z}
Dirichlet	$\frac{\sin(\frac{(2n+1)x}{2})}{\sin(\frac{x}{2})}$	$\{0, \pm 1, \pm 2, \dots, \pm n\}$
Fejér	$\frac{1}{n+1} \frac{\sin^2(\frac{n+1}{2}x)}{\sin^2(\frac{x}{2})}$	$\{0, \pm 1, \pm 2, \dots, \pm n\}$
Cosine	$\cos(\sigma x)$	$\{-\sigma, \sigma\}$

For $x \in \mathbb{R}^d$: $k_0(x) = \prod_{j=1}^d k_0(x_j)$, $\hat{k}_0(\omega) = \prod_{j=1}^d \hat{k}_0(\omega_j)$.

B-spline kernel type kernels

- Still k : continuous, bounded, shift-invariant.
- **B-spline kernel**: $\text{supp}(k_0)$ is compact \leftarrow practically relevant.
- Note: $\text{supp}(f) := \overline{\{x \in \mathbb{R}^d : f(x) \neq 0\}}$.

B-spline kernel type kernels

- Still k : continuous, bounded, shift-invariant.
- B-spline kernel: $\text{supp}(k_0)$ is compact \leftarrow practically relevant.
- Note: $\text{supp}(f) := \overline{\{x \in \mathbb{R}^d : f(x) \neq 0\}}$.
- More generally

Theorem ([Sriperumbudur et al., 2010b])

$\text{supp}(k_0)$: compact $\Rightarrow k$ is characteristic.

Construction of new characteristic kernels: +, \times

Theorem ([Sriperumbudur et al., 2010b])

If k, k_1, k_2 : continuous, bounded, shift-invariant; k : characteristic, $k_2 \neq 0$. Then $k + k_1, kk_2$ is also characteristic.

Construction of new characteristic kernels: +, ×

Theorem ([Sriperumbudur et al., 2010b])

If k, k_1, k_2 : continuous, bounded, shift-invariant; k : characteristic, $k_2 \neq 0$. Then $k + k_1, kk_2$ is also characteristic.

Proof.

We focus on $k + k_1$ (product: similarly):

$$\begin{aligned}(k + k_1)(x, y) &:= k(x, y) + k_1(x, y) = k_0(x - y) + (k_1)_0(x - y) \\ &= \int_{\mathbb{R}^d} e^{-i\langle x-y, \omega \rangle} d(\Lambda + \Lambda_1)(\omega).\end{aligned}$$



Construction of new characteristic kernels: +, \times

Theorem ([Sriperumbudur et al., 2010b])

If k, k_1, k_2 : continuous, bounded, shift-invariant; k : characteristic, $k_2 \neq 0$. Then $k + k_1, kk_2$ is also characteristic.

Proof.

We focus on $k + k_1$ (product: similarly):

$$\begin{aligned}(k + k_1)(x, y) &:= k(x, y) + k_1(x, y) = k_0(x - y) + (k_1)_0(x - y) \\ &= \int_{\mathbb{R}^d} e^{-i\langle x-y, \omega \rangle} d(\Lambda + \Lambda_1)(\omega).\end{aligned}$$

- k : characteristic $\Rightarrow \text{supp}(\Lambda) = \mathbb{R}^d$.



Construction of new characteristic kernels: $+$, \times

Theorem ([Sriperumbudur et al., 2010b])

If k, k_1, k_2 : continuous, bounded, shift-invariant; k : characteristic, $k_2 \neq 0$. Then $k + k_1$, kk_2 is also characteristic.

Proof.

We focus on $k + k_1$ (product: similarly):

$$\begin{aligned}(k + k_1)(x, y) &:= k(x, y) + k_1(x, y) = k_0(x - y) + (k_1)_0(x - y) \\&= \int_{\mathbb{R}^d} e^{-i\langle x-y, \omega \rangle} d(\Lambda + \Lambda_1)(\omega).\end{aligned}$$

- k : characteristic $\Rightarrow \text{supp}(\Lambda) = \mathbb{R}^d$.
- Since $\text{supp}(\Lambda) \subseteq \text{supp}(\Lambda + \Lambda_1)$, we get $\text{supp}(\Lambda + \Lambda_1) = \mathbb{R}^d$; hence $k + k_1$ is characteristic.



Recall (radial kernel):

$$k(x, y) = k_0(\|x - y\|_2), \quad k_0(z) = \int_{[0, \infty)} e^{-tz^2} d\nu(t).$$

Recall (radial kernel):

$$k(x, y) = k_0(\|x - y\|_2), \quad k_0(z) = \int_{[0, \infty)} e^{-tz^2} d\nu(t).$$

Theorem ([Sriperumbudur et al., 2010b])

k is characteristic iff. $\text{supp}(\nu) \neq \{0\}$.

More general spaces

- $\mathcal{M}_b(\mathcal{X})$: set of all finite signed (Radon) measures on \mathcal{X}
(Radon $\Rightarrow \exists supp$).

More general spaces

- $\mathcal{M}_b(\mathcal{X})$: set of all finite signed (Radon) measures on \mathcal{X} (Radon $\Rightarrow \exists supp$).
- Ulam's Theorem [Dudley, 2004]: On an \mathcal{X} Polish space \forall Borel measure is Radon.

More general spaces

- $\mathcal{M}_b(\mathcal{X})$: set of all finite signed (Radon) measures on \mathcal{X} (Radon $\Rightarrow \exists supp$).
- Ulam's Theorem [Dudley, 2004]: On an \mathcal{X} Polish space \forall Borel measure is Radon.

Definition

A $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ bounded, measurable kernel is called *integrally strictly positive definite (ispd)* if

$$\int_{\mathcal{X}} \int_{\mathcal{X}} k(x, y) d\mathbb{F}(x)\mathbb{F}(y) > 0 \quad \forall 0 \neq \mathbb{F} \in \mathcal{M}_b(\mathcal{X}).$$

Sufficient condition: ispd

Theorem ([Sriperumbudur et al., 2010b])

Ispld kernels are characteristic on an \mathcal{X} topological space.

- **ispd on \mathbb{R}^d :** Gaussian, Laplacian, inverse multiquadratics, Matérn kernels, B-splines.

Sufficient condition: ispd

Theorem ([Sriperumbudur et al., 2010b])

Ispd kernels are characteristic on an \mathcal{X} topological space.

- ispd on \mathbb{R}^d : Gaussian, Laplacian, inverse multiquadratics, Matérn kernels, B-splines.
- Dirichlet kernel: characteristic, though not ispd.

Sufficient condition: ispd

Theorem ([Sriperumbudur et al., 2010b])

Ispd kernels are characteristic on an \mathcal{X} topological space.

- ispd on \mathbb{R}^d : Gaussian, Laplacian, inverse multiquadratics, Matérn kernels, B-splines.
- Dirichlet kernel: characteristic, though not ispd.
- ispd property: checking might not be easy.

Shift-variant ispd from shift-invariant ispd kernel:

$$k_0(x, y) = f(x)k(x, y)f(y), \quad f \in C_b(\mathcal{X}).$$

Shift-variant ispd from shift-invariant ispd kernel:

$$k_0(x, y) = f(x)k(x, y)f(y), \quad f \in C_b(\mathcal{X}).$$

Example (exponential \leftarrow Gaussian): $k_0(x, y) = e^{\sigma\langle x, y \rangle}$, $\mathcal{X} \subset \mathbb{R}^d$ compact

$$k(x, y) = e^{-\sigma \frac{\|x-y\|^2}{2}}, \quad f(x) = e^{\sigma \frac{\|x\|^2}{2}}.$$

Theorem ([Fukumizu et al., 2008, Fukumizu et al., 2009a])

Let $r \geq 1$.

- Sufficient condition: A $k : (\mathcal{X}, \mathcal{A}) \times (\mathcal{X}, \mathcal{A}) \rightarrow \mathbb{R}$ bounded measurable kernel is characteristic if $\mathcal{H}_k + \mathbb{R}$ is dense in $L^r(\mathcal{X}, \mathcal{A}, \mathbb{P})$ for all $\mathbb{P} \in \mathcal{M}_1^+(\mathcal{X})$.

Denseness in L^r

Theorem ([Fukumizu et al., 2008, Fukumizu et al., 2009a])

Let $r \geq 1$.

- *Sufficient condition: A $k : (\mathcal{X}, \mathcal{A}) \times (\mathcal{X}, \mathcal{A}) \rightarrow \mathbb{R}$ bounded measurable kernel is characteristic if $\mathcal{H}_k + \mathbb{R}$ is dense in $L^r(\mathcal{X}, \mathcal{A}, \mathbb{P})$ for all $\mathbb{P} \in \mathcal{M}_1^+(\mathcal{X})$.*
- *With $r = 2$, it is also a necessary condition.*

Denseness in L^r

Theorem ([Fukumizu et al., 2008, Fukumizu et al., 2009a])

Let $r \geq 1$.

- *Sufficient condition: A $k : (\mathcal{X}, \mathcal{A}) \times (\mathcal{X}, \mathcal{A}) \rightarrow \mathbb{R}$ bounded measurable kernel is characteristic if $\mathcal{H}_k + \mathbb{R}$ is dense in $L^r(\mathcal{X}, \mathcal{A}, \mathbb{P})$ for all $\mathbb{P} \in \mathcal{M}_1^+(\mathcal{X})$.*
- *With $r = 2$, it is also a necessary condition.*

Note:

- For a **c-universal kernel** k : sufficient condition holds with $r = 2$.
- This gives the **3rd 'universal \Rightarrow characteristic' proof**.

Denseness in L^r

Theorem ([Fukumizu et al., 2008, Fukumizu et al., 2009a])

Let $r \geq 1$.

- *Sufficient condition: A $k : (\mathcal{X}, \mathcal{A}) \times (\mathcal{X}, \mathcal{A}) \rightarrow \mathbb{R}$ bounded measurable kernel is characteristic if $\mathcal{H}_k + \mathbb{R}$ is dense in $L^r(\mathcal{X}, \mathcal{A}, \mathbb{P})$ for all $\mathbb{P} \in \mathcal{M}_1^+(\mathcal{X})$.*
- *With $r = 2$, it is also a necessary condition.*

Note:

- For a **c-universal kernel** k : sufficient condition holds with $r = 2$.
- This gives the **3rd 'universal \Rightarrow characteristic' proof**.

Let us prove this theorem...

Denseness is sufficient: idea

- Goal: in this case, $\mu_{\mathbb{P}} = \mu_{\mathbb{Q}} \Rightarrow \mathbb{P}(A) = \mathbb{Q}(A)$ for any $A \in \mathcal{A}$.

Denseness is sufficient: idea

- Goal: in this case, $\mu_{\mathbb{P}} = \mu_{\mathbb{Q}} \Rightarrow \mathbb{P}(A) = \mathbb{Q}(A)$ for any $A \in \mathcal{A}$.
- Enough: $|\mathbb{P}(A) - \mathbb{Q}(A)| = |\mathbb{P}\chi_A - \mathbb{Q}\chi_A| \leq \epsilon$, $\forall A \in \mathcal{A}$, $\forall \epsilon > 0$.

Denseness is sufficient: idea

- Enough: $|\mathbb{P}(A) - \mathbb{Q}(A)| = |\mathbb{P}\chi_A - \mathbb{Q}\chi_A| \leq \epsilon$, $\forall A \in \mathcal{A}$, $\forall \epsilon > 0$.
- Idea:
 - 1 using the max. difference of \mathbb{P} and $\mathbb{Q} \Rightarrow \text{TV}$ of $\mathbb{P} - \mathbb{Q}$,

$$|\mathbb{P} - \mathbb{Q}|(\mathcal{X}) = 2 \sup_{A \in \mathcal{A}} |\mathbb{P}(A) - \mathbb{Q}(A)|.$$

Densemness is sufficient: idea

- Enough: $|\mathbb{P}(A) - \mathbb{Q}(A)| = |\mathbb{P}\chi_A - \mathbb{Q}\chi_A| \leq \epsilon, \forall A \in \mathcal{A}, \forall \epsilon > 0.$
- Idea:
 - using the max. difference of \mathbb{P} and $\mathbb{Q} \Rightarrow \text{TV}$ of $\mathbb{P} - \mathbb{Q}$,

$$|\mathbb{P} - \mathbb{Q}|(\mathcal{X}) = 2 \sup_{A \in \mathcal{A}} |\mathbb{P}(A) - \mathbb{Q}(A)|.$$

- exploit denseness for $\chi_A \in \underbrace{L^r(\mathcal{X}, \mathcal{A}, |\mathbb{P} - \mathbb{Q}|)}_{=: L^r(|\mathbb{P} - \mathbb{Q}|)}$.

Total variation: quick summary

Idea: $f = f^+ - f^- \rightarrow |f| = f^+ + f^-.$

Total variation: quick summary

Idea: $f = f^+ - f^- \rightarrow |f| = f^+ + f^-$. Analogously:

- $(\mathcal{X}, \mathcal{A})$: measurable space. μ : signed measure on it.

Total variation: quick summary

Idea: $f = f^+ - f^- \rightarrow |f| = f^+ + f^-$. Analogously:

- $(\mathcal{X}, \mathcal{A})$: measurable space. μ : signed measure on it.
- Hahn-Jordan decomposition of μ : $\mathcal{X} = \mathcal{P} \dot{\cup} \mathcal{N}$.

Total variation: quick summary

Idea: $f = f^+ - f^- \rightarrow |f| = f^+ + f^-$. Analogously:

- $(\mathcal{X}, \mathcal{A})$: measurable space. μ : signed measure on it.
- Hahn-Jordan decomposition of μ : $\mathcal{X} = \mathcal{P} \dot{\cup} \mathcal{N}$.
- Positive & negative part of μ ($= \mu^+ - \mu^-$):

$$\mu^+(A) = \mu(A \cap \mathcal{P}), \quad \mu^-(A) = \mu(A \cap \mathcal{N}) \quad \forall A \in \mathcal{A}.$$

Total variation: quick summary

Idea: $f = f^+ - f^- \rightarrow |f| = f^+ + f^-$. Analogously:

- $(\mathcal{X}, \mathcal{A})$: measurable space. μ : signed measure on it.
- Hahn-Jordan decomposition of μ : $\mathcal{X} = \mathcal{P} \dot{\cup} \mathcal{N}$.
- Positive & negative part of μ ($= \mu^+ - \mu^-$):

$$\mu^+(A) = \mu(A \cap \mathcal{P}), \quad \mu^-(A) = \mu(A \cap \mathcal{N}) \quad \forall A \in \mathcal{A}.$$

- TV of μ : $|\mu| := \mu^+ + \mu^-$.

Total variation: quick summary

Idea: $f = f^+ - f^- \rightarrow |f| = f^+ + f^-$. Analogously:

- $(\mathcal{X}, \mathcal{A})$: measurable space. μ : signed measure on it.
- Hahn-Jordan decomposition of μ : $\mathcal{X} = \mathcal{P} \dot{\cup} \mathcal{N}$.
- Positive & negative part of μ ($= \mu^+ - \mu^-$):

$$\mu^+(A) = \mu(A \cap \mathcal{P}), \quad \mu^-(A) = \mu(A \cap \mathcal{N}) \quad \forall A \in \mathcal{A}.$$

- TV of μ : $|\mu| := \mu^+ + \mu^-$.
- μ : finite $\Rightarrow \mu^+$, μ^- : finite.

Denseness is sufficient: proof

- Take: $A \in \mathcal{A}$, $\epsilon > 0$.

Denseness is sufficient: proof

- Take: $A \in \mathcal{A}$, $\epsilon > 0$.
- $\mathcal{H}_k + \mathbb{R}$ is dense in $L^r(|\mathbb{P} - \mathbb{Q}|) \Rightarrow \exists f \in \mathcal{H}_k + \mathbb{R}$

$$\|f - \chi_A\|_{L^r(|\mathbb{P} - \mathbb{Q}|)} \leq \epsilon.$$

Denseness is sufficient: proof

- Take: $A \in \mathcal{A}$, $\epsilon > 0$.
- $\mathcal{H}_k + \mathbb{R}$ is dense in $L^r(|\mathbb{P} - \mathbb{Q}|) \Rightarrow \exists f \in \mathcal{H}_k + \mathbb{R}$

$$\|f - \chi_A\|_{L^r(|\mathbb{P} - \mathbb{Q}|)} \leq \epsilon.$$

- Some lower bounding

$$\epsilon \geq \|f - \chi_A\|_{L^r(|\mathbb{P} - \mathbb{Q}|)} \stackrel{r \geq 1}{\gtrsim} \|\underbrace{f - \chi_A}_{=: g}\|_{L^1(|\mathbb{P} - \mathbb{Q}|)}$$

Denseness is sufficient: proof

- Take: $A \in \mathcal{A}$, $\epsilon > 0$.
- $\mathcal{H}_k + \mathbb{R}$ is dense in $L^r(|\mathbb{P} - \mathbb{Q}|) \Rightarrow \exists f \in \mathcal{H}_k + \mathbb{R}$

$$\|f - \chi_A\|_{L^r(|\mathbb{P} - \mathbb{Q}|)} \leq \epsilon.$$

- Some lower bounding

$$\epsilon \geq \|f - \chi_A\|_{L^r(|\mathbb{P} - \mathbb{Q}|)} \stackrel{r \geq 1}{\gtrsim} \|\underbrace{f - \chi_A}_{=: g}\|_{L^1(|\mathbb{P} - \mathbb{Q}|)} = |\mathbb{P} - \mathbb{Q}|(|g|)$$

Denseness is sufficient: proof

- Take: $A \in \mathcal{A}$, $\epsilon > 0$.
- $\mathcal{H}_k + \mathbb{R}$ is dense in $L^r(|\mathbb{P} - \mathbb{Q}|) \Rightarrow \exists f \in \mathcal{H}_k + \mathbb{R}$

$$\|f - \chi_A\|_{L^r(|\mathbb{P} - \mathbb{Q}|)} \leq \epsilon.$$

- Some lower bounding

$$\begin{aligned}\epsilon &\geq \|f - \chi_A\|_{L^r(|\mathbb{P} - \mathbb{Q}|)} \stackrel{r \geq 1}{\gtrsim} \|\underbrace{f - \chi_A}_{=: g}\|_{L^1(|\mathbb{P} - \mathbb{Q}|)} = |\mathbb{P} - \mathbb{Q}|(|g|) \\ &\geq |\mathbb{P} - \mathbb{Q}|(g)\end{aligned}$$

Denseness is sufficient: proof

- Take: $A \in \mathcal{A}$, $\epsilon > 0$.
- $\mathcal{H}_k + \mathbb{R}$ is dense in $L^r(|\mathbb{P} - \mathbb{Q}|) \Rightarrow \exists f \in \mathcal{H}_k + \mathbb{R}$

$$\|f - \chi_A\|_{L^r(|\mathbb{P} - \mathbb{Q}|)} \leq \epsilon.$$

- Some lower bounding

$$\begin{aligned}\epsilon &\geq \|f - \chi_A\|_{L^r(|\mathbb{P} - \mathbb{Q}|)} \stackrel{r \geq 1}{\gtrsim} \|\underbrace{f - \chi_A}_{=: g}\|_{L^1(|\mathbb{P} - \mathbb{Q}|)} = |\mathbb{P} - \mathbb{Q}|(|g|) \\ &\geq |\mathbb{P} - \mathbb{Q}|(g) \geq |(\mathbb{P} - \mathbb{Q})(g)|\end{aligned}$$

Denseness is sufficient: proof

- Take: $A \in \mathcal{A}$, $\epsilon > 0$.
- $\mathcal{H}_k + \mathbb{R}$ is dense in $L^r(|\mathbb{P} - \mathbb{Q}|) \Rightarrow \exists f \in \mathcal{H}_k + \mathbb{R}$

$$\|f - \chi_A\|_{L^r(|\mathbb{P} - \mathbb{Q}|)} \leq \epsilon.$$

- Some lower bounding

$$\begin{aligned}\epsilon &\geq \|f - \chi_A\|_{L^r(|\mathbb{P} - \mathbb{Q}|)} \stackrel{r \geq 1}{\gtrsim} \|\underbrace{f - \chi_A}_{=: g}\|_{L^1(|\mathbb{P} - \mathbb{Q}|)} = |\mathbb{P} - \mathbb{Q}|(|g|) \\ &\geq |\mathbb{P} - \mathbb{Q}|(g) \geq |(\mathbb{P} - \mathbb{Q})(g)| = |\mathbb{P}(f - \chi_A) - \mathbb{Q}(f - \chi_A)|\end{aligned}$$

Denseness is sufficient: proof

- Take: $A \in \mathcal{A}$, $\epsilon > 0$.
- $\mathcal{H}_k + \mathbb{R}$ is dense in $L^r(|\mathbb{P} - \mathbb{Q}|) \Rightarrow \exists f \in \mathcal{H}_k + \mathbb{R}$

$$\|f - \chi_A\|_{L^r(|\mathbb{P} - \mathbb{Q}|)} \leq \epsilon.$$

- Some lower bounding

$$\begin{aligned}\epsilon &\geq \|f - \chi_A\|_{L^r(|\mathbb{P} - \mathbb{Q}|)} \stackrel{r \geq 1}{\gtrsim} \|\underbrace{f - \chi_A}_{=: g}\|_{L^1(|\mathbb{P} - \mathbb{Q}|)} = |\mathbb{P} - \mathbb{Q}|(|g|) \\ &\geq |\mathbb{P} - \mathbb{Q}|(g) \geq |(\mathbb{P} - \mathbb{Q})(g)| = |\mathbb{P}(f - \chi_A) - \mathbb{Q}(f - \chi_A)| \\ &\stackrel{(*)}{=} |\mathbb{P}\chi_A - \mathbb{Q}\chi_A|.\end{aligned}$$

(*): $\mathbb{P}f = \mathbb{Q}f$ for any $f \in \mathcal{H}_k$ since $\mu_{\mathbb{P}} = \mu_{\mathbb{Q}}$.

Denseness in L^2 is necessary: proof

If $\mathcal{H}_k + \mathbb{R}$ is *not* dense in $L^2(\mathbb{P}) := L^2(\mathcal{X}, \mathcal{A}, \mathbb{P})$, then

- goal: $\underbrace{\exists \mathbb{Q}_1 \neq \mathbb{Q}_2 \in \mathcal{M}_1^+(\mathcal{X}) \text{ s.t. } \mu_{\mathbb{Q}_1} = \mu_{\mathbb{Q}_2}}_{\mu \text{ is not injective}}$.

Denseness in L^2 is necessary: proof

If $\mathcal{H}_k + \mathbb{R}$ is *not* dense in $L^2(\mathbb{P}) := L^2(\mathcal{X}, \mathcal{A}, \mathbb{P})$, then

- goal: $\underbrace{\exists \mathbb{Q}_1 \neq \mathbb{Q}_2 \in \mathcal{M}_1^+(\mathcal{X}) \text{ s.t. } \mu_{\mathbb{Q}_1} = \mu_{\mathbb{Q}_2}}_{\mu \text{ is not injective}}$.
- Hahn-Banach: $0 \neq f \in L^2(\mathbb{P})$ s.t. $f \perp \mathbf{1}_{\mathcal{H}_k}$, thus

$$\langle f, \mathbf{1} \rangle_{L^2(\mathbb{P})} = 0, \quad \langle f, h \rangle_{L^2(\mathbb{P})} = 0 \quad (\forall h \in \mathcal{H}_k).$$

Denseness in L^2 is necessary: proof

If $\mathcal{H}_k + \mathbb{R}$ is *not* dense in $L^2(\mathbb{P}) := L^2(\mathcal{X}, \mathcal{A}, \mathbb{P})$, then

- goal: $\underbrace{\exists \mathbb{Q}_1 \neq \mathbb{Q}_2 \in \mathcal{M}_1^+(\mathcal{X}) \text{ s.t. } \mu_{\mathbb{Q}_1} = \mu_{\mathbb{Q}_2}}_{\mu \text{ is not injective}}$.
- Hahn-Banach: $0 \neq f \in L^2(\mathbb{P})$ s.t. $f \perp 1, \mathcal{H}_k$, thus

$$\langle f, 1 \rangle_{L^2(\mathbb{P})} = 0, \quad \langle f, h \rangle_{L^2(\mathbb{P})} = 0 \quad (\forall h \in \mathcal{H}_k).$$

- We define $\mathbb{Q}_1, \mathbb{Q}_2 \in \mathcal{M}_1^+(\mathcal{X})$ from f ($f \neq 0 \Rightarrow \mathbb{Q}_1 \neq \mathbb{Q}_2$):

$$\mathbb{Q}_1(A) = c \int_A |f| d\mathbb{P}, \quad \mathbb{Q}_2(A) = c \int_A (\underbrace{|f| - f}_{\geq 0}) d\mathbb{P}, \quad c = \frac{1}{\int_{\mathcal{X}} |f| d\mathbb{P}}.$$

Denseness in L^2 is necessary: proof continued

We arrive at

$$\mu_{\mathbb{Q}_1} - \mu_{\mathbb{Q}_2} = \int_{\mathcal{X}} k(\cdot, x) d\mathbb{Q}_1(x) - \int_{\mathcal{X}} k(\cdot, x) d\mathbb{Q}_2(x)$$

Denseness in L^2 is necessary: proof continued

We arrive at

$$\begin{aligned}\mu_{\mathbb{Q}_1} - \mu_{\mathbb{Q}_2} &= \int_{\mathcal{X}} k(\cdot, x) d\mathbb{Q}_1(x) - \int_{\mathcal{X}} k(\cdot, x) d\mathbb{Q}_2(x) \\ &= \int_{\mathcal{X}} k(\cdot, x) d(\mathbb{Q}_1 - \mathbb{Q}_2)(x)\end{aligned}$$

Denseness in L^2 is necessary: proof continued

We arrive at

$$\begin{aligned}\mu_{\mathbb{Q}_1} - \mu_{\mathbb{Q}_2} &= \int_{\mathcal{X}} k(\cdot, x) d\mathbb{Q}_1(x) - \int_{\mathcal{X}} k(\cdot, x) d\mathbb{Q}_2(x) \\ &= \int_{\mathcal{X}} k(\cdot, x) d(\mathbb{Q}_1 - \mathbb{Q}_2)(x) = c \int_{\mathcal{X}} f(x) k(\cdot, x) d\mathbb{P}(x),\end{aligned}$$

Denseness in L^2 is necessary: proof continued

We arrive at

$$\begin{aligned}\mu_{\mathbb{Q}_1} - \mu_{\mathbb{Q}_2} &= \int_{\mathcal{X}} k(\cdot, x) d\mathbb{Q}_1(x) - \int_{\mathcal{X}} k(\cdot, x) d\mathbb{Q}_2(x) \\ &= \int_{\mathcal{X}} k(\cdot, x) d(\mathbb{Q}_1 - \mathbb{Q}_2)(x) = c \int_{\mathcal{X}} f(x) k(\cdot, x) d\mathbb{P}(x), \\ (\mu_{\mathbb{Q}_1} - \mu_{\mathbb{Q}_2})(y) &= c \int_{\mathcal{X}} f(x) k(y, x) d\mathbb{P}(x)\end{aligned}$$

Denseness in L^2 is necessary: proof continued

We arrive at

$$\begin{aligned}\mu_{\mathbb{Q}_1} - \mu_{\mathbb{Q}_2} &= \int_{\mathcal{X}} k(\cdot, x) d\mathbb{Q}_1(x) - \int_{\mathcal{X}} k(\cdot, x) d\mathbb{Q}_2(x) \\ &= \int_{\mathcal{X}} k(\cdot, x) d(\mathbb{Q}_1 - \mathbb{Q}_2)(x) = c \int_{\mathcal{X}} f(x) k(\cdot, x) d\mathbb{P}(x),\end{aligned}$$

$$\begin{aligned}(\mu_{\mathbb{Q}_1} - \mu_{\mathbb{Q}_2})(y) &= c \int_{\mathcal{X}} f(x) k(y, x) d\mathbb{P}(x) \\ &= c \langle f, \underbrace{k(y, \cdot)}_{\in \mathcal{H}_k} \rangle_{L^2(\mathbb{P})} = \mathbf{0} \quad (\forall y \in \mathcal{X}).\end{aligned}$$

Denseness in L^2 is necessary: proof continued

We arrive at

$$\begin{aligned}\mu_{\mathbb{Q}_1} - \mu_{\mathbb{Q}_2} &= \int_{\mathcal{X}} k(\cdot, x) d\mathbb{Q}_1(x) - \int_{\mathcal{X}} k(\cdot, x) d\mathbb{Q}_2(x) \\ &= \int_{\mathcal{X}} k(\cdot, x) d(\mathbb{Q}_1 - \mathbb{Q}_2)(x) = c \int_{\mathcal{X}} f(x) k(\cdot, x) d\mathbb{P}(x),\end{aligned}$$

$$\begin{aligned}(\mu_{\mathbb{Q}_1} - \mu_{\mathbb{Q}_2})(y) &= c \int_{\mathcal{X}} f(x) k(y, x) d\mathbb{P}(x) \\ &= c \langle f, \underbrace{k(y, \cdot)}_{\in \mathcal{H}_k} \rangle_{L^2(\mathbb{P})} = \mathbf{0} \quad (\forall y \in \mathcal{X}).\end{aligned}$$

Thus $\mu_{\mathbb{Q}_1} - \mu_{\mathbb{Q}_2} = 0$ despite $\mathbb{Q}_1 \neq \mathbb{Q}_2$.

Infinitely divisible distributions: quick summary

U : random variable.

Question

Can it be decomposed to the sum of 2 i.i.d. random variables?

Infinitely divisible distributions: quick summary

U : random variable.

Question

Can it be decomposed to the sum of 2 i.i.d. random variables?

Question

Can it be decomposed to the sum of 3 i.i.d. random variables?

Infinitely divisible distributions: quick summary

U : random variable.

Question

Can it be decomposed to the sum of 2 i.i.d. random variables?

Question

Can it be decomposed to the sum of 3 i.i.d. random variables?

Question

Can it be decomposed to the sum of 4 i.i.d. random variables?

Infinitely divisible distributions: quick summary

U : random variable.

Question

Can it be decomposed to the sum of 2 i.i.d. random variables?

Question

Can it be decomposed to the sum of 3 i.i.d. random variables?

Question

Can it be decomposed to the sum of 4 i.i.d. random variables?

Question

Can it be decomposed to the sum of n i.i.d. random variables for any $n \in \mathbb{Z}^+$?

Examples:

- Poisson, negative binomial, Gamma distribution, student t .

Examples:

- Poisson, negative binomial, Gamma distribution, student t .
- **normal**, Cauchy distribution

Examples:

- Poisson, negative binomial, Gamma distribution, student t .
- **normal** ($\alpha = 2$), Cauchy distribution ($\alpha = 1$) $\xleftarrow{\text{spec.}}$ $\forall \alpha$ -stable.

Examples:

- Poisson, negative binomial, Gamma distribution, student t .
- **normal** ($\alpha = 2$), Cauchy distribution ($\alpha = 1$) $\xleftarrow{\text{spec.}}$ $\forall \alpha$ -stable.

Counterexamples:

- uniform, binomial distribution

Examples:

- Poisson, negative binomial, Gamma distribution, student t .
- **normal** ($\alpha = 2$), Cauchy distribution ($\alpha = 1$) $\xleftarrow{\text{spec.}}$ $\forall \alpha$ -stable.

Counterexamples:

- uniform, binomial distribution $\xleftarrow{\text{spec.}}$ \forall any distribution with bounded (finite) support.

Theorem ([Nishiyama and Fukumizu, 2016])

Assume

- $k(x, y) = k_0(x - y)$, $k_0 \in C_b(\mathbb{R}^d)$, k_0 is the pdf of
- an infinitely divisible, symmetric distribution.

Then k is characteristic.

Theorem ([Nishiyama and Fukumizu, 2016])

Assume

- $k(x, y) = k_0(x - y)$, $k_0 \in C_b(\mathbb{R}^d)$, k_0 is the pdf of
- an infinitely divisible, symmetric distribution.

Then k is characteristic.

Examples: Gaussian, Matérn kernel, α -stable kernels, student t -kernels, . . .

Characteristic kernels: finished.

- Dependency measure applications.
- KCCA. Mean embedding: $\mu_{\mathbb{P}} = \int_X k(\cdot, x)d\mathbb{P}(x) \in \mathcal{H}_k$.
- Injectivity of μ on
 - probability distributions: characteristic property.
 - finite signed measures: universality (\mathcal{X} : compact metric).
- By definition: injectivity of $\mu \Leftrightarrow$

$$d_k(\mathbb{P}, \mathbb{Q}) := \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k}$$

is a metric.

Maximum mean discrepancy (MMD)

MMD is a specific integral probability metric (IPM)

- $\mathcal{F} = \left\{ f \in \mathcal{H}_k : \|f\|_{\mathcal{H}_k} = 1 \right\}$: unit ball in \mathcal{H}_k .

$$d_k(\mathbb{P}, \mathbb{Q}) = \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k}$$

MMD is a specific integral probability metric (IPM)

- $\mathcal{F} = \left\{ f \in \mathcal{H}_k : \|f\|_{\mathcal{H}_k} = 1 \right\}$: unit ball in \mathcal{H}_k .

$$\begin{aligned} d_k(\mathbb{P}, \mathbb{Q}) &= \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k} \\ &= \sup_{f \in \mathcal{F}} \langle f, \mu_{\mathbb{P}} - \mu_{\mathbb{Q}} \rangle_{\mathcal{H}_k} \end{aligned}$$

MMD is a specific integral probability metric (IPM)

- $\mathcal{F} = \left\{ f \in \mathcal{H}_k : \|f\|_{\mathcal{H}_k} = 1 \right\}$: unit ball in \mathcal{H}_k .

$$\begin{aligned} d_k(\mathbb{P}, \mathbb{Q}) &= \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k} \\ &= \sup_{f \in \mathcal{F}} \langle f, \mu_{\mathbb{P}} - \mu_{\mathbb{Q}} \rangle_{\mathcal{H}_k} \\ &= \sup_{f \in \mathcal{F}} (\mathbb{P}f - \mathbb{Q}f). \end{aligned}$$

MMD is a specific integral probability metric (IPM)

- $\mathcal{F} = \left\{ f \in \mathcal{H}_k : \|f\|_{\mathcal{H}_k} = 1 \right\}$: unit ball in \mathcal{H}_k .

$$\begin{aligned} d_k(\mathbb{P}, \mathbb{Q}) &= \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k} \\ &= \sup_{f \in \mathcal{F}} \langle f, \mu_{\mathbb{P}} - \mu_{\mathbb{Q}} \rangle_{\mathcal{H}_k} \\ &= \sup_{f \in \mathcal{F}} (\mathbb{P}f - \mathbb{Q}f). \end{aligned}$$

- IPMs [Zolotarev, 1983, Müller, 1997].

IPM: other \mathcal{F} examples giving metric

- $\mathcal{F} = C_b(\mathcal{X})$ with \mathcal{X} metric space.

IPM: other \mathcal{F} examples giving metric

- $\mathcal{F} = C_b(\mathcal{X})$ with \mathcal{X} metric space.
- $\mathcal{F} = \{f : \|f\|_\infty := \sup_{x \in \mathcal{X}} |f(x)| \leq 1\}$:
 - bounded functions.
 - total variation distance.

IPM: other \mathcal{F} examples giving metric

- $\mathcal{F} = C_b(\mathcal{X})$ with \mathcal{X} metric space.
- $\mathcal{F} = \{f : \|f\|_{\infty} := \sup_{x \in \mathcal{X}} |f(x)| \leq 1\}$:
 - bounded functions.
 - total variation distance.
- $\mathcal{F} = \left\{ f : \|f\|_L := \sup_{x \neq y} \frac{|f(x) - f(y)|}{\rho(x, y)} \leq 1 \right\}$:
 - Kantorovich metric $\xrightarrow{\mathcal{X}: \text{separable metric}}$ Wasserstein distance.

IPM: other \mathcal{F} examples giving metric

- $\mathcal{F} = \{f : \|f\|_\infty := \sup_{x \in \mathcal{X}} |f(x)| \leq 1\}$:
 - bounded functions.
 - total variation distance.

TV upper bounds MMD [Sriperumbudur et al., 2010b]:

$$d_k(\mathbb{P}, \mathbb{Q}) \leq \sup_{x \in \mathcal{X}} \sqrt{k(x, x)} \textcolor{blue}{TV}(\mathbb{P}, \mathbb{Q}).$$

- $\mathcal{F} = \{f : \|f\|_{BL} := \|f\|_\infty + \|f\|_L \leq 1\}$
 - bounded Lipschitz functions,
 - Dudley metric.

- $\mathcal{F} = \{f : \|f\|_{BL} := \|f\|_\infty + \|f\|_L \leq 1\}$
 - bounded Lipschitz functions,
 - Dudley metric.
- $\mathcal{F} = \{\chi_{(-\infty, t]} : t \in \mathbb{R}^d\}$:
 - indicator functions of half-intervals.
 - Kolmogorov distance.

[Sriperumbudur et al., 2012]:

- Kantorovich, Dudley metric: linear programming task.
- MMD (d_k): easier.

MMD estimators

MMD estimator: intuition

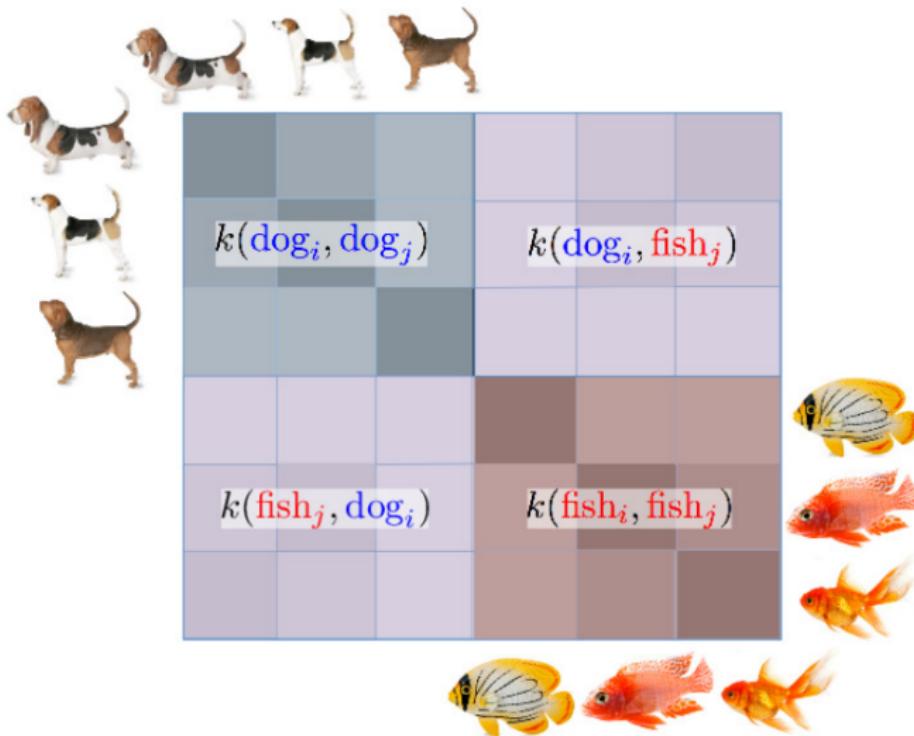


$\sim P$

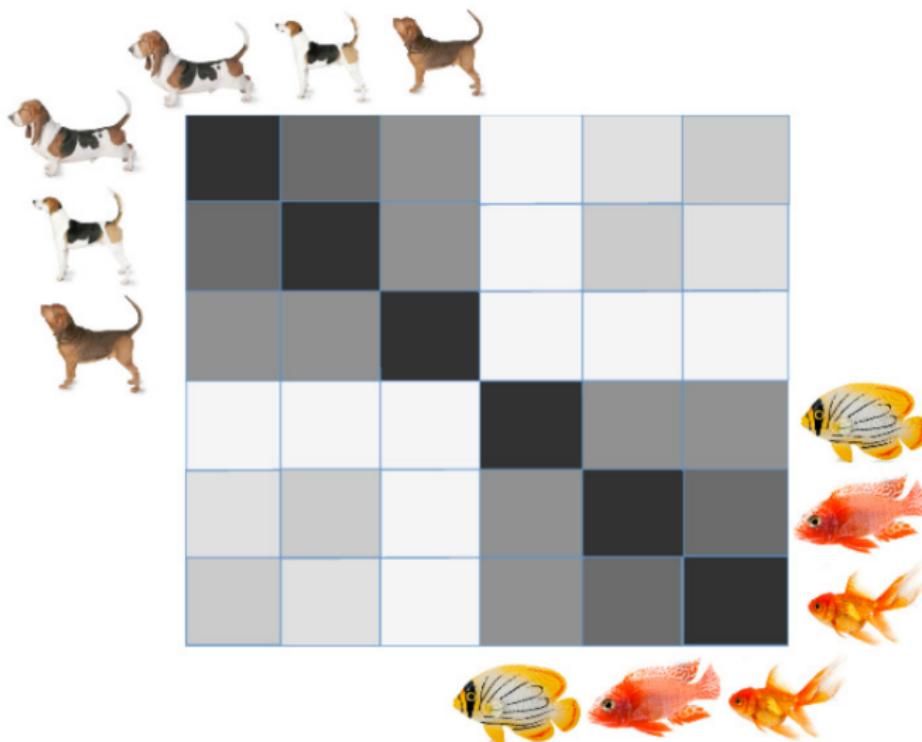


$\sim Q$

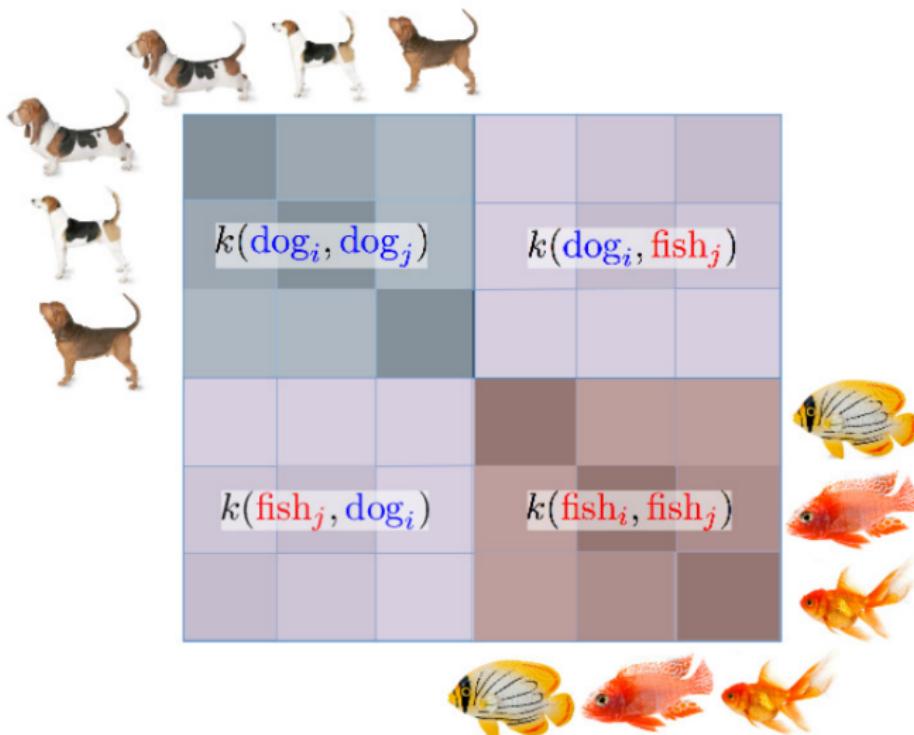
MMD estimator: intuition



MMD estimator: intuition



MMD estimator: intuition



$$\widehat{\text{MMD}}^2(\mathbb{P}, \mathbb{Q}) = \overline{G_{\mathbb{P}, \mathbb{P}}} + \overline{G_{\mathbb{Q}, \mathbb{Q}}} - 2\overline{G_{\mathbb{P}, \mathbb{Q}}} \quad (\text{without diagonals in } \overline{G_{\mathbb{P}, \mathbb{P}}}, \overline{G_{\mathbb{Q}, \mathbb{Q}}})$$

[†] MMD & HSIC illustration credit: Arthur Gretton

MMD estimator-1

Recall: MMD = squared difference between feature means:

$$\begin{aligned}\text{MMD}^2(\mathbb{P}, \mathbb{Q}) &:= d_k^2(\mathbb{P}, \mathbb{Q}) = \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k}^2 = \\ &= \mathbb{E}_{x \sim \mathbb{P}, x' \sim \mathbb{P}} k(x, x') + \mathbb{E}_{y \sim \mathbb{Q}, y' \sim \mathbb{Q}} k(y, y') \\ &\quad - 2\mathbb{E}_{x \sim \mathbb{P}, y \sim \mathbb{Q}} k(x, y).\end{aligned}$$

MMD estimator-1

Recall: MMD = squared difference between feature means:

$$\begin{aligned}\text{MMD}^2(\mathbb{P}, \mathbb{Q}) &:= d_k^2(\mathbb{P}, \mathbb{Q}) = \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k}^2 = \\ &= \mathbb{E}_{x \sim \mathbb{P}, x' \sim \mathbb{P}} k(x, x') + \mathbb{E}_{y \sim \mathbb{Q}, y' \sim \mathbb{Q}} k(y, y') \\ &\quad - 2\mathbb{E}_{x \sim \mathbb{P}, y \sim \mathbb{Q}} k(x, y).\end{aligned}$$

Unbiased empirical estimator using $\{x_i\}_{i=1}^m \sim \mathbb{P}$, $\{y_j\}_{j=1}^n \sim \mathbb{Q}$:

$$\widehat{\text{MMD}}_u^2(\mathbb{P}, \mathbb{Q}) = \overline{G_{\mathbb{P}, \mathbb{P}}} + \overline{G_{\mathbb{Q}, \mathbb{Q}}} - 2\overline{G_{\mathbb{P}, \mathbb{Q}}}$$

MMD estimator-1

Recall: MMD = squared difference between feature means:

$$\begin{aligned}\text{MMD}^2(\mathbb{P}, \mathbb{Q}) &:= d_k^2(\mathbb{P}, \mathbb{Q}) = \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k}^2 = \\ &= \mathbb{E}_{x \sim \mathbb{P}, x' \sim \mathbb{P}} k(x, x') + \mathbb{E}_{y \sim \mathbb{Q}, y' \sim \mathbb{Q}} k(y, y') \\ &\quad - 2\mathbb{E}_{x \sim \mathbb{P}, y \sim \mathbb{Q}} k(x, y).\end{aligned}$$

Unbiased empirical estimator using $\{x_i\}_{i=1}^m \sim \mathbb{P}$, $\{y_j\}_{j=1}^n \sim \mathbb{Q}$:

$$\begin{aligned}\widehat{\text{MMD}}_u^2(\mathbb{P}, \mathbb{Q}) &= \overline{G_{\mathbb{P}, \mathbb{P}}} + \overline{G_{\mathbb{Q}, \mathbb{Q}}} - 2\overline{G_{\mathbb{P}, \mathbb{Q}}} \\ &= \underbrace{\frac{1}{m(m-1)} \sum_{i=1}^m \sum_{j \neq i}^m k(\mathbf{x}_i, \mathbf{x}_j)}_{\text{U-statistic-1}} + \underbrace{\frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j \neq i}^n k(\mathbf{y}_i, \mathbf{y}_j)}_{\text{U-statistic-2}} \\ &\quad - \underbrace{\frac{2}{mn} \sum_{i=1}^m \sum_{j=1}^n k(\mathbf{x}_i, \mathbf{y}_j)}_{\text{sample average}}.\end{aligned}$$

MMD estimator-2

We plug in the empirical measures $(\mathbb{P}_m, \mathbb{Q}_n)$:

$$\text{MMD}^2(\mathbb{P}, \mathbb{Q}) = \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k}^2,$$

$$\widehat{\text{MMD}}_b^2(\mathbb{P}, \mathbb{Q}) = \|\mu_{\mathbb{P}_m} - \mu_{\mathbb{Q}_n}\|_{\mathcal{H}_k}^2$$

MMD estimator-2

We plug in the empirical measures $(\mathbb{P}_m, \mathbb{Q}_n)$:

$$\text{MMD}^2(\mathbb{P}, \mathbb{Q}) = \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k}^2,$$

$$\begin{aligned}\widehat{\text{MMD}}_b^2(\mathbb{P}, \mathbb{Q}) &= \|\mu_{\mathbb{P}_m} - \mu_{\mathbb{Q}_n}\|_{\mathcal{H}_k}^2 \\ &= \|\mu_{\mathbb{P}_m}\|_{\mathcal{H}_k}^2 + \|\mu_{\mathbb{Q}_n}\|_{\mathcal{H}_k}^2 - 2\langle \mu_{\mathbb{P}_m}, \mu_{\mathbb{Q}_n} \rangle_{\mathcal{H}_k}.\end{aligned}$$

MMD estimator-2

We plug in the empirical measures $(\mathbb{P}_m, \mathbb{Q}_n)$:

$$\begin{aligned}\text{MMD}^2(\mathbb{P}, \mathbb{Q}) &= \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k}^2, \\ \widehat{\text{MMD}}_b^2(\mathbb{P}, \mathbb{Q}) &= \|\mu_{\mathbb{P}_m} - \mu_{\mathbb{Q}_n}\|_{\mathcal{H}_k}^2 \\ &= \|\mu_{\mathbb{P}_m}\|_{\mathcal{H}_k}^2 + \|\mu_{\mathbb{Q}_n}\|_{\mathcal{H}_k}^2 - 2\langle \mu_{\mathbb{P}_m}, \mu_{\mathbb{Q}_n} \rangle_{\mathcal{H}_k}.\end{aligned}$$

Enough:

$$\langle \mu_{\mathbb{P}_m}, \mu_{\mathbb{Q}_n} \rangle_{\mathcal{H}_k} = \left\langle \frac{1}{m} \sum_{i=1}^m k(\cdot, x_i), \frac{1}{n} \sum_{j=1}^n k(\cdot, y_j) \right\rangle_{\mathcal{H}_k}$$

MMD estimator-2

We plug in the empirical measures $(\mathbb{P}_m, \mathbb{Q}_n)$:

$$\begin{aligned}\text{MMD}^2(\mathbb{P}, \mathbb{Q}) &= \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k}^2, \\ \widehat{\text{MMD}}_b^2(\mathbb{P}, \mathbb{Q}) &= \|\mu_{\mathbb{P}_m} - \mu_{\mathbb{Q}_n}\|_{\mathcal{H}_k}^2 \\ &= \|\mu_{\mathbb{P}_m}\|_{\mathcal{H}_k}^2 + \|\mu_{\mathbb{Q}_n}\|_{\mathcal{H}_k}^2 - 2\langle \mu_{\mathbb{P}_m}, \mu_{\mathbb{Q}_n} \rangle_{\mathcal{H}_k}.\end{aligned}$$

Enough:

$$\begin{aligned}\langle \mu_{\mathbb{P}_m}, \mu_{\mathbb{Q}_n} \rangle_{\mathcal{H}_k} &= \left\langle \frac{1}{m} \sum_{i=1}^m k(\cdot, x_i), \frac{1}{n} \sum_{j=1}^n k(\cdot, y_j) \right\rangle_{\mathcal{H}_k} \\ &= \frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n \underbrace{\langle k(\cdot, x_i), k(\cdot, y_j) \rangle_{\mathcal{H}_k}}_{k(x_i, y_j)}.\end{aligned}$$

MMD estimator-2: continued

$$\widehat{\text{MMD}}_b^2(\mathbb{P}, \mathbb{Q}) = \underbrace{\frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m k(x_i, x_j)}_{\text{V-statistic-1}} + \underbrace{\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n k(y_i, y_j)}_{\text{V-statistic-2}} - \underbrace{\frac{2}{mn} \sum_{i=1}^m \sum_{j=1}^n k(x_i, y_j)}_{\text{sample average}}.$$

MMD estimator-2: continued

$$\widehat{\text{MMD}}_b^2(\mathbb{P}, \mathbb{Q}) = \underbrace{\frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m k(x_i, x_j)}_{\text{V-statistic-1}} + \underbrace{\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n k(y_i, y_j)}_{\text{V-statistic-2}} - \underbrace{\frac{2}{mn} \sum_{i=1}^m \sum_{j=1}^n k(x_i, y_j)}_{\text{sample average}}.$$

Notes:

- $\widehat{\text{MMD}}_u^2(\mathbb{P}, \mathbb{Q})$: unbiased; it might be negative.

MMD estimator-2: continued

$$\widehat{\text{MMD}}_b^2(\mathbb{P}, \mathbb{Q}) = \underbrace{\frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m k(x_i, x_j)}_{\text{V-statistic-1}} + \underbrace{\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n k(y_i, y_j)}_{\text{V-statistic-2}} - \underbrace{\frac{2}{mn} \sum_{i=1}^m \sum_{j=1}^n k(x_i, y_j)}_{\text{sample average}}.$$

Notes:

- $\widehat{\text{MMD}}_u^2(\mathbb{P}, \mathbb{Q})$: unbiased; it might be negative.
- $\widehat{\text{MMD}}_b^2(\mathbb{P}, \mathbb{Q}) = \|\mu_{\mathbb{P}_m} - \mu_{\mathbb{Q}_n}\|_{\mathcal{H}_k}^2 \geq 0$.

MMD estimator-2: continued

$$\widehat{\text{MMD}}_b^2(\mathbb{P}, \mathbb{Q}) = \underbrace{\frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m k(x_i, x_j)}_{\text{V-statistic-1}} + \underbrace{\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n k(y_i, y_j)}_{\text{V-statistic-2}} - \underbrace{\frac{2}{mn} \sum_{i=1}^m \sum_{j=1}^n k(x_i, y_j)}_{\text{sample average}}.$$

Notes:

- $\widehat{\text{MMD}}_u^2(\mathbb{P}, \mathbb{Q})$: unbiased; it might be negative.
- $\widehat{\text{MMD}}_b^2(\mathbb{P}, \mathbb{Q}) = \|\mu_{\mathbb{P}_m} - \mu_{\mathbb{Q}_n}\|_{\mathcal{H}_k}^2 \geq 0$.
- Computational complexity: $\mathcal{O}((m+n)^2)$, quadratic.

- Set kernel, convolution kernel.

- Set kernel, convolution kernel.
- Other valid $K(\mu_{\mathbb{P}}, \mu_{\mathbb{Q}})$ examples → distribution classification [Póczos et al., 2012, Muandet et al., 2011] / distribution regression [Szabó et al., 2016].

- Set kernel, convolution kernel.
- Other valid $K(\mu_{\mathbb{P}}, \mu_{\mathbb{Q}})$ examples → distribution classification [Póczos et al., 2012, Muandet et al., 2011] / distribution regression [Szabó et al., 2016].
- Few analytic expressions exist for MMD.

- Set kernel, convolution kernel.
- Other valid $K(\mu_{\mathbb{P}}, \mu_{\mathbb{Q}})$ examples → distribution classification [Póczos et al., 2012, Muandet et al., 2011] / distribution regression [Szabó et al., 2016].
- Few analytic expressions exist for MMD.
- Embedding to RKBS.

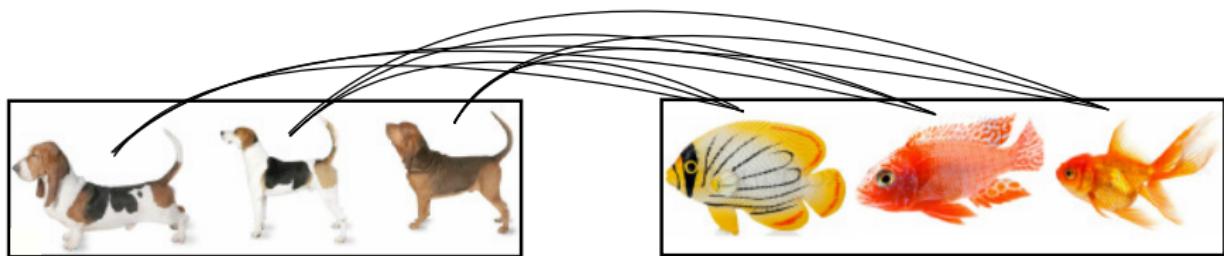
- Set kernel, convolution kernel.
- Other valid $K(\mu_{\mathbb{P}}, \mu_{\mathbb{Q}})$ examples → distribution classification [Póczos et al., 2012, Muandet et al., 2011] / distribution regression [Szabó et al., 2016].
- Few analytic expressions exist for MMD.
- Embedding to RKBS.

Let us see the details.

Set kernel

Convolution kernels [Haussler, 1999] \ni set kernel [Gärtner et al., 2002]:

$$K(\mathbb{P}_m, \mathbb{Q}_n) := \langle \mu_{\mathbb{P}_m}, \mu_{\mathbb{Q}_n} \rangle_{\mathcal{H}_k} = \frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n k(x_i, y_j).$$



Other valid K examples [Christmann and Steinwart, 2010],
[Szabó et al., 2015] → distribution regression

Recall: $K(\mathbb{P}, \mathbb{Q}) = \langle \mu_{\mathbb{P}}, \mu_{\mathbb{Q}} \rangle_{\mathcal{H}_k}$, linear kernel.

K_G	K_e	K_C
$e^{-\frac{\ \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\ _{\mathcal{H}_k}^2}{2\theta^2}}$	$e^{-\frac{\ \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\ _{\mathcal{H}_k}}{2\theta^2}}$	$\left(1 + \ \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\ _{\mathcal{H}_k}^2 / \theta^2\right)^{-1}$

K_t	K_i
$\left(1 + \ \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\ _{\mathcal{H}_k}^\theta\right)^{-1}$	$\left(\ \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\ _{\mathcal{H}_k}^2 + \theta^2\right)^{-\frac{1}{2}}$

Other valid K examples [Christmann and Steinwart, 2010],
[Szabó et al., 2015] → distribution regression

Recall: $K(\mathbb{P}, \mathbb{Q}) = \langle \mu_{\mathbb{P}}, \mu_{\mathbb{Q}} \rangle_{\mathcal{H}_k}$, linear kernel.

K_G	K_e	K_C
$e^{-\frac{\ \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\ _{\mathcal{H}_k}^2}{2\theta^2}}$	$e^{-\frac{\ \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\ _{\mathcal{H}_k}}{2\theta^2}}$	$\left(1 + \ \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\ _{\mathcal{H}_k}^2 / \theta^2\right)^{-1}$

K_t	K_i
$\left(1 + \ \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\ _{\mathcal{H}_k}^\theta\right)^{-1}$	$\left(\ \mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\ _{\mathcal{H}_k}^2 + \theta^2\right)^{-\frac{1}{2}}$

Functions of $\|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k}$ ⇒ computation: similar to set kernel.

Few analytic expressions exist: examples
[Gretton et al., 2007, Muandet et al., 2011]

Assume: $\mathbb{P} = N(m_1, \Sigma_1)$, $\mathbb{Q} = N(m_2, \Sigma_2)$.

$k(x, y)$	$K(\mu_{\mathbb{P}}, \mu_{\mathbb{Q}}) = \langle \mu_{\mathbb{P}}, \mu_{\mathbb{Q}} \rangle_{\mathcal{H}_k}$
$e^{-\frac{\gamma}{2}\ x-y\ _2^2}$	$\frac{e^{-\frac{1}{2}(m_1-m_2)^T(\Sigma_1+\Sigma_2+\gamma I)^{-1}(m_1-m_2)}}{ \gamma\Sigma_1+\gamma\Sigma_2+I ^{\frac{1}{2}}}$

Few analytic expressions exist: examples
[Gretton et al., 2007, Muandet et al., 2011]

Assume: $\mathbb{P} = N(m_1, \Sigma_1)$, $\mathbb{Q} = N(m_2, \Sigma_2)$.

$k(x, y)$	$K(\mu_{\mathbb{P}}, \mu_{\mathbb{Q}}) = \langle \mu_{\mathbb{P}}, \mu_{\mathbb{Q}} \rangle_{\mathcal{H}_k}$
$e^{-\frac{\gamma}{2}\ x-y\ _2^2}$	$\frac{e^{-\frac{1}{2}(m_1-m_2)^T(\Sigma_1+\Sigma_2+\gamma I)^{-1}(m_1-m_2)}}{ \gamma\Sigma_1+\gamma\Sigma_2+I ^{\frac{1}{2}}}$
$(1 + \langle x, y \rangle)^2$	$(1 + \langle m_1, m_2 \rangle)^2 + \text{tr}(\Sigma_1\Sigma_2) + m_1\Sigma_2m_1 + m_2\Sigma_1m_2$

Few analytic expressions exist: examples
[Gretton et al., 2007, Muandet et al., 2011]

Assume: $\mathbb{P} = N(m_1, \Sigma_1)$, $\mathbb{Q} = N(m_2, \Sigma_2)$.

$k(x, y)$	$K(\mu_{\mathbb{P}}, \mu_{\mathbb{Q}}) = \langle \mu_{\mathbb{P}}, \mu_{\mathbb{Q}} \rangle_{\mathcal{H}_k}$
$e^{-\frac{\gamma}{2}\ x-y\ _2^2}$	$\frac{e^{-\frac{1}{2}(m_1-m_2)^T(\Sigma_1+\Sigma_2+\gamma I)^{-1}(m_1-m_2)}}{ \gamma\Sigma_1+\gamma\Sigma_2+I ^{\frac{1}{2}}}$
$(1 + \langle x, y \rangle)^2$	$(1 + \langle m_1, m_2 \rangle)^2 + \text{tr}(\Sigma_1 \Sigma_2) + m_1 \Sigma_2 m_1 + m_2 \Sigma_1 m_2$
$(1 + \langle x, y \rangle)^3$	$(1 + \langle m_1, m_2 \rangle)^3 + 6m_1^T \Sigma_1 \Sigma_2 m_2 + 3(1 + \langle m_1, m_2 \rangle) \times [\text{tr}(\Sigma_1 \Sigma_2) + m_1 \Sigma_2 m_1 + m_2 \Sigma_1 m_2]$

Definition ([Zhang et al., 2009])

A reflexive \mathcal{B} Banach space $[(\mathcal{B}')' = \mathcal{B}]$ of functions on \mathcal{X} is called RKBS if

Definition ([Zhang et al., 2009])

A reflexive \mathcal{B} Banach space $[(\mathcal{B}')' = \mathcal{B}]$ of functions on \mathcal{X} is called RKBS if

- \mathcal{B}' is isometric to a Banach space of functions on \mathcal{X} , and

Definition ([Zhang et al., 2009])

A reflexive \mathcal{B} Banach space $[(\mathcal{B}')' = \mathcal{B}]$ of functions on \mathcal{X} is called RKBS if

- \mathcal{B}' is isometric to a Banach space of functions on \mathcal{X} , and
- the evaluation functional is continuous on both \mathcal{B} and \mathcal{B}' .

Definition ([Zhang et al., 2009])

A reflexive \mathcal{B} Banach space $[(\mathcal{B}')' = \mathcal{B}]$ of functions on \mathcal{X} is called RKBS if

- \mathcal{B}' is isometric to a Banach space of functions on \mathcal{X} , and
- the evaluation functional is continuous on both \mathcal{B} and \mathcal{B}' .

Notes:

- Generally, $\mathcal{B} \subseteq \mathcal{B}''$.

Definition ([Zhang et al., 2009])

A reflexive \mathcal{B} Banach space $[(\mathcal{B}')' = \mathcal{B}]$ of functions on \mathcal{X} is called RKBS if

- \mathcal{B}' is isometric to a Banach space of functions on \mathcal{X} , and
- the evaluation functional is continuous on both \mathcal{B} and \mathcal{B}' .

Notes:

- Generally, $\mathcal{B} \subseteq \mathcal{B}''$.
- For $\mathcal{B} = \mathcal{H}$ Hilbert: $(\mathcal{H}')' = \mathcal{H}$ (Riesz representation theorem).

RKBS properties

Using the

$$\langle f, g' \rangle_{\mathcal{B}} := g'(f), \quad (f \in \mathcal{B}, g' \in \mathcal{B}')$$

notation

RKBS properties

Using the

$$\langle f, g' \rangle_{\mathcal{B}} := g'(f), \quad (f \in \mathcal{B}, g' \in \mathcal{B}')$$

notation

$$k(\cdot, x) \in \mathcal{B}' \quad (\forall x \in \mathcal{X}), \qquad f(x) = \langle f, k(\cdot, x) \rangle_{\mathcal{B}} \quad \forall f \in \mathcal{B},$$

An RKBS has exactly one reproducing kernel, but...

RKBS properties

Using the

$$\langle f, g' \rangle_{\mathcal{B}} := g'(f), \quad (f \in \mathcal{B}, g' \in \mathcal{B}')$$

notation

$$k(\cdot, x) \in \mathcal{B}' \quad (\forall x \in \mathcal{X}),$$

$$k(x, \cdot) \in \mathcal{B} \quad (\forall x \in \mathcal{X}),$$

$$f(x) = \langle f, k(\cdot, x) \rangle_{\mathcal{B}} \quad \forall f \in \mathcal{B},$$

$$f'(x) = \langle k(x, \cdot), f' \rangle_{\mathcal{B}} \quad \forall f' \in \mathcal{B}',$$

An RKBS has exactly one reproducing kernel, but...

RKBS properties

Using the

$$\langle f, g' \rangle_{\mathcal{B}} := g'(f), \quad (f \in \mathcal{B}, g' \in \mathcal{B}')$$

notation

$$k(\cdot, x) \in \mathcal{B}' \quad (\forall x \in \mathcal{X}),$$

$$k(x, \cdot) \in \mathcal{B} \quad (\forall x \in \mathcal{X}),$$

$$\mathcal{B} = \overline{\text{span}\{k(x, \cdot), x \in \mathcal{X}\}},$$

$$f(x) = \langle f, k(\cdot, x) \rangle_{\mathcal{B}} \quad \forall f \in \mathcal{B},$$

$$f'(x) = \langle k(x, \cdot), f' \rangle_{\mathcal{B}} \quad \forall f' \in \mathcal{B}',$$

$$\mathcal{B}' = \overline{\text{span}\{k(\cdot, x), x \in \mathcal{X}\}},$$

An RKBS has exactly one reproducing kernel, but...

RKBS properties

Using the

$$\langle f, g' \rangle_{\mathcal{B}} := g'(f), \quad (f \in \mathcal{B}, g' \in \mathcal{B}')$$

notation

$$k(\cdot, x) \in \mathcal{B}' \quad (\forall x \in \mathcal{X}),$$

$$k(x, \cdot) \in \mathcal{B} \quad (\forall x \in \mathcal{X}),$$

$$\mathcal{B} = \overline{\text{span}\{k(x, \cdot), x \in \mathcal{X}\}},$$

$$k(x, y) = \langle k(x, \cdot), k(\cdot, y) \rangle_{\mathcal{B}} \quad \forall x, y \in \mathcal{X}.$$

$$f(x) = \langle f, k(\cdot, x) \rangle_{\mathcal{B}} \quad \forall f \in \mathcal{B},$$

$$f'(x) = \langle k(x, \cdot), f' \rangle_{\mathcal{B}} \quad \forall f' \in \mathcal{B}',$$

$$\mathcal{B}' = \overline{\text{span}\{k(\cdot, x), x \in \mathcal{X}\}},$$

An RKBS has exactly one reproducing kernel, but...

Peculiarities of RKBS-s

- Different RKBSs can have the same r.k.

Peculiarities of RKBS-s

- Different RKBSs can have the same r.k.
- No inner product on \mathcal{B} \Rightarrow an r.k. can be an **arbitrary** function.

Peculiarities of RKBS-s

- Different RKBSs can have the same r.k.
- No inner product on \mathcal{B} \Rightarrow an r.k. can be an **arbitrary** function.
- For specific RKBSs[†]:

[†]Uniformly Fréchet differentiable and uniformly convex, e.g. $L^p(\mathcal{X}, \mathcal{A}, \mu)$, $p \in (1, \infty)$.

Peculiarities of RKBS-s

- Different RKBSs can have the same r.k.
- No inner product on \mathcal{B} \Rightarrow an r.k. can be an **arbitrary** function.
- For specific RKBSs[†]:
 - 'Riesz representation theorem' exists, ...

[†]Uniformly Fréchet differentiable and uniformly convex, e.g. $L^p(\mathcal{X}, \mathcal{A}, \mu)$, $p \in (1, \infty)$.

Peculiarities of RKBS-s

- Different RKBSs can have the same r.k.
- No inner product on \mathcal{B} \Rightarrow an r.k. can be an **arbitrary** function.
- For specific RKBSs[†]:
 - 'Riesz representation theorem' exists, . . .
 - $\mu_{\mathbb{P}} = \int_{\mathcal{X}} \underbrace{k(\cdot, x)}_{\in \mathcal{B}'} d\mathbb{P}(x) \in \mathcal{B}'$ [Sriperumbudur et al., 2011].

[†]Uniformly Fréchet differentiable and uniformly convex, e.g. $L^p(\mathcal{X}, \mathcal{A}, \mu)$, $p \in (1, \infty)$.

Key for RKHS \mathcal{H}_k :

$$d_k^2(\mathbb{P}, \mathbb{Q}) = \int_{\mathcal{X}} \int_{\mathcal{X}} k(x, y) d(\mathbb{P} - \mathbb{Q})(x) d(\mathbb{P} - \mathbb{Q})(y).$$

Key for RKHS \mathcal{H}_k :

$$d_k^2(\mathbb{P}, \mathbb{Q}) = \int_{\mathcal{X}} \int_{\mathcal{X}} k(x, y) d(\mathbb{P} - \mathbb{Q})(x) d(\mathbb{P} - \mathbb{Q})(y).$$

For RKBS \mathcal{B} :

- d_k : **not expressible** in terms of $k(x, y)$,

Key for RKHS \mathcal{H}_k :

$$d_k^2(\mathbb{P}, \mathbb{Q}) = \int_{\mathcal{X}} \int_{\mathcal{X}} k(x, y) d(\mathbb{P} - \mathbb{Q})(x) d(\mathbb{P} - \mathbb{Q})(y).$$

For RKBS \mathcal{B} :

- d_k : **not expressible** in terms of $k(x, y)$,
- associated distances and estimators: **no closed form expressions**.

MMD: finished

Covariance operator

Idea: (un)centered cross-covariance

$$C_{xy}^{\textcolor{blue}{u}} = \mathbb{E}_{xy} [xy^T],$$

u: uncentered, **c**: centered.

Idea: (un)centered cross-covariance

$$C_{xy}^{\textcolor{blue}{u}} = \mathbb{E}_{xy} [xy^T], \quad C_{xy}^{\textcolor{red}{c}} = \mathbb{E}_{xy} [(x - \mathbb{E}x)(y - \mathbb{E}y)^T],$$

u: uncentered, **c**: centered.

Idea: (un)centered cross-covariance

$$C_{xy}^{\textcolor{blue}{u}} = \mathbb{E}_{xy} [xy^T], \quad C_{xy}^{\textcolor{red}{c}} = \mathbb{E}_{xy} [(x - \mathbb{E}x)(y - \mathbb{E}y)^T],$$
$$C_{xy}^{\textcolor{blue}{u}} = \mathbb{E}_{xy} [\varphi(x) \otimes \psi(y)],$$

u: uncentered, **c**: centered.

Idea: (un)centered cross-covariance

$$C_{xy}^{\textcolor{blue}{u}} = \mathbb{E}_{xy} [xy^T], \quad C_{xy}^{\textcolor{red}{c}} = \mathbb{E}_{xy} [(x - \mathbb{E}x)(y - \mathbb{E}y)^T],$$
$$C_{xy}^{\textcolor{blue}{u}} = \mathbb{E}_{xy} [\varphi(x) \otimes \psi(y)], \quad C_{xy}^{\textcolor{red}{c}} = \mathbb{E}_{xy} [(\varphi(x) - \mathbb{E}_x \varphi(x)) \otimes (\psi(y) - \mathbb{E}_y \psi(y))]$$

u: uncentered, **c**: centered.

Idea: (un)centered cross-covariance

$$C_{xy}^{\textcolor{blue}{u}} = \mathbb{E}_{xy} [xy^T], \quad C_{xy}^{\textcolor{red}{c}} = \mathbb{E}_{xy} [(x - \mathbb{E}x)(y - \mathbb{E}y)^T],$$
$$C_{xy}^{\textcolor{blue}{u}} = \mathbb{E}_{xy} [\varphi(x) \otimes \psi(y)], \quad C_{xy}^{\textcolor{red}{c}} = \mathbb{E}_{xy} [(\varphi(x) - \mathbb{E}_x \varphi(x)) \otimes (\psi(y) - \mathbb{E}_y \psi(y))]$$

u: uncentered, **c**: centered. In short, $xy^T \rightarrow \varphi(x) \otimes \psi(y)$.

$$C_{xy}^c = \mathbb{E}_{xy} [(\varphi(x) - \mathbb{E}_x \varphi(x)) \otimes (\psi(y) - \mathbb{E}_y \psi(y))]$$

encodes the dependency of x and y .

- $C_{xy}^c = 0 \Leftrightarrow x \perp y$ for 'rich' $\mathcal{H}_k, \mathcal{H}_\ell$.

$$C_{xy}^c = \mathbb{E}_{xy} [(\varphi(x) - \mathbb{E}_x \varphi(x)) \otimes (\psi(y) - \mathbb{E}_y \psi(y))]$$

encodes the dependency of x and y .

- $C_{xy}^c = 0 \Leftrightarrow x \perp y$ for 'rich' $\mathcal{H}_k, \mathcal{H}_\ell$.
- $\text{HSIC}(x, y) = \|C_{xy}^c\|_{HS}$. $\|\cdot\|_{HS}$: extension of Frobenius norm.

$$C_{xy}^c = \mathbb{E}_{xy} [(\varphi(x) - \mathbb{E}_x \varphi(x)) \otimes (\psi(y) - \mathbb{E}_y \psi(y))]$$

encodes the dependency of x and y .

- $C_{xy}^c = 0 \Leftrightarrow x \perp y$ for 'rich' $\mathcal{H}_k, \mathcal{H}_\ell$.
- $\text{HSIC}(x, y) = \|C_{xy}^c\|_{HS}$. $\|\cdot\|_{HS}$: extension of Frobenius norm.
- $\text{HSIC}(x, y)$: It will be easy to estimate. KCCA alternative.

$$C_{xy}^c = \mathbb{E}_{xy} [(\varphi(x) - \mathbb{E}_x \varphi(x)) \otimes (\psi(y) - \mathbb{E}_y \psi(y))]$$

encodes the dependency of x and y .

- $C_{xy}^c = 0 \Leftrightarrow x \perp y$ for 'rich' $\mathcal{H}_k, \mathcal{H}_\ell$.
- $\text{HSIC}(x, y) = \|C_{xy}^c\|_{HS}$. $\|\cdot\|_{HS}$: extension of Frobenius norm.
- $\text{HSIC}(x, y)$: It will be easy to estimate. KCCA alternative.

Question

What is $\varphi(x) \otimes \psi(y)$ and $\|\cdot\|_{HS}$?

Intuition of $a \otimes b$, goal: $a := \varphi(x) \in \mathcal{H}_k$, $b := \psi(y) \in \mathcal{H}_\ell$

- If $a \in \mathbb{R}^{d_1}$, $b \in \mathbb{R}^{d_2}$, then $ab^T \in \mathbb{R}^{d_1 \times d_2}$.

Intuition of $a \otimes b$, goal: $a := \varphi(x) \in \mathcal{H}_k$, $b := \psi(y) \in \mathcal{H}_\ell$

- If $a \in \mathbb{R}^{d_1}$, $b \in \mathbb{R}^{d_2}$, then $\textcolor{blue}{ab^T} \in \mathbb{R}^{d_1 \times d_2}$.
- For $g \in \mathbb{R}^{d_2}$

$$\left(\textcolor{blue}{ab^T} \right) g = a \left(b^T g \right)$$

Intuition of $a \otimes b$, goal: $a := \varphi(x) \in \mathcal{H}_k$, $b := \psi(y) \in \mathcal{H}_\ell$

- If $a \in \mathbb{R}^{d_1}$, $b \in \mathbb{R}^{d_2}$, then $ab^T \in \mathbb{R}^{d_1 \times d_2}$.
- For $g \in \mathbb{R}^{d_2}$

$$(ab^T)g = a(b^Tg) = a \langle b, g \rangle \in \mathbb{R}^{d_1},$$

$ab^T : \mathbb{R}^{d_2} \rightarrow \mathbb{R}^{d_1}$ linear mapping.

Intuition of $a \otimes b$, goal: $a := \varphi(x) \in \mathcal{H}_k$, $b := \psi(y) \in \mathcal{H}_\ell$

- If $a \in \mathbb{R}^{d_1}$, $b \in \mathbb{R}^{d_2}$, then $\textcolor{blue}{ab^T} \in \mathbb{R}^{d_1 \times d_2}$.
- For $g \in \mathbb{R}^{d_2}$

$$\left(\textcolor{blue}{ab^T} \right) g = a \left(b^T g \right) = a \langle b, g \rangle \in \mathbb{R}^{d_1},$$

$ab^T : \mathbb{R}^{d_2} \rightarrow \mathbb{R}^{d_1}$ linear mapping.

- Alternatively

$$\mathbb{R} \ni f^T \left(\textcolor{blue}{ab^T} \right) g = \left(f^T a \right) \left(b^T g \right)$$

$ab^T : \mathbb{R}^{d_1} \times \mathbb{R}^{d_2} \rightarrow \mathbb{R}$ bilinear form.

Intuition of $a \otimes b$, goal: $a := \varphi(x) \in \mathcal{H}_k$, $b := \psi(y) \in \mathcal{H}_\ell$

- If $a \in \mathbb{R}^{d_1}$, $b \in \mathbb{R}^{d_2}$, then $\textcolor{blue}{ab^T} \in \mathbb{R}^{d_1 \times d_2}$.
- For $g \in \mathbb{R}^{d_2}$

$$\left(\textcolor{blue}{ab^T} \right) g = a \left(b^T g \right) = a \langle b, g \rangle \in \mathbb{R}^{d_1},$$

$ab^T : \mathbb{R}^{d_2} \rightarrow \mathbb{R}^{d_1}$ linear mapping.

- Alternatively

$$\mathbb{R} \ni f^T \left(\textcolor{blue}{ab^T} \right) g = \left(f^T a \right) \left(b^T g \right) = \langle f, a \rangle \langle g, b \rangle$$

$ab^T : \mathbb{R}^{d_1} \times \mathbb{R}^{d_2} \rightarrow \mathbb{R}$ bilinear form.

Definition of $a \otimes b$, $\mathcal{H}_1 \otimes \mathcal{H}_2$

- Given: $\mathcal{H}_1, \mathcal{H}_2$ Hilbert spaces.
- $a \otimes b : (f, g) \in \mathcal{H}_1 \times \mathcal{H}_2 \rightarrow \mathbb{R}$ is the bilinear form:

$$(a \otimes b)(f, g) := \langle f, a \rangle_{\mathcal{H}_1} \langle g, b \rangle_{\mathcal{H}_2}.$$

Definition of $a \otimes b$, $\mathcal{H}_1 \otimes \mathcal{H}_2$

- Given: $\mathcal{H}_1, \mathcal{H}_2$ Hilbert spaces.
- $a \otimes b : (f, g) \in \mathcal{H}_1 \times \mathcal{H}_2 \rightarrow \mathbb{R}$ is the bilinear form:

$$(a \otimes b)(f, g) := \langle f, a \rangle_{\mathcal{H}_1} \langle g, b \rangle_{\mathcal{H}_2}.$$

- Finite linear combinations of $a \otimes b$ -s:

$$\mathcal{L} := \left\{ \sum_{i=1}^n c_i (a_i \otimes b_i), c_i \in \mathbb{R}, a_i \in \mathcal{H}_1, b_i \in \mathcal{H}_2, n \in \mathbb{Z}^+ \right\}.$$

Definition of $a \otimes b$, $\mathcal{H}_1 \otimes \mathcal{H}_2$

- Given: $\mathcal{H}_1, \mathcal{H}_2$ Hilbert spaces.
- $a \otimes b : (f, g) \in \mathcal{H}_1 \times \mathcal{H}_2 \rightarrow \mathbb{R}$ is the bilinear form:

$$(a \otimes b)(f, g) := \langle f, a \rangle_{\mathcal{H}_1} \langle g, b \rangle_{\mathcal{H}_2}.$$

- Finite linear combinations of $a \otimes b$ -s:

$$\mathcal{L} := \left\{ \sum_{i=1}^n c_i (a_i \otimes b_i), c_i \in \mathbb{R}, a_i \in \mathcal{H}_1, b_i \in \mathcal{H}_2, n \in \mathbb{Z}^+ \right\}.$$

- Define inner product on \mathcal{L} , and extend by linearity

$$\langle a_1 \otimes b_1, a_2 \otimes b_2 \rangle := \langle a_1, a_2 \rangle_{\mathcal{H}_1} \langle b_1, b_2 \rangle_{\mathcal{H}_2}.$$

Definition of $a \otimes b$, $\mathcal{H}_1 \otimes \mathcal{H}_2$

- Given: $\mathcal{H}_1, \mathcal{H}_2$ Hilbert spaces.
- $a \otimes b : (f, g) \in \mathcal{H}_1 \times \mathcal{H}_2 \rightarrow \mathbb{R}$ is the bilinear form:

$$(a \otimes b)(f, g) := \langle f, a \rangle_{\mathcal{H}_1} \langle g, b \rangle_{\mathcal{H}_2}.$$

- Finite linear combinations of $a \otimes b$ -s:

$$\mathcal{L} := \left\{ \sum_{i=1}^n c_i (a_i \otimes b_i), c_i \in \mathbb{R}, a_i \in \mathcal{H}_1, b_i \in \mathcal{H}_2, n \in \mathbb{Z}^+ \right\}.$$

- Define inner product on \mathcal{L} , and extend by linearity

$$\langle a_1 \otimes b_1, a_2 \otimes b_2 \rangle := \langle a_1, a_2 \rangle_{\mathcal{H}_1} \langle b_1, b_2 \rangle_{\mathcal{H}_2}.$$

- $\mathcal{H}_1 \otimes \mathcal{H}_2$: completion of \mathcal{L} .

$a_1 \otimes \dots \otimes a_M$, $\mathcal{H}_1 \otimes \dots \otimes \mathcal{H}_M$ would work similarly

Tensor product of M Hilbert spaces:

$$(a_1 \otimes \dots \otimes a_M) (h_1, \dots, h_M) = \prod_{m=1}^M \langle a_m, h_m \rangle_{\mathcal{H}_m},$$

$a_1 \otimes \dots \otimes a_M$, $\mathcal{H}_1 \otimes \dots \otimes \mathcal{H}_M$ would work similarly

Tensor product of M Hilbert spaces:

$$(a_1 \otimes \dots \otimes a_M) (h_1, \dots, h_M) = \prod_{m=1}^M \langle a_m, h_m \rangle_{\mathcal{H}_m},$$
$$\left\langle \otimes_{m=1}^M a_m, \otimes_{m=1}^M h_m \right\rangle = \prod_{m=1}^M \langle a_m, h_m \rangle_{\mathcal{H}_m}.$$

\Rightarrow HSIC for M -variables: \checkmark

$\langle \cdot, \cdot \rangle$: well-defined & pos. definite [Reed and Simon, 1980]

Well-definedness: $\langle \lambda, \lambda' \rangle$ is expansion-independent, i.e.

$$\lambda_1 = \sum_{i=1}^{n_1} c_i a_i \otimes b_i = \lambda_2 = \sum_{j=1}^{n_2} c'_j a'_j \otimes b'_j,$$

$$\langle \lambda_1, \lambda' \rangle \stackrel{?}{=} \langle \lambda_2, \lambda' \rangle$$

$\langle \cdot, \cdot \rangle$: well-defined & pos. definite [Reed and Simon, 1980]

Well-definedness: $\langle \lambda, \lambda' \rangle$ is expansion-independent, i.e.

$$\lambda_1 = \sum_{i=1}^{n_1} c_i a_i \otimes b_i = \lambda_2 = \sum_{j=1}^{n_2} c'_j a'_j \otimes b'_j,$$

$$\langle \lambda_1, \lambda' \rangle \stackrel{?}{=} \langle \lambda_2, \lambda' \rangle \Leftrightarrow \underbrace{\langle \lambda_1 - \lambda_2, \lambda' \rangle}_{=0} \stackrel{?}{=} 0 \quad (\forall \lambda' \in \mathcal{L}).$$

$\langle \cdot, \cdot \rangle$: well-defined & pos. definite [Reed and Simon, 1980]

Well-definedness: $\langle \lambda, \lambda' \rangle$ is expansion-independent, i.e.

$$\lambda_1 = \sum_{i=1}^{n_1} c_i a_i \otimes b_i = \lambda_2 = \sum_{j=1}^{n_2} c'_j a'_j \otimes b'_j,$$

$$\langle \lambda_1, \lambda' \rangle \stackrel{?}{=} \langle \lambda_2, \lambda' \rangle \Leftrightarrow \underbrace{\langle \lambda_1 - \lambda_2, \lambda' \rangle}_{=0} \stackrel{?}{=} 0 \quad (\forall \lambda' \in \mathcal{L}).$$

- In other words: $v = 0 \stackrel{?}{\Rightarrow} \langle v, \lambda' \rangle = 0, \forall \lambda' \in \mathcal{L}$.

$\langle \cdot, \cdot \rangle$: well-defined & pos. definite [Reed and Simon, 1980]

Well-definedness: $\langle \lambda, \lambda' \rangle$ is expansion-independent, i.e.

$$\lambda_1 = \sum_{i=1}^{n_1} c_i a_i \otimes b_i = \lambda_2 = \sum_{j=1}^{n_2} c'_j a'_j \otimes b'_j,$$

$$\langle \lambda_1, \lambda' \rangle \stackrel{?}{=} \langle \lambda_2, \lambda' \rangle \Leftrightarrow \underbrace{\langle \lambda_1 - \lambda_2, \lambda' \rangle}_{=0} \stackrel{?}{=} 0 \quad (\forall \lambda' \in \mathcal{L}).$$

- In other words: $v = 0 \stackrel{?}{\Rightarrow} \langle v, \lambda' \rangle = 0, \forall \lambda' \in \mathcal{L}$.
- $\lambda' := \sum_{i=1}^{n'} d_i e_i \otimes f_i$,

$$\begin{aligned} \langle 0, \lambda' \rangle &= \left\langle 0, \sum_{i=1}^{n'} d_i e_i \otimes f_i \right\rangle = \sum_{i=1}^{n'} d_i \underbrace{\langle 0, e_i \otimes f_i \rangle}_{\substack{=0 \\ (e_i, f_i)=0}} = 0. \end{aligned}$$

0 form

$\langle \cdot, \cdot \rangle$ is positive definite

- Goal: $\langle \lambda, \lambda \rangle = 0 \Rightarrow \lambda = 0.$

$\langle \cdot, \cdot \rangle$ is positive definite

- Goal: $\langle \lambda, \lambda \rangle = 0 \Rightarrow \lambda = 0.$
- $\lambda := \sum_{i=1}^n c_i \color{red}{a_i} \otimes \color{blue}{b_i}, \color{red}{A} := \text{span}\{(a_i)\} \subset \mathcal{H}_1, \color{blue}{B} := \text{span}\{(b_i)\} \subset \mathcal{H}_2.$

$\langle \cdot, \cdot \rangle$ is positive definite

- Goal: $\langle \lambda, \lambda \rangle = 0 \Rightarrow \lambda = 0.$
- $\lambda := \sum_{i=1}^n c_i \textcolor{red}{a}_i \otimes \textcolor{blue}{b}_i,$ $A := \text{span}\{(a_i)\} \subset \mathcal{H}_1,$ $B := \text{span}\{(b_i)\} \subset \mathcal{H}_2.$
- $(\alpha_i) := \text{ONB for } A,$ $(\beta_j) := \text{ONB for } B.$

$\langle \cdot, \cdot \rangle$ is positive definite

- Goal: $\langle \lambda, \lambda \rangle = 0 \Rightarrow \lambda = 0$.
- $\lambda := \sum_{i=1}^n c_i \mathbf{a}_i \otimes \mathbf{b}_i$, $A := \text{span}\{\mathbf{a}_i\} \subset \mathcal{H}_1$, $B := \text{span}\{\mathbf{b}_i\} \subset \mathcal{H}_2$.
- (α_i) := ONB for A , (β_j) := ONB for B .
- $a_i \in A$, $b_i \in B$ hence

$$\lambda = \sum_{i,j} c_{ij} \alpha_i \otimes \beta_j,$$

$$\langle \lambda, \lambda \rangle = \left\langle \sum_{i,j} c_{ij} \alpha_i \otimes \beta_j, \sum_{u,v} c_{uv} \alpha_u \otimes \beta_v \right\rangle$$

$\langle \cdot, \cdot \rangle$ is positive definite

- Goal: $\langle \lambda, \lambda \rangle = 0 \Rightarrow \lambda = 0$.
- $\lambda := \sum_{i=1}^n c_i \mathbf{a}_i \otimes \mathbf{b}_i$, $A := \text{span}\{\mathbf{a}_i\} \subset \mathcal{H}_1$, $B := \text{span}\{\mathbf{b}_i\} \subset \mathcal{H}_2$.
- (α_i) := ONB for A , (β_j) := ONB for B .
- $a_i \in A$, $b_i \in B$ hence

$$\lambda = \sum_{i,j} c_{ij} \alpha_i \otimes \beta_j,$$

$$\begin{aligned}\langle \lambda, \lambda \rangle &= \left\langle \sum_{i,j} c_{ij} \alpha_i \otimes \beta_j, \sum_{u,v} c_{uv} \alpha_u \otimes \beta_v \right\rangle \\ &= \sum_{i,j,u,v} c_{ij} c_{uv} \underbrace{\langle \alpha_i \otimes \beta_j, \alpha_u \otimes \beta_v \rangle}_{\langle \alpha_i, \alpha_u \rangle_{\mathcal{H}_1} \langle \beta_j, \beta_v \rangle_{\mathcal{H}_2}} = \delta_{iu} \delta_{jv}\end{aligned}$$

$\langle \cdot, \cdot \rangle$ is positive definite

- Goal: $\langle \lambda, \lambda \rangle = 0 \Rightarrow \lambda = 0.$
- $\lambda := \sum_{i=1}^n c_i \mathbf{a}_i \otimes \mathbf{b}_i$, $A := \text{span}\{(a_i)\} \subset \mathcal{H}_1$, $B := \text{span}\{(b_i)\} \subset \mathcal{H}_2$.
- (α_i) := ONB for A , (β_j) := ONB for B .
- $a_i \in A$, $b_i \in B$ hence

$$\lambda = \sum_{i,j} c_{ij} \alpha_i \otimes \beta_j,$$

$$\langle \lambda, \lambda \rangle = \left\langle \sum_{i,j} c_{ij} \alpha_i \otimes \beta_j, \sum_{u,v} c_{uv} \alpha_u \otimes \beta_v \right\rangle$$

$$= \sum_{i,j,u,v} c_{ij} c_{uv} \underbrace{\langle \alpha_i \otimes \beta_j, \alpha_u \otimes \beta_v \rangle}_{\langle \alpha_i, \alpha_u \rangle_{\mathcal{H}_1} \langle \beta_j, \beta_v \rangle_{\mathcal{H}_2} = \delta_{iu} \delta_{jv}} = \sum_{i,j} c_{ij}^2.$$

$\langle \cdot, \cdot \rangle$ is positive definite

- Goal: $\langle \lambda, \lambda \rangle = 0 \Rightarrow \lambda = 0.$
- $\lambda := \sum_{i=1}^n c_i \mathbf{a}_i \otimes \mathbf{b}_i$, $A := \text{span}\{\mathbf{a}_i\} \subset \mathcal{H}_1$, $B := \text{span}\{\mathbf{b}_i\} \subset \mathcal{H}_2$.
- $(\alpha_i) := \text{ONB for } A$, $(\beta_j) := \text{ONB for } B$.
- $a_i \in A$, $b_i \in B$ hence

$$\lambda = \sum_{i,j} c_{ij} \alpha_i \otimes \beta_j,$$

$$\begin{aligned}\langle \lambda, \lambda \rangle &= \left\langle \sum_{i,j} c_{ij} \alpha_i \otimes \beta_j, \sum_{u,v} c_{uv} \alpha_u \otimes \beta_v \right\rangle \\ &= \sum_{i,j,u,v} c_{ij} c_{uv} \underbrace{\langle \alpha_i \otimes \beta_j, \alpha_u \otimes \beta_v \rangle}_{\langle \alpha_i, \alpha_u \rangle_{\mathcal{H}_1} \langle \beta_j, \beta_v \rangle_{\mathcal{H}_2} = \delta_{iu} \delta_{jv}} = \sum_{i,j} c_{ij}^2.\end{aligned}$$

- In short, $\langle \lambda, \lambda \rangle = 0 \Rightarrow c_{ij} = 0 \ (\forall i, j)$, i.e. $\lambda = 0.$

Theorem ([Berlinet and Thomas-Agnan, 2004])

- Given: $\mathcal{H}_1 = \mathcal{H}_k$, $\mathcal{H}_2 = \mathcal{H}_\ell$ RKHSs with kernel k and ℓ .
- Then $\mathcal{H}_1 \otimes \mathcal{H}_2$ is RKHS with kernel

$$k \otimes \ell : (\mathcal{X} \times \mathcal{Y}) \times (\mathcal{X} \times \mathcal{Y}) \rightarrow \mathbb{R},$$

$$(k \otimes \ell)((x_1, y_1), (x_2, y_2)) := k(x_1, x_2)\ell(y_1, y_2).$$

Tensor product of RKHSs

Theorem ([Berlinet and Thomas-Agnan, 2004])

- Given: $\mathcal{H}_1 = \mathcal{H}_k$, $\mathcal{H}_2 = \mathcal{H}_\ell$ RKHSs with kernel k and ℓ .
- Then $\mathcal{H}_1 \otimes \mathcal{H}_2$ is RKHS with kernel

$$\begin{aligned} k \otimes \ell : (\mathcal{X} \times \mathcal{Y}) \times (\mathcal{X} \times \mathcal{Y}) &\rightarrow \mathbb{R}, \\ (k \otimes \ell)((x_1, y_1), (x_2, y_2)) &:= k(x_1, x_2)\ell(y_1, y_2). \end{aligned}$$

Intuition:

- inner product on \mathcal{X} and $\mathcal{Y} \rightarrow$ inner product on $\mathcal{X} \times \mathcal{Y}$.
- $\mathcal{X} =$ animal images, $\mathcal{Y} =$ descriptions of animals.

Until now

- $a \otimes b$: defined; 'nice' operator (HS:=Hilbert-Schmidt).

Until now

- $a \otimes b$: defined; 'nice' operator (HS:=Hilbert-Schmidt).
- It will descend to its expectation ($C_{xy}^u \Rightarrow$ HSIC).

Until now

- $a \otimes b$: defined; 'nice' operator (HS:=Hilbert-Schmidt).
- It will descend to its expectation ($C_{xy}^u \Rightarrow$ HSIC).
- $(e_i), (f_j)$: canonical basis in $\mathbb{R}^{d_1}, \mathbb{R}^{d_2}$.
- HS operators: extensions of $L \in \mathbb{R}^{d_2 \times d_1}$ with

$$\|L\|_F^2 = \sum_i \| \underbrace{Le_i}_{i^{th} \text{ column of } L} \|_2^2$$

Until now

- $a \otimes b$: defined; 'nice' operator (HS:=Hilbert-Schmidt).
- It will descend to its expectation ($C_{xy}^u \Rightarrow$ HSIC).
- $(e_i), (f_j)$: canonical basis in $\mathbb{R}^{d_1}, \mathbb{R}^{d_2}$.
- HS operators: extensions of $L \in \mathbb{R}^{d_2 \times d_1}$ with

$$\|L\|_F^2 = \sum_i \| \underbrace{Le_i}_{i^{th} \text{ column of } L} \|_2^2 = \sum_i \sum_j (Le_i)_j^2$$

Until now

- $a \otimes b$: defined; 'nice' operator (HS:=Hilbert-Schmidt).
- It will descend to its expectation ($C_{xy}^u \Rightarrow$ HSIC).
- $(e_i), (f_j)$: canonical basis in $\mathbb{R}^{d_1}, \mathbb{R}^{d_2}$.
- HS operators: extensions of $L \in \mathbb{R}^{d_2 \times d_1}$ with

$$\|L\|_F^2 = \sum_i \| \underbrace{Le_i}_{i^{th} \text{ column of } L} \|_2^2 = \sum_i \sum_j (Le_i)_j^2 = \sum_{i,j} L_{ji}^2.$$

Hilbert-Schmidt operators: quick summary

- An $L : \mathcal{H}_1 \rightarrow \mathcal{H}_2$ bounded linear operator is called Hilbert-Schmidt if

$$\|L\|_{HS}^2 := \sum_i \underbrace{\|Le_i\|_{\mathcal{H}_2}^2}_{=\sum_j \langle Le_i, f_j \rangle_{\mathcal{H}_2}^2} < \infty.$$

Hilbert-Schmidt operators: quick summary

- $\mathcal{H}_1, \mathcal{H}_2$: separable Hilbert spaces. $(e_i)_{i \in I}, (f_j)_{j \in J}$: ONB in $\mathcal{H}_1, \mathcal{H}_2$.
- An $L : \mathcal{H}_1 \rightarrow \mathcal{H}_2$ bounded linear operator is called **Hilbert-Schmidt** if

$$\|L\|_{HS}^2 := \sum_i \underbrace{\|Le_i\|_{\mathcal{H}_2}^2}_{=\sum_j \langle Le_i, f_j \rangle_{\mathcal{H}_2}^2} < \infty.$$

Hilbert-Schmidt operators: quick summary

- $\mathcal{H}_1, \mathcal{H}_2$: separable Hilbert spaces. $(e_i)_{i \in I}, (f_j)_{j \in J}$: ONB in $\mathcal{H}_1, \mathcal{H}_2$.
- An $L : \mathcal{H}_1 \rightarrow \mathcal{H}_2$ bounded linear **operator** is called **Hilbert-Schmidt** if

$$\|L\|_{HS}^2 := \sum_i \underbrace{\|Le_i\|_{\mathcal{H}_2}^2}_{=\sum_j \langle Le_i, f_j \rangle_{\mathcal{H}_2}^2} < \infty.$$

- Inner product on $L_1, L_2 \in HS(\mathcal{H}_1, \mathcal{H}_2)$

$$\langle L_1, L_2 \rangle_{HS} := \sum_i \langle L_1 e_i, L_2 e_i \rangle_{\mathcal{H}_2}.$$

Hilbert-Schmidt operators: quick summary

- $\mathcal{H}_1, \mathcal{H}_2$: separable Hilbert spaces. $(e_i)_{i \in I}, (f_j)_{j \in J}$: ONB in $\mathcal{H}_1, \mathcal{H}_2$.
- An $L : \mathcal{H}_1 \rightarrow \mathcal{H}_2$ bounded linear operator is called **Hilbert-Schmidt** if

$$\|L\|_{HS}^2 := \sum_i \underbrace{\|Le_i\|_{\mathcal{H}_2}^2}_{=\sum_j \langle Le_i, f_j \rangle_{\mathcal{H}_2}^2} < \infty.$$

- Inner product on $L_1, L_2 \in HS(\mathcal{H}_1, \mathcal{H}_2)$

$$\langle L_1, L_2 \rangle_{HS} := \sum_i \langle L_1 e_i, L_2 e_i \rangle_{\mathcal{H}_2}.$$

- $HS(\mathcal{H}_1, \mathcal{H}_2)$: **Hilbert space**.

Hilbert-Schmidt operators: notes

- $\mathcal{H}_1, \mathcal{H}_2$: separable $\Rightarrow I, J$: countable, i.e. 'sums'.

Hilbert-Schmidt operators: notes

- $\mathcal{H}_1, \mathcal{H}_2$: separable $\Rightarrow I, J$: countable, i.e. 'sums'.
- $\langle L_1, L_2 \rangle_{HS}$: well-defined (independent of the chosen basis).

Hilbert-Schmidt operators: notes

- $\mathcal{H}_1, \mathcal{H}_2$: separable $\Rightarrow I, J$: countable, i.e. 'sums'.
- $\langle L_1, L_2 \rangle_{HS}$: well-defined (independent of the chosen basis).
- For RKHSs (\mathcal{H}_i): \mathcal{X} : separable, k : continuous $\Rightarrow \mathcal{H}_k$: separable [Steinwart and Christmann, 2008].

For $a \otimes b$ with $a \in \mathcal{H}_1$, $b \in \mathcal{H}_2$:

- linearity: ✓
- boundedness ($c \in \mathcal{H}_2$):

$$\|(a \otimes b)c\|_{\mathcal{H}_1} = \|a \langle b, c \rangle_{\mathcal{H}_2}\|_{\mathcal{H}_1}$$

For $a \otimes b$ with $a \in \mathcal{H}_1$, $b \in \mathcal{H}_2$:

- linearity: ✓
- boundedness ($c \in \mathcal{H}_2$):

$$\|(a \otimes b)c\|_{\mathcal{H}_1} = \|a \langle b, c \rangle_{\mathcal{H}_2}\|_{\mathcal{H}_1} = |\langle b, c \rangle_{\mathcal{H}_2}| \|a\|_{\mathcal{H}_1}$$

$a \otimes b$ is Hilbert-Schmidt: linear & bounded

For $a \otimes b$ with $a \in \mathcal{H}_1$, $b \in \mathcal{H}_2$:

- linearity: ✓
- boundedness ($c \in \mathcal{H}_2$):

$$\begin{aligned}\|(a \otimes b)c\|_{\mathcal{H}_1} &= \|a \langle b, c \rangle_{\mathcal{H}_2}\|_{\mathcal{H}_1} = |\langle b, c \rangle_{\mathcal{H}_2}| \|a\|_{\mathcal{H}_1} \\ &\stackrel{\text{CBS}}{\leqslant} \|b\|_{\mathcal{H}_2} \|c\|_{\mathcal{H}_2} \|a\|_{\mathcal{H}_1}.\end{aligned}$$

For $a \otimes b$ with $a \in \mathcal{H}_1$, $b \in \mathcal{H}_2$:

- linearity: ✓
- boundedness ($c \in \mathcal{H}_2$):

$$\begin{aligned}\|(a \otimes b)c\|_{\mathcal{H}_1} &= \|a \langle b, c \rangle_{\mathcal{H}_2}\|_{\mathcal{H}_1} = |\langle b, c \rangle_{\mathcal{H}_2}| \|a\|_{\mathcal{H}_1} \\ &\stackrel{\text{CBS}}{\leqslant} \|b\|_{\mathcal{H}_2} \|c\|_{\mathcal{H}_2} \|a\|_{\mathcal{H}_1}.\end{aligned}$$

Thus $\|a \otimes b\| \leq \|a\|_{\mathcal{H}_1} \|b\|_{\mathcal{H}_2} < \infty$.

$a \otimes b$ is a Hilbert-Schmidt operator

Let $(e_i)_{i \in I} \subset \mathcal{H}_2$ ONB,

$$\|a \otimes b\|_{HS}^2 = \sum_i \left\| (a \otimes b)e_i \right\|_{\mathcal{H}_1}^2$$

$a \otimes b$ is a Hilbert-Schmidt operator

Let $(e_i)_{i \in I} \subset \mathcal{H}_2$ ONB,

$$\|a \otimes b\|_{HS}^2 = \sum_i \left\| (a \otimes b)e_i \right\|_{\mathcal{H}_1}^2 = \sum_i \underbrace{\left\| a \langle b, e_i \rangle_{\mathcal{H}_2} \right\|_{\mathcal{H}_1}^2}_{\|a\|_{\mathcal{H}_1}^2 |\langle b, e_i \rangle_{\mathcal{H}_2}|^2}$$

$a \otimes b$ is a Hilbert-Schmidt operator

Let $(e_i)_{i \in I} \subset \mathcal{H}_2$ ONB,

$$\begin{aligned}\|a \otimes b\|_{HS}^2 &= \sum_i \left\| (a \otimes b)e_i \right\|_{\mathcal{H}_1}^2 = \sum_i \underbrace{\left\| a \langle b, e_i \rangle_{\mathcal{H}_2} \right\|_{\mathcal{H}_1}^2}_{\|a\|_{\mathcal{H}_1}^2 |\langle b, e_i \rangle_{\mathcal{H}_2}|^2} \\ &= \|a\|_{\mathcal{H}_1}^2 \underbrace{\sum_i |\langle b, e_i \rangle_{\mathcal{H}_2}|^2}_{\|b\|_{\mathcal{H}_2}^2 \text{ (Parseval equality)}} < \infty.\end{aligned}$$

$a \otimes b$ is a Hilbert-Schmidt operator

Let $(e_i)_{i \in I} \subset \mathcal{H}_2$ ONB,

$$\begin{aligned}\|a \otimes b\|_{HS}^2 &= \sum_i \left\| (a \otimes b)e_i \right\|_{\mathcal{H}_1}^2 = \sum_i \underbrace{\left\| a \langle b, e_i \rangle_{\mathcal{H}_2} \right\|_{\mathcal{H}_1}^2}_{\|a\|_{\mathcal{H}_1}^2 |\langle b, e_i \rangle_{\mathcal{H}_2}|^2} \\ &= \|a\|_{\mathcal{H}_1}^2 \underbrace{\sum_i |\langle b, e_i \rangle_{\mathcal{H}_2}|^2}_{\|b\|_{\mathcal{H}_2}^2 \text{ (Parseval equality)}} < \infty.\end{aligned}$$

In short

$$\|a \otimes b\|_{HS}^2 = \|a\|_{\mathcal{H}_1}^2 \|b\|_{\mathcal{H}_2}^2.$$

Uncentered cross-covariance operator

$$C_{xy}^u := \mathbb{E}_{xy} \left[\underbrace{\varphi(x) \otimes \psi(y)}_{\in HS(\mathcal{H}_\ell, \mathcal{H}_k)} \right] \in HS(\mathcal{H}_\ell, \mathcal{H}_k).$$

Uncentered cross-covariance operator

$$C_{xy}^u := \mathbb{E}_{xy} \left[\underbrace{\varphi(x) \otimes \psi(y)}_{\in HS(\mathcal{H}_\ell, \mathcal{H}_k)} \right] \in HS(\mathcal{H}_\ell, \mathcal{H}_k).$$

- 'Same' construction as $\mu_{\mathbb{P}}$: we changed \mathcal{H}_k to $HS(\mathcal{H}_\ell, \mathcal{H}_k)$.

Uncentered cross-covariance operator

$$C_{xy}^u := \mathbb{E}_{xy} \left[\underbrace{\varphi(x) \otimes \psi(y)}_{\in HS(\mathcal{H}_\ell, \mathcal{H}_k)} \right] \in HS(\mathcal{H}_\ell, \mathcal{H}_k).$$

- 'Same' construction as $\mu_{\mathbb{P}}$: we changed \mathcal{H}_k to $HS(\mathcal{H}_\ell, \mathcal{H}_k)$.
- Bochner integral: $\exists C_{xy}^u \Leftrightarrow \mathbb{E}_{xy} \|\varphi(x) \otimes \psi(y)\|_{HS} < \infty$.

Uncentered cross-covariance operator

$$C_{xy}^u := \mathbb{E}_{xy} \left[\underbrace{\varphi(x) \otimes \psi(y)}_{\in HS(\mathcal{H}_\ell, \mathcal{H}_k)} \right] \in HS(\mathcal{H}_\ell, \mathcal{H}_k).$$

- 'Same' construction as $\mu_{\mathbb{P}}$: we changed \mathcal{H}_k to $HS(\mathcal{H}_\ell, \mathcal{H}_k)$.
- Bochner integral: $\exists C_{xy}^u \Leftrightarrow \mathbb{E}_{xy} \|\varphi(x) \otimes \psi(y)\|_{HS} < \infty$.
- $\|\varphi(x) \otimes \psi(y)\|_{HS} = \|\varphi(x)\|_{\mathcal{H}_k} \|\psi(y)\|_{\mathcal{H}_\ell}$

Uncentered cross-covariance operator

$$C_{xy}^u := \mathbb{E}_{xy} \left[\underbrace{\varphi(x) \otimes \psi(y)}_{\in HS(\mathcal{H}_\ell, \mathcal{H}_k)} \right] \in HS(\mathcal{H}_\ell, \mathcal{H}_k).$$

- 'Same' construction as $\mu_{\mathbb{P}}$: we changed \mathcal{H}_k to $HS(\mathcal{H}_\ell, \mathcal{H}_k)$.
- Bochner integral: $\exists C_{xy}^u \Leftrightarrow \mathbb{E}_{xy} \|\varphi(x) \otimes \psi(y)\|_{HS} < \infty$.
- $\|\varphi(x) \otimes \psi(y)\|_{HS} = \|\varphi(x)\|_{\mathcal{H}_k} \|\psi(y)\|_{\mathcal{H}_\ell} = \sqrt{k(x, x)} \sqrt{\ell(y, y)}$.

Uncentered cross-covariance operator

$$C_{xy}^u := \mathbb{E}_{xy} \left[\underbrace{\varphi(x) \otimes \psi(y)}_{\in HS(\mathcal{H}_\ell, \mathcal{H}_k)} \right] \in HS(\mathcal{H}_\ell, \mathcal{H}_k).$$

- 'Same' construction as $\mu_{\mathbb{P}}$: we changed \mathcal{H}_k to $HS(\mathcal{H}_\ell, \mathcal{H}_k)$.
- Bochner integral: $\exists C_{xy}^u \Leftrightarrow \mathbb{E}_{xy} \|\varphi(x) \otimes \psi(y)\|_{HS} < \infty$.
- $\|\varphi(x) \otimes \psi(y)\|_{HS} = \|\varphi(x)\|_{\mathcal{H}_k} \|\psi(y)\|_{\mathcal{H}_\ell} = \sqrt{k(x, x)} \sqrt{\ell(y, y)}$.
- Sufficient condition: k and ℓ are bounded.

Centered covariance operator [Baker, 1973]

Let $\mu_x := \mu_{\mathbb{P}_x}$, $\mu_y := \mu_{\mathbb{P}_y}$. $\mathbb{P}_x, \mathbb{P}_y$: marginals of \mathbb{P}_{xy} .

$$C_{xy}^c = \mathbb{E}_{xy} \left[\left(\varphi(x) - \underbrace{\mathbb{E}_x \varphi(x)}_{\mu_x} \right) \otimes \left(\psi(y) - \underbrace{\mathbb{E}_y \psi(y)}_{\mu_y} \right) \right]$$

Centered covariance operator [Baker, 1973]

Let $\mu_x := \mu_{\mathbb{P}_x}$, $\mu_y := \mu_{\mathbb{P}_y}$. $\mathbb{P}_x, \mathbb{P}_y$: marginals of \mathbb{P}_{xy} .

$$\begin{aligned} C_{xy}^c &= \mathbb{E}_{xy} \left[\left(\varphi(x) - \underbrace{\mathbb{E}_x \varphi(x)}_{\mu_x} \right) \otimes \left(\psi(y) - \underbrace{\mathbb{E}_y \psi(y)}_{\mu_y} \right) \right] \\ &= \underbrace{\mathbb{E}_{xy} [\varphi(x) \otimes \psi(y)]}_{C_{xy}^u \in HS(\mathcal{H}_\ell, \mathcal{H}_k)} - \underbrace{\mu_x \otimes \mu_y}_{\in HS(\mathcal{H}_\ell, \mathcal{H}_k)} \in HS(\mathcal{H}_\ell, \mathcal{H}_k). \end{aligned}$$

Hilbert-Schmidt independence criterion (HSIC)

HSIC [Fukumizu et al., 2004, Gretton et al., 2005a]:

$$\text{HSIC}(x, y; \mathcal{H}_k, \mathcal{H}_\ell) := \|C_{xy}^c\|_{HS}.$$

Hilbert-Schmidt independence criterion (HSIC)

HSIC [Fukumizu et al., 2004, Gretton et al., 2005a]:

$$\text{HSIC}(x, y; \mathcal{H}_k, \mathcal{H}_\ell) := \|C_{xy}^c\|_{HS}.$$

Question

When does HSIC characterize independence?

Hilbert-Schmidt independence criterion (HSIC)

HSIC [Fukumizu et al., 2004, Gretton et al., 2005a]:

$$\text{HSIC}(x, y; \mathcal{H}_k, \mathcal{H}_\ell) := \|C_{xy}^c\|_{HS}.$$

Question

When does HSIC characterize independence?

We will discuss it later (after $\text{HSIC} \Leftrightarrow$ distance covariance).

How do covariance operators encode covariance?

Let $g \in \mathcal{H}_\ell$, $f \in \mathcal{H}_k$, $HS := HS(\mathcal{H}_\ell, \mathcal{H}_k)$.

$$\langle f, C_{xy}^u g \rangle_{\mathcal{H}_k} = \langle C_{xy}^u, f \otimes g \rangle_{HS}$$

Cheating:

- next slide.
- Enough $f \in \mathcal{H}_1$, $g \in \mathcal{H}_2$, $L \in HS(\mathcal{H}_2, \mathcal{H}_1)$

$$\langle f, Lg \rangle_{\mathcal{H}_1} = \langle L, f \otimes g \rangle_{HS(\mathcal{H}_2, \mathcal{H}_1)}$$

How do covariance operators encode covariance?

Let $g \in \mathcal{H}_\ell$, $f \in \mathcal{H}_k$, $HS := HS(\mathcal{H}_\ell, \mathcal{H}_k)$.

$$\langle f, C_{xy}^u g \rangle_{\mathcal{H}_k} = \langle C_{xy}^u, f \otimes g \rangle_{HS} = \langle \mathbb{E}_{xy}[\varphi(x) \otimes \psi(y)], f \otimes g \rangle_{HS}$$

Cheating:

- next slide.
- Enough $f \in \mathcal{H}_1$, $g \in \mathcal{H}_2$, $L \in HS(\mathcal{H}_2, \mathcal{H}_1)$

$$\langle f, Lg \rangle_{\mathcal{H}_1} = \langle L, f \otimes g \rangle_{HS(\mathcal{H}_2, \mathcal{H}_1)}$$

How do covariance operators encode covariance?

Let $g \in \mathcal{H}_\ell$, $f \in \mathcal{H}_k$, $HS := HS(\mathcal{H}_\ell, \mathcal{H}_k)$.

$$\begin{aligned}\langle f, C_{xy}^u g \rangle_{\mathcal{H}_k} &= \langle C_{xy}^u, f \otimes g \rangle_{HS} = \langle \mathbb{E}_{xy}[\varphi(x) \otimes \psi(y)], f \otimes g \rangle_{HS} \\ &= \mathbb{E}_{xy} \underbrace{\langle \varphi(x) \otimes \psi(y), f \otimes g \rangle_{HS}}_{=f(x)g(y)}\end{aligned}$$

Cheating:

- next slide.
- Enough $f \in \mathcal{H}_1$, $g \in \mathcal{H}_2$, $L \in HS(\mathcal{H}_2, \mathcal{H}_1)$

$$\langle f, Lg \rangle_{\mathcal{H}_1} = \langle L, f \otimes g \rangle_{HS(\mathcal{H}_2, \mathcal{H}_1)}$$

How do covariance operators encode covariance?

Let $g \in \mathcal{H}_\ell$, $f \in \mathcal{H}_k$, $HS := HS(\mathcal{H}_\ell, \mathcal{H}_k)$.

$$\begin{aligned}\langle f, C_{xy}^u g \rangle_{\mathcal{H}_k} &= \langle C_{xy}^u, f \otimes g \rangle_{HS} = \langle \mathbb{E}_{xy}[\varphi(x) \otimes \psi(y)], f \otimes g \rangle_{HS} \\ &= \mathbb{E}_{xy} \underbrace{\langle \varphi(x) \otimes \psi(y), f \otimes g \rangle_{HS}}_{=f(x)g(y)} = \mathbb{E}_{xy} [f(x)g(y)].\end{aligned}$$

Cheating:

- next slide.
- Enough $f \in \mathcal{H}_1$, $g \in \mathcal{H}_2$, $L \in HS(\mathcal{H}_2, \mathcal{H}_1)$

$$\langle f, Lg \rangle_{\mathcal{H}_1} = \langle L, f \otimes g \rangle_{HS(\mathcal{H}_2, \mathcal{H}_1)}$$

Cheating

Statement: with $f \in \mathcal{H}_1$, $g \in \mathcal{H}_2$, $L \in HS(\mathcal{H}_2, \mathcal{H}_1)$

$$\langle f, Lg \rangle_{\mathcal{H}_1} = \langle L, f \otimes g \rangle_{HS(\mathcal{H}_2, \mathcal{H}_1)}.$$

Proof: $(b_i)_{i \in I}$ ONB in \mathcal{H}_2 ,

$$\langle f, Lg \rangle_{\mathcal{H}_1} = \left\langle f, L \sum_i \langle g, b_i \rangle_{\mathcal{H}_2} b_i \right\rangle_{\mathcal{H}_1}$$

Cheating

Statement: with $f \in \mathcal{H}_1$, $g \in \mathcal{H}_2$, $L \in HS(\mathcal{H}_2, \mathcal{H}_1)$

$$\langle f, Lg \rangle_{\mathcal{H}_1} = \langle L, f \otimes g \rangle_{HS(\mathcal{H}_2, \mathcal{H}_1)}.$$

Proof: $(b_i)_{i \in I}$ ONB in \mathcal{H}_2 ,

$$\langle f, Lg \rangle_{\mathcal{H}_1} = \left\langle f, L \sum_i \langle g, b_i \rangle_{\mathcal{H}_2} b_i \right\rangle_{\mathcal{H}_1} = \sum_i \langle g, b_i \rangle_{\mathcal{H}_2} \langle f, Lb_i \rangle_{\mathcal{H}_1}$$

Cheating

Statement: with $f \in \mathcal{H}_1$, $g \in \mathcal{H}_2$, $L \in HS(\mathcal{H}_2, \mathcal{H}_1)$

$$\langle f, Lg \rangle_{\mathcal{H}_1} = \langle L, f \otimes g \rangle_{HS(\mathcal{H}_2, \mathcal{H}_1)}.$$

Proof: $(b_i)_{i \in I}$ ONB in \mathcal{H}_2 ,

$$\begin{aligned}\langle f, Lg \rangle_{\mathcal{H}_1} &= \left\langle f, L \sum_i \langle g, b_i \rangle_{\mathcal{H}_2} b_i \right\rangle_{\mathcal{H}_1} = \sum_i \langle g, b_i \rangle_{\mathcal{H}_2} \langle f, Lb_i \rangle_{\mathcal{H}_1} \\ &\stackrel{\otimes}{=} \sum_i \langle Lb_i, (f \otimes g)b_i \rangle_{\mathcal{H}_1}\end{aligned}$$

Cheating

Statement: with $f \in \mathcal{H}_1$, $g \in \mathcal{H}_2$, $L \in HS(\mathcal{H}_2, \mathcal{H}_1)$

$$\langle f, Lg \rangle_{\mathcal{H}_1} = \langle L, f \otimes g \rangle_{HS(\mathcal{H}_2, \mathcal{H}_1)}.$$

Proof: $(b_i)_{i \in I}$ ONB in \mathcal{H}_2 ,

$$\begin{aligned}\langle f, Lg \rangle_{\mathcal{H}_1} &= \left\langle f, L \sum_i \langle g, b_i \rangle_{\mathcal{H}_2} b_i \right\rangle_{\mathcal{H}_1} = \sum_i \langle g, b_i \rangle_{\mathcal{H}_2} \langle f, Lb_i \rangle_{\mathcal{H}_1} \\ &\stackrel{\otimes}{=} \sum_i \langle Lb_i, (f \otimes g)b_i \rangle_{\mathcal{H}_1} \stackrel{\langle \cdot, \cdot \rangle_{HS}}{=} \langle L, f \otimes g \rangle_{HS}.\end{aligned}$$

Cheating

Statement: with $f \in \mathcal{H}_1$, $g \in \mathcal{H}_2$, $L \in HS(\mathcal{H}_2, \mathcal{H}_1)$

$$\langle f, Lg \rangle_{\mathcal{H}_1} = \langle L, f \otimes g \rangle_{HS(\mathcal{H}_2, \mathcal{H}_1)}.$$

With $L := a \otimes b$

$$\langle a \otimes b, f \otimes g \rangle_{HS(\mathcal{H}_2, \mathcal{H}_1)} = \langle f, (a \otimes b)g \rangle_{\mathcal{H}_1}$$

Cheating

Statement: with $f \in \mathcal{H}_1$, $g \in \mathcal{H}_2$, $L \in HS(\mathcal{H}_2, \mathcal{H}_1)$

$$\langle f, Lg \rangle_{\mathcal{H}_1} = \langle L, f \otimes g \rangle_{HS(\mathcal{H}_2, \mathcal{H}_1)}.$$

With $L := a \otimes b$

$$\langle a \otimes b, f \otimes g \rangle_{HS(\mathcal{H}_2, \mathcal{H}_1)} = \langle f, (a \otimes b)g \rangle_{\mathcal{H}_1} \stackrel{\otimes}{=} \langle a, f \rangle_{\mathcal{H}_1} \langle b, g \rangle_{\mathcal{H}_2}.$$

Remember: we have seen this for $a = f$, $b = g$ (when proving $a \otimes b$ is HS).

Effect of the centered cross-covariance operator

Using that $C_{xy}^c = C_{xy}^u - \mu_x \otimes \mu_y$

$$\langle f, C_{xy}^c g \rangle_{\mathcal{H}_k} = \langle f, C_{xy}^u g \rangle_{\mathcal{H}_k} - \langle f, (\mu_x \otimes \mu_y) g \rangle_{\mathcal{H}_k}$$

Effect of the centered cross-covariance operator

Using that $C_{xy}^c = C_{xy}^u - \mu_x \otimes \mu_y$

$$\begin{aligned}\langle f, C_{xy}^c g \rangle_{\mathcal{H}_k} &= \langle f, C_{xy}^u g \rangle_{\mathcal{H}_k} - \langle f, (\mu_x \otimes \mu_y) g \rangle_{\mathcal{H}_k} \\ &\stackrel{\otimes}{=} \mathbb{E}_{xy}[f(x)g(y)] - \underbrace{\langle f, \mu_x \rangle_{\mathcal{H}_k}}_{\mathbb{E}_x f(x)} \underbrace{\langle g, \mu_y \rangle_{\mathcal{H}_\ell}}_{\mathbb{E}_y g(y)}\end{aligned}$$

Effect of the centered cross-covariance operator

Using that $C_{xy}^c = C_{xy}^u - \mu_x \otimes \mu_y$

$$\begin{aligned}\langle f, C_{xy}^c g \rangle_{\mathcal{H}_k} &= \langle f, C_{xy}^u g \rangle_{\mathcal{H}_k} - \langle f, (\mu_x \otimes \mu_y) g \rangle_{\mathcal{H}_k} \\ &\stackrel{\otimes}{=} \mathbb{E}_{xy}[f(x)g(y)] - \underbrace{\langle f, \mu_x \rangle_{\mathcal{H}_k}}_{\mathbb{E}_x f(x)} \underbrace{\langle g, \mu_y \rangle_{\mathcal{H}_\ell}}_{\mathbb{E}_y g(y)} \\ &= \text{cov}(f(x), g(y)).\end{aligned}$$

Three notes

- KCCA formulation: using C_{xy}^c , C_{xx}^c , C_{yy}^c .

Three notes

- KCCA formulation: using C_{xy}^c , C_{xx}^c , C_{yy}^c .
- HSIC: captures $\mathbb{P}_{xy} \stackrel{?}{=} \mathbb{P}_x \otimes \mathbb{P}_y$ in $\mathcal{H}_k \otimes \mathcal{H}_\ell$.

Three notes

- KCCA formulation: using C_{xy}^c , C_{xx}^c , C_{yy}^c .
- HSIC: captures $\mathbb{P}_{xy} \stackrel{?}{=} \mathbb{P}_x \otimes \mathbb{P}_y$ in $\mathcal{H}_k \otimes \mathcal{H}_\ell$.
- Link to distance covariance, energy distance.

In other words, ...

KCCA formulation with cross-covariance operators

$$\rho_{\text{KCCA}}(x, y) = \sup_{f \in \mathcal{H}_k, g \in \mathcal{H}_\ell} \text{corr}(f(x), g(y)) \Leftrightarrow$$
$$\sup_{f \in \mathcal{H}_k, g \in \mathcal{H}_\ell} \langle f, \mathcal{C}_{xy}^c g \rangle_{\mathcal{H}_k} \text{ s.t. } \begin{cases} \langle f, \mathcal{C}_{xx}^c f \rangle_{\mathcal{H}_k} &= 1, \\ \langle g, \mathcal{C}_{yy}^c g \rangle_{\mathcal{H}_\ell} &= 1. \end{cases}$$

KCCA: with κ -regularization

$$\rho_{\text{KCCA}}(x, y, \kappa) = \sup_{f \in \mathcal{H}_k, g \in \mathcal{H}_\ell} \text{corr}(f(x), g(y); \kappa),$$

$$\text{corr}(f(x), g(y); \kappa) = \frac{\text{cov}_{xy}(f(x), g(y))}{\sqrt{\text{var}_x f(x) + \kappa \|f\|_{\mathcal{H}_k}^2} \sqrt{\text{var}_y g(y) + \kappa \|g\|_{\mathcal{H}_\ell}^2}}.$$

KCCA: with κ -regularization

$$\rho_{\text{KCCA}}(x, y, \kappa) = \sup_{f \in \mathcal{H}_k, g \in \mathcal{H}_\ell} \text{corr}(f(x), g(y); \kappa),$$

$$\text{corr}(f(x), g(y); \kappa) = \frac{\text{cov}_{xy}(f(x), g(y))}{\sqrt{\text{var}_x f(x) + \kappa \|f\|_{\mathcal{H}_k}^2} \sqrt{\text{var}_y g(y) + \kappa \|g\|_{\mathcal{H}_\ell}^2}}.$$

Empirically,

$$\sup_{f \in \mathcal{H}_k, g \in \mathcal{H}_\ell} \left\langle f, \widehat{C_{xy}^c} g \right\rangle_{\mathcal{H}_k} \text{ s.t. } \begin{cases} \left\langle f, \left(\widehat{C_{xx}^c} + \kappa I \right) f \right\rangle_{\mathcal{H}_k} = 1, \\ \left\langle g, \left(\widehat{C_{yy}^c} + \kappa I \right) g \right\rangle_{\mathcal{H}_\ell} = 1. \end{cases}$$

KCCA: with κ -regularization

$$\rho_{\text{KCCA}}(x, y, \kappa) = \sup_{f \in \mathcal{H}_k, g \in \mathcal{H}_\ell} \text{corr}(f(x), g(y); \kappa),$$

$$\text{corr}(f(x), g(y); \kappa) = \frac{\text{cov}_{xy}(f(x), g(y))}{\sqrt{\text{var}_x f(x) + \kappa \|f\|_{\mathcal{H}_k}^2} \sqrt{\text{var}_y g(y) + \kappa \|g\|_{\mathcal{H}_\ell}^2}}.$$

Empirically,

$$\sup_{f \in \mathcal{H}_k, g \in \mathcal{H}_\ell} \left\langle f, \widehat{C_{xy}^c} g \right\rangle_{\mathcal{H}_k} \text{ s.t. } \begin{cases} \left\langle f, \left(\widehat{C_{xx}^c} + \kappa I \right) f \right\rangle_{\mathcal{H}_k} = 1, \\ \left\langle g, \left(\widehat{C_{yy}^c} + \kappa I \right) g \right\rangle_{\mathcal{H}_\ell} = 1. \end{cases}$$

KCCA consistency analysis [Fukumizu et al., 2007]

using this formulation & the convergence of $\widehat{C_{xy}^c}$, $\widehat{C_{xx}^c}$, $\widehat{C_{yy}^c}$.

HSIC: $\mathbb{P}_{xy} \stackrel{?}{=} \mathbb{P}_x \otimes \mathbb{P}_y$ in $\mathcal{H}_k \otimes \mathcal{H}_\ell$

We saw $h((x, y), (x', y')) = k(x, x')\ell(y, y')$ is the kernel of $\mathcal{H}_k \otimes \mathcal{H}_\ell$. Let

$$\|\mu_{\mathbb{P}_{xy}} - \mu_{\mathbb{P}_x \otimes \mathbb{P}_y}\|_{\mathcal{H}_h}$$

HSIC: $\mathbb{P}_{xy} \stackrel{?}{=} \mathbb{P}_x \otimes \mathbb{P}_y$ in $\mathcal{H}_k \otimes \mathcal{H}_\ell$

We saw $h((x, y), (x', y')) = k(x, x')\ell(y, y')$ is the kernel of $\mathcal{H}_k \otimes \mathcal{H}_\ell$. Let

$$\|\mu_{\mathbb{P}_{xy}} - \mu_{\mathbb{P}_x \otimes \mathbb{P}_y}\|_{\mathcal{H}_h} = \left\| \underbrace{\mathbb{E}_{\mathbb{P}_{xy}} [k(\cdot_1, x)\ell(\cdot_2, y)]}_{\mathbb{E}_{xy}[\varphi(x) \otimes \psi(y)]} - \underbrace{\mathbb{E}_{\mathbb{P}_x \otimes \mathbb{P}_y} [k(\cdot_1, x)\ell(\cdot_2, y)]}_{\mathbb{E}_x k(\cdot, x) \otimes \mathbb{E}_y \ell(\cdot, y)} \right\|_{\mathcal{H}_h}$$

HSIC: $\mathbb{P}_{xy} \stackrel{?}{=} \mathbb{P}_x \otimes \mathbb{P}_y$ in $\mathcal{H}_k \otimes \mathcal{H}_\ell$

We saw $h((x, y), (x', y')) = k(x, x')\ell(y, y')$ is the kernel of $\mathcal{H}_k \otimes \mathcal{H}_\ell$. Let

$$\begin{aligned}\|\mu_{\mathbb{P}_{xy}} - \mu_{\mathbb{P}_x \otimes \mathbb{P}_y}\|_{\mathcal{H}_h} &= \left\| \underbrace{\mathbb{E}_{\mathbb{P}_{xy}} [k(\cdot_1, x)\ell(\cdot_2, y)]}_{\mathbb{E}_{xy}[\varphi(x) \otimes \psi(y)]} - \underbrace{\mathbb{E}_{\mathbb{P}_x \otimes \mathbb{P}_y} [k(\cdot_1, x)\ell(\cdot_2, y)]}_{\mathbb{E}_x k(\cdot, x) \otimes \mathbb{E}_y \ell(\cdot, y)} \right\|_{\mathcal{H}_h} \\ &= \|\mathbb{E}_{xy} [\varphi(x) \otimes \psi(y)] - \mu_x \otimes \mu_y\|_{\mathcal{H}_h}\end{aligned}$$

HSIC: $\mathbb{P}_{xy} \stackrel{?}{=} \mathbb{P}_x \otimes \mathbb{P}_y$ in $\mathcal{H}_k \otimes \mathcal{H}_\ell$

We saw $h((x, y), (x', y')) = k(x, x')\ell(y, y')$ is the kernel of $\mathcal{H}_k \otimes \mathcal{H}_\ell$. Let

$$\begin{aligned}\|\mu_{\mathbb{P}_{xy}} - \mu_{\mathbb{P}_x \otimes \mathbb{P}_y}\|_{\mathcal{H}_h} &= \left\| \underbrace{\mathbb{E}_{\mathbb{P}_{xy}} [k(\cdot_1, x)\ell(\cdot_2, y)]}_{\mathbb{E}_{xy}[\varphi(x) \otimes \psi(y)]} - \underbrace{\mathbb{E}_{\mathbb{P}_x \otimes \mathbb{P}_y} [k(\cdot_1, x)\ell(\cdot_2, y)]}_{\mathbb{E}_x k(\cdot, x) \otimes \mathbb{E}_y \ell(\cdot, y)} \right\|_{\mathcal{H}_h} \\ &= \|\mathbb{E}_{xy} [\varphi(x) \otimes \psi(y)] - \mu_x \otimes \mu_y\|_{\mathcal{H}_h} = \text{HSIC}(x, y)\end{aligned}$$

using $\mathcal{H}_1 \otimes \mathcal{H}_2 \simeq HS(\mathcal{H}_2, \mathcal{H}_1)$.

- Characteristic functions: ϕ_{xy} , ϕ_x , ϕ_y .

Distance covariance

- Characteristic functions: ϕ_{xy} , ϕ_x , ϕ_y .
- Idea [Székely et al., 2007, Székely and Rizzo, 2009]:

$$x \perp y \Leftrightarrow \phi_{xy} = \phi_x \phi_y, \quad (x \in \mathbb{R}^{d_1}, y \in \mathbb{R}^{d_2}).$$

Distance covariance

- Characteristic functions: ϕ_{xy} , ϕ_x , ϕ_y .
- Idea [Székely et al., 2007, Székely and Rizzo, 2009]:

$$x \perp y \Leftrightarrow \phi_{xy} = \phi_x \phi_y, \quad (x \in \mathbb{R}^{d_1}, y \in \mathbb{R}^{d_2}).$$

- L_w^2 norm of ϕ_{xy} and $\phi_x \phi_y$:

$$dCov(x, y) = \|\phi_{xy} - \phi_x \phi_y\|_{L_w^2}$$

Distance covariance

- Characteristic functions: ϕ_{xy} , ϕ_x , ϕ_y .
- Idea [Székely et al., 2007, Székely and Rizzo, 2009]:

$$x \perp y \Leftrightarrow \phi_{xy} = \phi_x \phi_y, \quad (x \in \mathbb{R}^{d_1}, y \in \mathbb{R}^{d_2}).$$

- L_w^2 norm of ϕ_{xy} and $\phi_x \phi_y$, $\alpha \in (0, 2)$:

$$dCov(x, y) = \|\phi_{xy} - \phi_x \phi_y\|_{L_w^2}$$
$$w(a, b) = \frac{1}{c(d_1, \alpha)c(d_2, \alpha) \|a\|_2^{d_1+\alpha} \|b\|_2^{d_2+\alpha}},$$

Distance covariance

- Characteristic functions: ϕ_{xy} , ϕ_x , ϕ_y .
- Idea [Székely et al., 2007, Székely and Rizzo, 2009]:

$$x \perp y \Leftrightarrow \phi_{xy} = \phi_x \phi_y, \quad (x \in \mathbb{R}^{d_1}, y \in \mathbb{R}^{d_2}).$$

- L_w^2 norm of ϕ_{xy} and $\phi_x \phi_y$, $\alpha \in (0, 2)$:

$$\begin{aligned} dCov(x, y) &= \|\phi_{xy} - \phi_x \phi_y\|_{L_w^2} \\ w(a, b) &= \frac{1}{c(d_1, \alpha)c(d_2, \alpha)\|a\|_2^{d_1+\alpha}\|b\|_2^{d_2+\alpha}}, \\ c(d, \alpha) &= \frac{2\pi^{\frac{d}{2}}\Gamma(1 - \frac{\alpha}{2})}{\alpha 2^\alpha \Gamma(\frac{d+\alpha}{2})}. \end{aligned}$$

Distance covariance

- Characteristic functions: ϕ_{xy} , ϕ_x , ϕ_y .
- Idea [Székely et al., 2007, Székely and Rizzo, 2009]:

$$x \perp y \Leftrightarrow \phi_{xy} = \phi_x \phi_y, \quad (x \in \mathbb{R}^{d_1}, y \in \mathbb{R}^{d_2}).$$

- L_w^2 norm of ϕ_{xy} and $\phi_x \phi_y$, $\alpha \in (0, 2)$:

$$\begin{aligned} dCov(x, y) &= \|\phi_{xy} - \phi_x \phi_y\|_{L_w^2} \\ w(a, b) &= \frac{1}{c(d_1, \alpha)c(d_2, \alpha)\|a\|_2^{d_1+\alpha}\|b\|_2^{d_2+\alpha}}, \\ c(d, \alpha) &= \frac{2\pi^{\frac{d}{2}}\Gamma(1 - \frac{\alpha}{2})}{\alpha 2^\alpha \Gamma(\frac{d+\alpha}{2})}. \end{aligned}$$

- $x \perp y$ iff. $dCov(x, y) = 0$.

Distance covariance: $\alpha = 1$

Alternative form in terms of pairwise distances:

$$\begin{aligned} dCov^2(x, y) &= \mathbb{E}_{xy} \mathbb{E}_{x'y'} \|x - x'\|_2 \|y - y'\|_2 + \mathbb{E}_{xx'} \|x - x'\|_2 \mathbb{E}_{yy'} \|y - y'\|_2 \\ &\quad - 2\mathbb{E}_{xy} [\mathbb{E}_{x'} \|x - x'\|_2 \mathbb{E}_{y'} \|y - y'\|_2]. \end{aligned}$$

Distance covariance: $\alpha = 1$

Alternative form in terms of pairwise distances:

$$dCov^2(x, y) = \mathbb{E}_{xy} \mathbb{E}_{x'y'} \|x - x'\|_2 \|y - y'\|_2 + \mathbb{E}_{xx'} \|x - x'\|_2 \mathbb{E}_{yy'} \|y - y'\|_2 \\ - 2\mathbb{E}_{xy} [\mathbb{E}_{x'} \|x - x'\|_2 \mathbb{E}_{y'} \|y - y'\|_2].$$

Extension [Lyons, 2013]:

$$dCov^2(x, y) = \mathbb{E}_{xy} \mathbb{E}_{x'y'} \rho_1(x, x') \rho_2(y, y') + \mathbb{E}_{xx'}(x, x') \mathbb{E}_{yy'}(y, y') \\ - 2\mathbb{E}_{xy} [\mathbb{E}_{x'} \rho_1(x, x') \mathbb{E}_{y'} \rho_2(y, y')],$$

Distance covariance: $\alpha = 1$

Alternative form in terms of pairwise distances:

$$dCov^2(x, y) = \mathbb{E}_{xy} \mathbb{E}_{x'y'} \|x - x'\|_2 \|y - y'\|_2 + \mathbb{E}_{xx'} \|x - x'\|_2 \mathbb{E}_{yy'} \|y - y'\|_2 \\ - 2\mathbb{E}_{xy} [\mathbb{E}_{x'} \|x - x'\|_2 \mathbb{E}_{y'} \|y - y'\|_2].$$

Extension [Lyons, 2013]:

$$dCov^2(x, y) = \mathbb{E}_{xy} \mathbb{E}_{x'y'} \rho_1(x, x') \rho_2(y, y') + \mathbb{E}_{xx'}(x, x') \mathbb{E}_{yy'}(y, y') \\ - 2\mathbb{E}_{xy} [\mathbb{E}_{x'} \rho_1(x, x') \mathbb{E}_{y'} \rho_2(y, y')],$$

$(\mathcal{X}, \rho_1), (\mathcal{Y}, \rho_2)$: metric spaces of negative type (def & examples: in a moment).

Distance covariance vs. HSIC

$$\begin{aligned} dCov^2(x, y) = & \mathbb{E}_{xy} \mathbb{E}_{x'y'} \rho_1(x, x') \rho_2(y, y') + \mathbb{E}_{xx'} \rho_1(x, x') \mathbb{E}_{yy'} \rho_2(y, y') \\ & - 2\mathbb{E}_{xy} [\mathbb{E}_{x'} \rho_1(x, x') \mathbb{E}_{y'} \rho_2(y, y')] . \end{aligned}$$

Distance covariance vs. HSIC

$$\textcolor{red}{dCov}^2(x, y) = \mathbb{E}_{xy} \mathbb{E}_{x'y'} \rho_1(x, x') \rho_2(y, y') + \mathbb{E}_{xx'} \rho_1(x, x') \mathbb{E}_{yy'} \rho_2(y, y') \\ - 2\mathbb{E}_{xy} [\mathbb{E}_{x'} \rho_1(x, x') \mathbb{E}_{y'} \rho_2(y, y')].$$

Similarly to MMD (see later at $\widehat{\text{HSIC}}$):

$$\text{HSIC}^2(x, y) = \mathbb{E}_{xy} \mathbb{E}_{x'y'} k(x, x') \ell(y, y') + \mathbb{E}_{xx'} k(x, x') \mathbb{E}_{yy'} \ell(y, y') \\ - 2\mathbb{E}_{xy} [\mathbb{E}_{x'} k(x, x') \mathbb{E}_{y'} \ell(y, y')].$$

HSIC \Leftrightarrow distance covariance

+extension to semi-metric spaces of negative type:

Theorem ([Sejdinovic et al., 2013b])

$dCov^2(x, y; \rho_1, \rho_2) = 4\text{HSIC}^2(x, y; \mathcal{H}_k, \mathcal{H}_\ell)$, where

$$\begin{aligned}\rho_1(x, x') &= k(x, x) + k(x', x') - 2k(x, x'), \\ \rho_2(y, y') &= \ell(y, y) + \ell(y', y') - 2\ell(y, y').\end{aligned}$$

Definition

$\rho : \mathcal{X} \times \mathcal{X} \rightarrow [0, +\infty)$ is a **metric** on \mathcal{X} if

- $\rho(x, y) = 0 \Leftrightarrow x = y.$

Definition

$\rho : \mathcal{X} \times \mathcal{X} \rightarrow [0, +\infty)$ is a **metric** on \mathcal{X} if

- $\rho(x, y) = 0 \Leftrightarrow x = y.$
- Symmetry: $\rho(x, y) = \rho(y, x)$ for $\forall x, y \in \mathcal{X}.$

Definition

$\rho : \mathcal{X} \times \mathcal{X} \rightarrow [0, +\infty)$ is a **metric** on \mathcal{X} if

- $\rho(x, y) = 0 \Leftrightarrow x = y.$
- Symmetry: $\rho(x, y) = \rho(y, x)$ for $\forall x, y \in \mathcal{X}.$
- Triangle inequality: $\rho(x, z) \leq \rho(x, y) + \rho(y, z)$ for $\forall x, y, z \in \mathcal{X}.$

Definition

$\rho : \mathcal{X} \times \mathcal{X} \rightarrow [0, +\infty)$ is a **metric** on \mathcal{X} if

- $\rho(x, y) = 0 \Leftrightarrow x = y.$
- Symmetry: $\rho(x, y) = \rho(y, x)$ for $\forall x, y \in \mathcal{X}.$
- Triangle inequality: $\rho(x, z) \leq \rho(x, y) + \rho(y, z)$ for $\forall x, y, z \in \mathcal{X}.$

Examples:

- $\mathcal{X} = \mathbb{R}^d$, $\rho(x, y) = \|x - y\|_p = \left(\sum_{i=1}^d |x_i - y_i|^p \right)^{\frac{1}{p}}$, $p \geq 1$.

Definition

$\rho : \mathcal{X} \times \mathcal{X} \rightarrow [0, +\infty)$ is a **metric** on \mathcal{X} if

- $\rho(x, y) = 0 \Leftrightarrow x = y.$
- Symmetry: $\rho(x, y) = \rho(y, x)$ for $\forall x, y \in \mathcal{X}.$
- Triangle inequality: $\rho(x, z) \leq \rho(x, y) + \rho(y, z)$ for $\forall x, y, z \in \mathcal{X}.$

Examples:

- $\mathcal{X} = \mathbb{R}^d$, $\rho(x, y) = \|x - y\|_p = \left(\sum_{i=1}^d |x_i - y_i|^p \right)^{\frac{1}{p}}$, $p \geq 1$.
- $\mathcal{X} = C[a, b]$, $\rho(x, y) = \max_{z \in [a, b]} |x(z) - y(z)|.$

Definition

$\rho : \mathcal{X} \times \mathcal{X} \rightarrow [0, +\infty)$ is a **metric** on \mathcal{X} if

- $\rho(x, y) = 0 \Leftrightarrow x = y.$
- Symmetry: $\rho(x, y) = \rho(y, x)$ for $\forall x, y \in \mathcal{X}.$
- Triangle inequality: $\rho(x, z) \leq \rho(x, y) + \rho(y, z)$ for $\forall x, y, z \in \mathcal{X}.$

Examples:

- $\mathcal{X} = \mathbb{R}^d$, $\rho(x, y) = \|x - y\|_p = \left(\sum_{i=1}^d |x_i - y_i|^p \right)^{\frac{1}{p}}$, $p \geq 1$.
- $\mathcal{X} = C[a, b]$, $\rho(x, y) = \max_{z \in [a, b]} |x(z) - y(z)|.$
- \mathcal{X} any set. $\rho(x, y) = \delta_{x=y}.$

Semi-metric space: no triangle inequality

Definition

$\rho : \mathcal{X} \times \mathcal{X} \rightarrow [0, +\infty)$ is a **semi-metric** on \mathcal{X} if

- $\rho(x, y) = 0 \Leftrightarrow x = y$.
- symmetry: $\rho(x, y) = \rho(y, x)$, for $\forall x, y \in \mathcal{X}$.

Semi-metric space: no triangle inequality

Definition

$\rho : \mathcal{X} \times \mathcal{X} \rightarrow [0, +\infty)$ is a **semi-metric** on \mathcal{X} if

- $\rho(x, y) = 0 \Leftrightarrow x = y$.
- symmetry: $\rho(x, y) = \rho(y, x)$, for $\forall x, y \in \mathcal{X}$.

It is called **negative type** if in addition

$$\sum_{i=1}^n \sum_{j=1}^n a_i a_j \rho(x_i, x_j) \leq 0$$

for $\forall n \geq 2$, $\forall x_1, \dots, x_n \in \mathcal{X}$ and $\forall a_1, \dots, a_n \in \mathbb{R}$ with $\sum_{i=1}^n a_i = 0$.

Semi-metric space of negative type

[Berg et al., 1984]:

- $\rho : \checkmark \Rightarrow \rho^a : \checkmark$ for $\forall a \in (0, 1)$.

Semi-metric space of negative type

[Berg et al., 1984]:

- $\rho : \checkmark \Rightarrow \rho^a : \checkmark$ for $\forall a \in (0, 1)$.
- \Leftrightarrow description: $\exists m : \mathcal{X} \rightarrow \mathcal{H}$ (ilbert) injective mapping such that

$$\rho(x, y) = \|m(x) - m(y)\|_{\mathcal{H}}^2.$$

Semi-metric space of negative type

[Berg et al., 1984]:

- $\rho : \checkmark \Rightarrow \rho^a : \checkmark$ for $\forall a \in (0, 1)$.
- \Leftrightarrow description: $\exists m : \mathcal{X} \rightarrow \mathcal{H}$ (ilbert) injective mapping such that

$$\rho(x, y) = \|m(x) - m(y)\|_{\mathcal{H}}^2.$$

Thus,

- 2nd part $\Rightarrow \left(\mathbb{R}^d, \|\cdot\|_2^2\right) \checkmark$

Semi-metric space of negative type

[Berg et al., 1984]:

- $\rho : \checkmark \Rightarrow \rho^a : \checkmark$ for $\forall a \in (0, 1)$.
- \Leftrightarrow description: $\exists m : \mathcal{X} \rightarrow \mathcal{H}$ (ilbert) injective mapping such that

$$\rho(x, y) = \|m(x) - m(y)\|_{\mathcal{H}}^2.$$

Thus,

- 2nd part $\Rightarrow (\mathbb{R}^d, \|\cdot\|_2^2) \checkmark$
- +1st part $\Rightarrow \rho(x, y) = \|x - y\|_2^q \checkmark$ with $q \in (0, 2)$.

Semi-metric space of negative type

[Berg et al., 1984]:

- $\rho : \checkmark \Rightarrow \rho^a : \checkmark$ for $\forall a \in (0, 1)$.
- \Leftrightarrow description: $\exists m : \mathcal{X} \rightarrow \mathcal{H}$ (ilbert) injective mapping such that

$$\rho(x, y) = \|m(x) - m(y)\|_{\mathcal{H}}^2.$$

Thus,

- 2nd part $\Rightarrow (\mathbb{R}^d, \|\cdot\|_2^2) \checkmark$
- +1st part $\Rightarrow \rho(x, y) = \|x - y\|_2^q \checkmark$ with $q \in (0, 2)$.
- Specifically: $\rho(x, y) = \|x - y\|_2$ is OK.

Energy distance [Székely and Rizzo, 2004, Baringhaus and Franz, 2004, Székely and Rizzo, 2005]

$x, x' \sim \mathbb{P}, y, y' \sim \mathbb{Q}$:

$$EnDist(\mathbb{P}, \mathbb{Q}) = 2\mathbb{E}_{xy} \|\textcolor{blue}{x} - y\|_2 - \mathbb{E}_{xx'} \|x - x'\|_2 - \mathbb{E}_{yy'} \|y - y'\|_2,$$

Energy distance [Székely and Rizzo, 2004, Baringhaus and Franz, 2004, Székely and Rizzo, 2005]

$x, x' \sim \mathbb{P}, y, y' \sim \mathbb{Q}$:

$$EnDist(\mathbb{P}, \mathbb{Q}) = 2\mathbb{E}_{xy} \|\textcolor{blue}{x} - y\|_2 - \mathbb{E}_{xx'} \|x - x'\|_2 - \mathbb{E}_{yy'} \|y - y'\|_2,$$

$$EnDist(\mathbb{P}, \mathbb{Q}) = 2\mathbb{E}_{xy} \rho(\textcolor{blue}{x}, y) - \mathbb{E}_{xx'} \rho(x, x') - \mathbb{E}_{yy'} \rho(y, y').$$

Energy distance [Székely and Rizzo, 2004, Baringhaus and Franz, 2004, Székely and Rizzo, 2005]

$x, x' \sim \mathbb{P}, y, y' \sim \mathbb{Q}$:

$$EnDist(\mathbb{P}, \mathbb{Q}) = 2\mathbb{E}_{xy} \|x - y\|_2 - \mathbb{E}_{xx'} \|x - x'\|_2 - \mathbb{E}_{yy'} \|y - y'\|_2,$$

$$EnDist(\mathbb{P}, \mathbb{Q}) = 2\mathbb{E}_{xy} \rho(x, y) - \mathbb{E}_{xx'} \rho(x, x') - \mathbb{E}_{yy'} \rho(y, y').$$

Properties:

- $EnDist(\mathbb{P}, \mathbb{Q}) \geq 0$ with ρ metric of negative-type.

Energy distance [Székely and Rizzo, 2004, Baringhaus and Franz, 2004, Székely and Rizzo, 2005]

$x, x' \sim \mathbb{P}, y, y' \sim \mathbb{Q}$:

$$EnDist(\mathbb{P}, \mathbb{Q}) = 2\mathbb{E}_{xy} \|x - y\|_2 - \mathbb{E}_{xx'} \|x - x'\|_2 - \mathbb{E}_{yy'} \|y - y'\|_2,$$

$$EnDist(\mathbb{P}, \mathbb{Q}) = 2\mathbb{E}_{xy} \rho(x, y) - \mathbb{E}_{xx'} \rho(x, x') - \mathbb{E}_{yy'} \rho(y, y').$$

Properties:

- $EnDist(\mathbb{P}, \mathbb{Q}) \geq 0$ with ρ metric of negative-type.
- $EnDist(\mathbb{P}, \mathbb{Q}) = 0 \Leftrightarrow \mathbb{P} = \mathbb{Q}$ for (\mathcal{X}, ρ) strictly negative spaces; example: $(\mathbb{R}^d, \|\cdot\|_2)$.

Strict negativity

In addition:

$$\sum_{i=1}^n \sum_{j=1}^n a_i a_j \rho(x_i, x_j) < 0$$

if x_i -s are distinct and $\exists a_i \neq 0$.

Energy distance vs. MMD

Energy distance: also called N-distance
[Zinger et al., 1992, Klebanov, 2005],

$$EnDist(\mathbb{P}, \mathbb{Q}) = 2\mathbb{E}_{xy}\rho(x, y) - \mathbb{E}_{xx'}\rho(x, x') - \mathbb{E}_{yy'}\rho(y, y').$$

Energy distance vs. MMD

Energy distance: also called N-distance
[Zinger et al., 1992, Klebanov, 2005],

$$EnDist(\mathbb{P}, \mathbb{Q}) = 2\mathbb{E}_{xy}\rho(x, y) - \mathbb{E}_{xx'}\rho(x, x') - \mathbb{E}_{yy'}\rho(y, y').$$

MMD (recall):

$$\text{MMD}^2(\mathbb{P}, \mathbb{Q}) = \mathbb{E}_{x,x'}k(x, x') + \mathbb{E}_{y,y'}k(y, y') - 2\mathbb{E}_{xy}k(x, y).$$

MMD \Leftrightarrow energy distance

Theorem ([Sejdinovic et al., 2013b])

$$EnDist(\mathbb{P}, \mathbb{Q}; \rho) = 2\text{MMD}^2(\mathbb{P}, \mathbb{Q}; \mathcal{H}_k),$$

where

$$\rho(x, y) = k(x, x) + k(y, y) - 2k(x, y).$$

Central in applications: characteristic property

- HSIC, $\mathbf{k} = \otimes_{m=1}^M k_m$, $x = (x_m)_{m=1}^M$:

$$\text{HSIC}_k(\mathbb{P}) := \text{MMD}_{\mathbf{k}}\left(\mathbb{P}, \otimes_{m=1}^M \mathbb{P}_m\right), \quad \mathbf{k}(x, x') := \prod_{m=1}^M k_m(x_m, x'_m).$$

Central in applications: characteristic property

- HSIC, $\mathbf{k} = \otimes_{m=1}^M k_m$, $x = (x_m)_{m=1}^M$:

$$\text{HSIC}_k(\mathbb{P}) := \text{MMD}_{\mathbf{k}} \left(\mathbb{P}, \otimes_{m=1}^M \mathbb{P}_m \right), \quad \mathbf{k}(x, x') := \prod_{m=1}^M k_m(x_m, x'_m).$$

$k = \otimes_{m=1}^M k_m$ will be called **I-characteristic** if

$$\text{HSIC}_k(\mathbb{P}) = 0 \Leftrightarrow \mathbb{P} = \otimes_{m=1}^M \mathbb{P}_m.$$

Central in applications: characteristic property

- HSIC, $k = \otimes_{m=1}^M k_m$, $x = (x_m)_{m=1}^M$:

$$\text{HSIC}_k(\mathbb{P}) := \text{MMD}_k\left(\mathbb{P}, \otimes_{m=1}^M \mathbb{P}_m\right), \quad k(x, x') := \prod_{m=1}^M k_m(x_m, x'_m).$$

$k = \otimes_{m=1}^M k_m$ will be called **I-characteristic** if

$$\text{HSIC}_k(\mathbb{P}) = 0 \Leftrightarrow \mathbb{P} = \otimes_{m=1}^M \mathbb{P}_m.$$

Recall (MMD): k is called **characteristic** if

$$\text{MMD}_k(\mathbb{P}, \mathbb{Q}) = 0 \Leftrightarrow \mathbb{P} = \mathbb{Q}.$$

Central in applications: characteristic property

- HSIC, $\mathbf{k} = \otimes_{m=1}^M k_m$, $x = (x_m)_{m=1}^M$:

$$\text{HSIC}_k(\mathbb{P}) := \text{MMD}_{\mathbf{k}}\left(\mathbb{P}, \otimes_{m=1}^M \mathbb{P}_m\right), \quad \mathbf{k}(x, x') := \prod_{m=1}^M k_m(x_m, x'_m).$$

$k = \otimes_{m=1}^M k_m$ will be called **\mathcal{I} -characteristic** if

$$\text{HSIC}_k(\mathbb{P}) = 0 \Leftrightarrow \mathbb{P} = \otimes_{m=1}^M \mathbb{P}_m.$$

Recall (MMD): k is called **characteristic** if

$$\text{MMD}_k(\mathbb{P}, \mathbb{Q}) = 0 \Leftrightarrow \mathbb{P} = \mathbb{Q}.$$

$\otimes_{m=1}^M k_m$: universal \Rightarrow characteristic \Rightarrow \mathcal{I} -characteristic.

Relation? Conditions in terms of k_m -s?

$\otimes_{m=1}^M k_m :$

$\mathcal{I}\text{-char}$ \longleftrightarrow char \longleftrightarrow universal



$(k_m)_{m=1}^M :$

char $\xrightarrow{\text{[Sriperumbudur et al., 2011]}}$ -universal
 $\xleftarrow{\text{[Sriperumbudur et al., 2011]}}$

Existing Results, $M = 2$

- [Blanchard et al., 2011, Waegeman et al., 2012, Gretton, 2015]:
 $k_1 \& k_2$: universal $\Rightarrow k_1 \otimes k_2$: universal ($\Rightarrow \mathcal{I}$ -characteristic).

Existing Results, $M = 2$

- [Blanchard et al., 2011, Waegeman et al., 2012, Gretton, 2015]:
 $k_1 \& k_2$: universal $\Rightarrow k_1 \otimes k_2$: universal ($\Rightarrow \mathcal{I}$ -characteristic).
- Distance covariance [Lyons, 2013, Sejdinovic et al., 2013b]:
 $k_1 \& k_2$: characteristic $\Leftrightarrow k_1 \otimes k_2$: \mathcal{I} -characteristic.

Existing Results, $M = 2$

- [Blanchard et al., 2011, Waegeman et al., 2012, Gretton, 2015]:
 $k_1 \& k_2$: universal $\Rightarrow k_1 \otimes k_2$: universal ($\Rightarrow \mathcal{I}$ -characteristic).
- Distance covariance [Lyons, 2013, Sejdinovic et al., 2013b]:
 $k_1 \& k_2$: characteristic $\Leftrightarrow k_1 \otimes k_2$: \mathcal{I} -characteristic.

Goal

Extension to $M \geq 2$.

Existing Results, $M = 2$

- [Blanchard et al., 2011, Waegeman et al., 2012, Gretton, 2015]:
 $k_1 \& k_2$: universal $\Rightarrow k_1 \otimes k_2$: universal ($\Rightarrow \mathcal{I}$ -characteristic).
- Distance covariance [Lyons, 2013, Sejdinovic et al., 2013b]:
 $k_1 \& k_2$: characteristic $\Leftrightarrow k_1 \otimes k_2$: \mathcal{I} -characteristic.

Goal

Extension to $M \geq 2$.

Main Challenge

' $\otimes k_m$: \mathcal{I} -characteristic $\Leftrightarrow k_m$: characteristic ($\forall m$)' does NOT hold.

Proposition (characteristic property)

- $\otimes_{m=1}^M k_m$: characteristic $\Rightarrow (k_m)_{m=1}^M$ are characteristic.
- $\nLeftarrow [|\mathcal{X}_m| = 2, k_m(x, x') = 2\delta_{x,x'} - 1]$

Proposition (characteristic property)

- $\otimes_{m=1}^M k_m$: characteristic $\Rightarrow (k_m)_{m=1}^M$ are characteristic.
- $\nLeftarrow [|\mathcal{X}_m| = 2, k_m(x, x') = 2\delta_{x,x'} - 1]$

Proposition (\mathcal{I} -characteristic property)

- k_1, k_2 : characteristic $\Rightarrow k_1 \otimes k_2$: \mathcal{I} -characteristic.
- \Leftarrow : for $\forall M \geq 2$.

Results [Szabó and Sriperumbudur, 2017]

Proposition (characteristic property)

- $\otimes_{m=1}^M k_m$: characteristic $\Rightarrow (k_m)_{m=1}^M$ are characteristic.
- $\nLeftarrow [|\mathcal{X}_m| = 2, k_m(x, x') = 2\delta_{x,x'} - 1]$

Proposition (\mathcal{I} -characteristic property)

- k_1, k_2 : characteristic $\Rightarrow k_1 \otimes k_2$: \mathcal{I} -characteristic.
- \Leftarrow : for $\forall M \geq 2$.
- k_1, k_2, k_3 : characteristic $\nRightarrow \otimes_{m=1}^3 k_m$: \mathcal{I} -characteristic [Ex].

Results [Szabó and Sriperumbudur, 2017]

Proposition (characteristic property)

- $\otimes_{m=1}^M k_m$: characteristic $\Rightarrow (k_m)_{m=1}^M$ are characteristic.
- $\nLeftarrow [|\mathcal{X}_m| = 2, k_m(x, x') = 2\delta_{x,x'} - 1]$

Proposition (\mathcal{I} -characteristic property)

- k_1, k_2 : characteristic $\Rightarrow k_1 \otimes k_2$: \mathcal{I} -characteristic.
- \Leftarrow : for $\forall M \geq 2$.
- k_1, k_2, k_3 : characteristic $\nRightarrow \otimes_{m=1}^3 k_m$: \mathcal{I} -characteristic [Ex].
- k_1, k_2 : universal, k_3 : characteristic $\nRightarrow \otimes_{m=1}^3 k_m$: \mathcal{I} -char [Ex].

Results: continued

Proposition ($\mathcal{X}_m = \mathbb{R}^{d_m}$, k_m : continuous, shift-invariant, bounded)

$(k_m)_{m=1}^M$ -s are characteristic $\Leftrightarrow \otimes_{m=1}^M k_m$: \mathcal{I} -characteristic \Leftrightarrow
 $\otimes_{m=1}^M k_m$: characteristic.

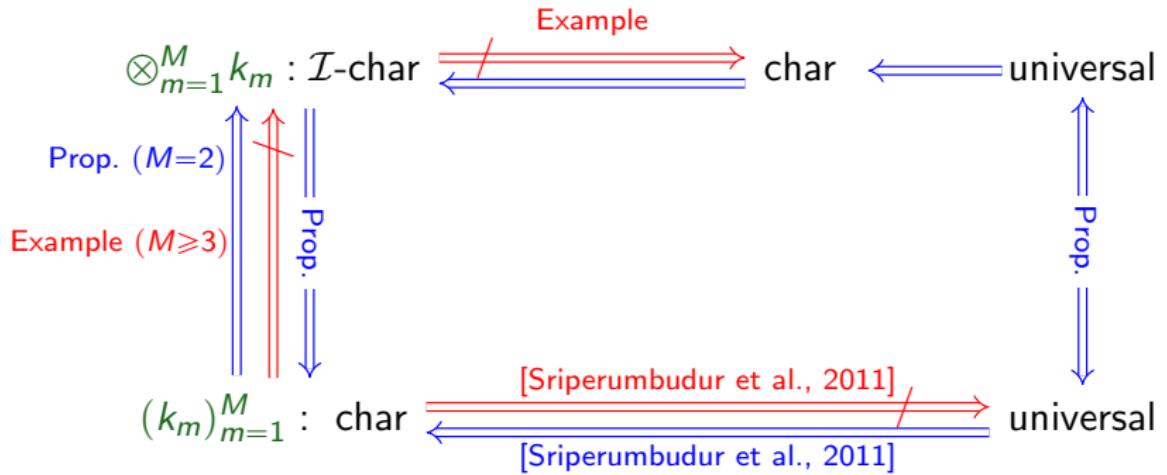
Results: continued

Proposition ($\mathcal{X}_m = \mathbb{R}^{d_m}$, k_m : continuous, shift-invariant, bounded)

$(k_m)_{m=1}^M$ -s are characteristic $\Leftrightarrow \otimes_{m=1}^M k_m$: \mathcal{I} -characteristic \Leftrightarrow
 $\otimes_{m=1}^M k_m$: characteristic.

Proposition (Universality)

$\otimes_{m=1}^M k_m$: universal $\Leftrightarrow (k_m)_{m=1}^M$ are universal.



Covariance operator: finished.

Recall

- KCCA: independence measure,

$$\rho_{\text{KCCA}}(x, y) = \sup_{f \in \mathcal{H}_k, g \in \mathcal{H}_\ell} \text{corr}(f(x), g(y)).$$

Recall

- KCCA: independence measure,

$$\rho_{\text{KCCA}}(x, y) = \sup_{f \in \mathcal{H}_k, g \in \mathcal{H}_\ell} \text{corr}(f(x), g(y)).$$

- Mean embedding: distribution representation,

$$\mu_{\mathbb{P}} = \int_{\mathcal{X}} k(\cdot, x) d\mathbb{P}(x).$$

Recall

- KCCA: independence measure,

$$\rho_{\text{KCCA}}(x, y) = \sup_{f \in \mathcal{H}_k, g \in \mathcal{H}_\ell} \text{corr}(f(x), g(y)).$$

- Mean embedding: distribution representation,

$$\mu_{\mathbb{P}} = \int_{\mathcal{X}} k(\cdot, x) d\mathbb{P}(x).$$

- MMD: (semi)-metric defined by mean embedding,

$$\text{MMD}(\mathbb{P}, \mathbb{Q}) = \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k}.$$

Recall

- KCCA: independence measure,

$$\rho_{\text{KCCA}}(x, y) = \sup_{f \in \mathcal{H}_k, g \in \mathcal{H}_\ell} \text{corr}(f(x), g(y)).$$

- Mean embedding: distribution representation,

$$\mu_{\mathbb{P}} = \int_{\mathcal{X}} k(\cdot, x) d\mathbb{P}(x).$$

- MMD: (semi)-metric defined by mean embedding,

$$\text{MMD}(\mathbb{P}, \mathbb{Q}) = \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k}.$$

- Cross-covariance operator:

$$C_{xy}^c = \mathbb{E}_{xy} [\varphi(x) \otimes \psi(y)] - \mu_x \otimes \mu_y.$$

Recall

- KCCA: independence measure,

$$\rho_{\text{KCCA}}(x, y) = \sup_{f \in \mathcal{H}_k, g \in \mathcal{H}_\ell} \text{corr}(f(x), g(y)).$$

- Mean embedding: distribution representation,

$$\mu_{\mathbb{P}} = \int_{\mathcal{X}} k(\cdot, x) d\mathbb{P}(x).$$

- MMD: (semi)-metric defined by mean embedding,

$$\text{MMD}(\mathbb{P}, \mathbb{Q}) = \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k}.$$

- Cross-covariance operator:

$$C_{xy}^c = \mathbb{E}_{xy} [\varphi(x) \otimes \psi(y)] - \mu_x \otimes \mu_y.$$

- HSIC: independence measure,

$$\text{HSIC}(x, y) = \|C_{xy}^c\|_{HS}.$$

No density estimation

Thus,

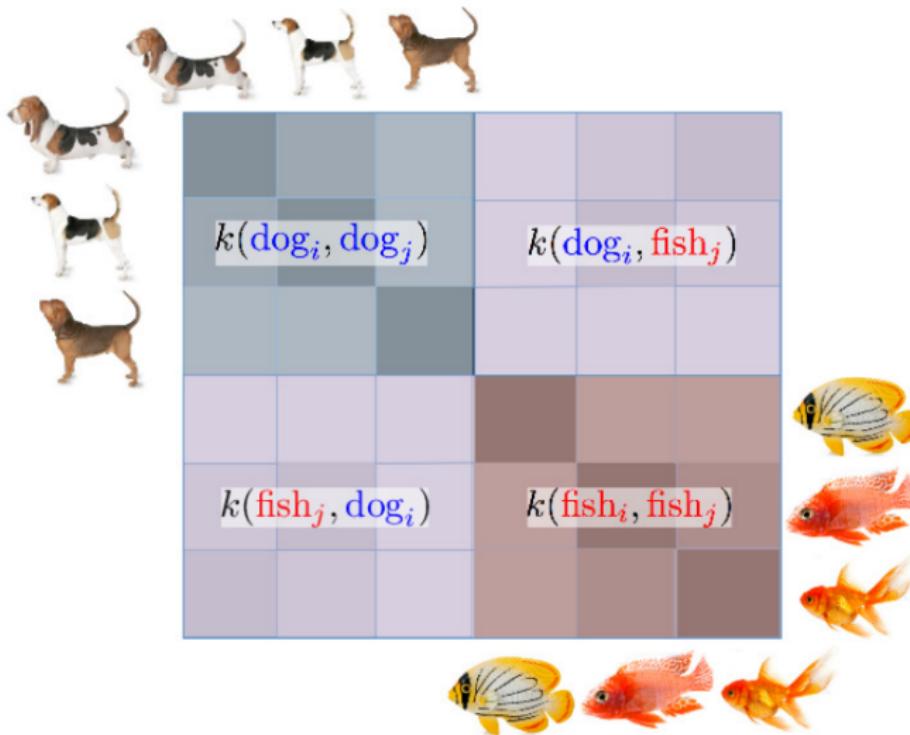
- independence measure,
- distance,
- inner product

measures/estimates on probability distributions

without density estimation!

HSIC estimators

Recall: MMD estimator



$$\widehat{\text{MMD}}_u^2(\mathbb{P}, \mathbb{Q}) = \overline{G_{\mathbb{P}, \mathbb{P}}} + \overline{G_{\mathbb{Q}, \mathbb{Q}}} - 2\overline{G_{\mathbb{P}, \mathbb{Q}}} \quad (\text{without diagonals in } \overline{G_{\mathbb{P}, \mathbb{P}}}, \overline{G_{\mathbb{Q}, \mathbb{Q}}})$$

HSIC: intuition. \mathcal{X} : images, \mathcal{Y} : descriptions.



Their noses guide them through life, and they're never happier than when following an interesting scent. They need plenty of exercise, about an hour a day if possible.



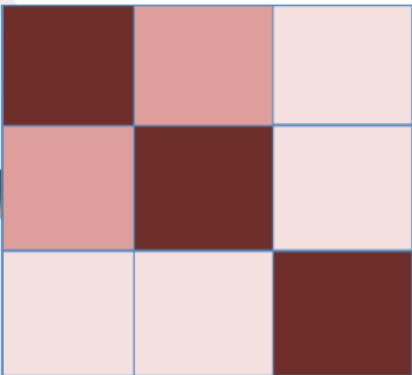
A large animal who slings slobber, exudes a distinctive houndy odor, and wants nothing more than to follow his nose. They need a significant amount of exercise and mental stimulation.



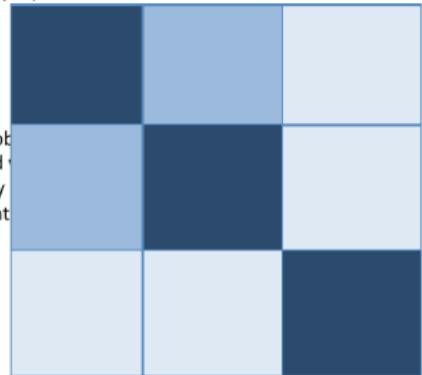
Known for their curiosity, intelligence, and excellent communication skills, the Javanese breed is perfect if you want a responsive, interactive pet, one that will blow in your ear and follow you everywhere.

Text from dogtime.com and petfinder.com

HSIC intuition: Gram matrices

 $\tilde{\mathbf{G}}_x$ 

Their noses guide them through life, and they're never happier than when following an interesting scent. They need plenty of exercise, about an hour a day if possible.

 $\tilde{\mathbf{G}}_y$ 

A large animal who slings slob distinctive houndy odor, and than to follow his nose. They amount of exercise and ment

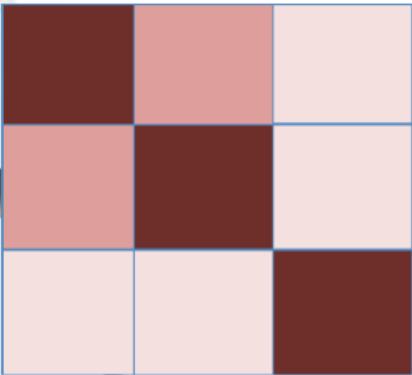


Known for their curiosity, intelligence, and excellent communication skills, the Javanese breed is perfect if you want a responsive, interactive pet, one that will blow in your ear and follow you everywhere.

HSIC intuition: Gram matrices



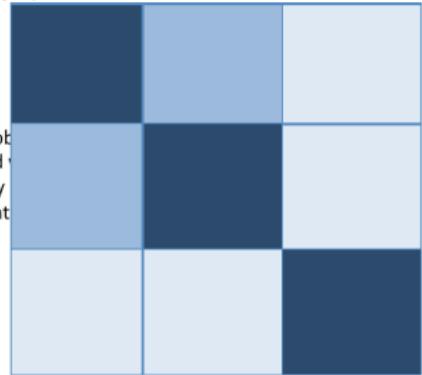
$\tilde{\mathbf{G}}_x$



Their noses guide them through life, and they're never happier than when following an interesting scent. They need plenty of exercise, about an hour a day if possible.



$\tilde{\mathbf{G}}_y$



A large animal who slings slob distinctive houndy odor, and than to follow his nose. They amount of exercise and ment



Known for their curiosity, intelligence, and excellent communication skills, the Javanese breed is perfect if you want a responsive, interactive pet, one that will blow in your ear and follow you everywhere.

Empirical estimate:

$$\widehat{\text{HSIC}^2} = \frac{1}{n^2} \left\langle \tilde{\mathbf{G}}_x, \tilde{\mathbf{G}}_y \right\rangle_F.$$

Cocktail party: HSIC demo



$$\mathbf{x} = \mathbf{A}\mathbf{s}, \quad \mathbf{s} = \left[\mathbf{s}^1; \dots; \mathbf{s}^M \right],$$

where \mathbf{s}^m -s are non-Gaussian & independent.

- Goal: $\{\mathbf{x}_t\}_{t=1}^T \rightarrow \mathbf{W} = \mathbf{A}^{-1}, \{\mathbf{s}_t\}_{t=1}^T$,

$$\mathbf{x} = \mathbf{A}\mathbf{s}, \quad \mathbf{s} = [\mathbf{s}^1; \dots; \mathbf{s}^M],$$

where \mathbf{s}^m -s are non-Gaussian & independent.

- Goal: $\{\mathbf{x}_t\}_{t=1}^T \rightarrow \mathbf{W} = \mathbf{A}^{-1}, \{\mathbf{s}_t\}_{t=1}^T$,
- Objective function:

$$\hat{\mathbf{s}} = \mathbf{W}\mathbf{x},$$

$$J(\mathbf{W}) = I\left(\hat{\mathbf{s}}^1, \dots, \hat{\mathbf{s}}^M\right) \rightarrow \min_{\mathbf{W}}.$$

- Hidden sources (s):

A B C D E F

ISA: source, observation

- Hidden sources (s):

A B C D E F



- Observation (x):



- Estimated sources (\hat{s}):



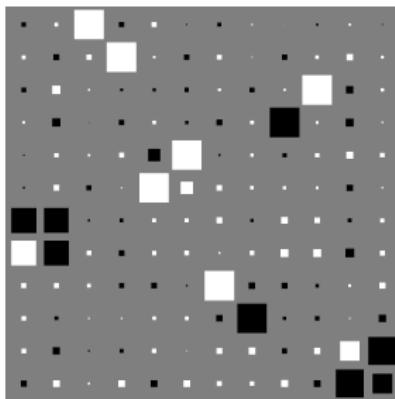
The image displays the word "BROADWAY" in a bold, sans-serif font. The letters are constructed from numerous small, dark gray or black dots, giving them a granular, point-based appearance. The letters are slightly overlapping, with some dots appearing in multiple locations to represent the estimated sources.

ISA: estimated sources using HSIC, ambiguity

- Estimated sources (\hat{s}):



- Performance ($\hat{W}A$), ambiguity:

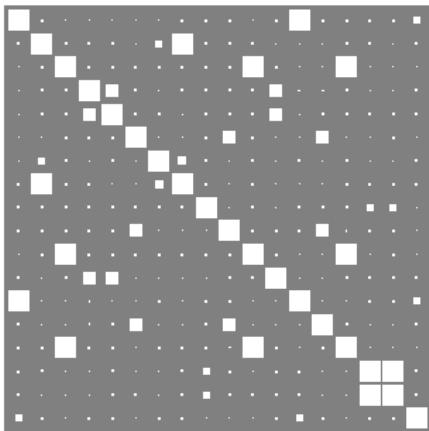


Conjecture: ISA separation theorem [Cardoso, 1998]

- $\text{ISA} = \text{ICA} + \text{permutation.}$

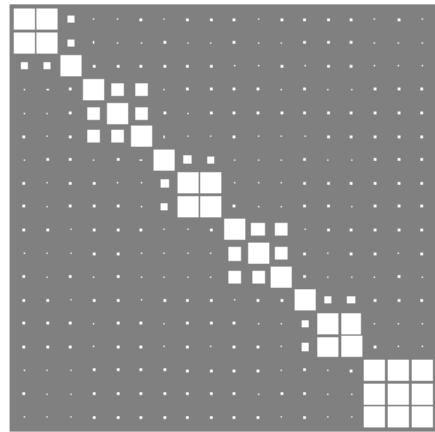
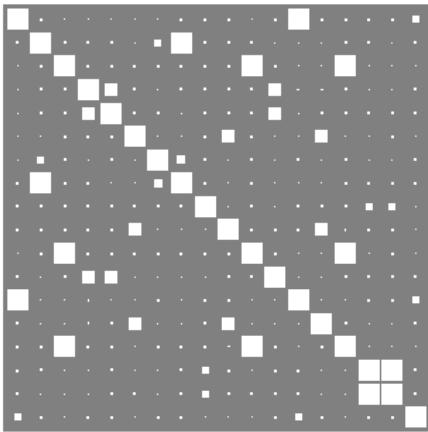
Conjecture: ISA separation theorem [Cardoso, 1998]

- ISA = ICA + permutation. $\widehat{\text{HSIC}}(\hat{s}_i, \hat{s}_j)$,



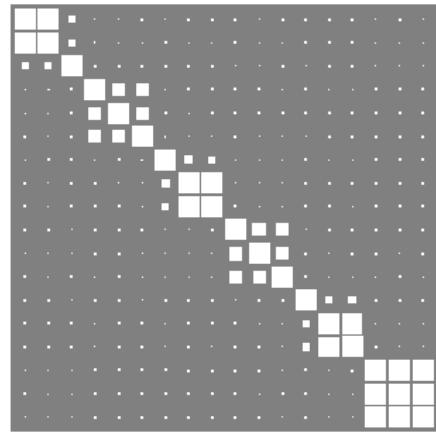
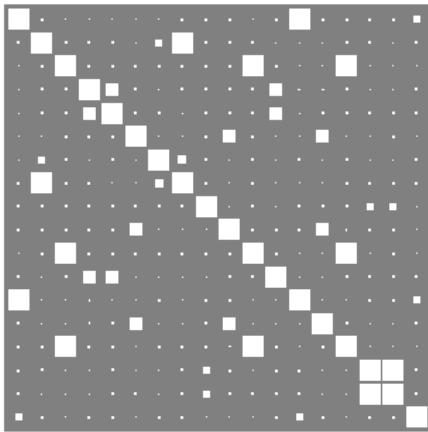
Conjecture: ISA separation theorem [Cardoso, 1998]

- ISA = ICA + permutation. $\widehat{\text{HSIC}}(\hat{s}_i, \hat{s}_j)$,



Conjecture: ISA separation theorem [Cardoso, 1998]

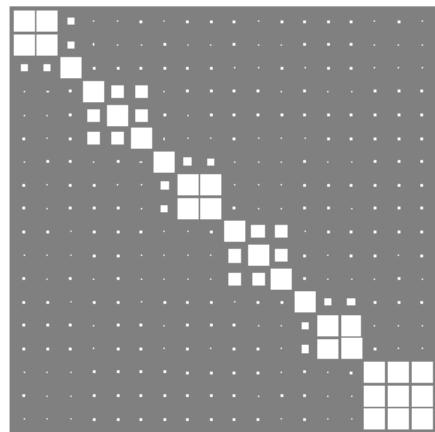
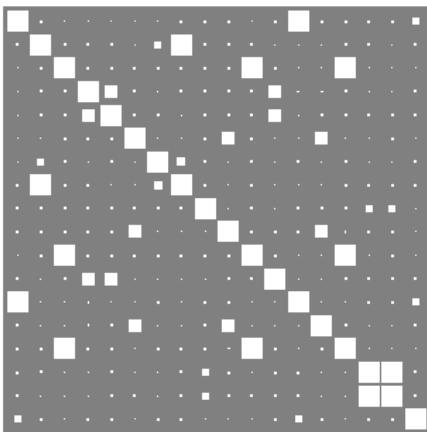
- ISA = ICA + permutation. $\widehat{\text{HSIC}}(\hat{s}_i, \hat{s}_j)$,



- Basis of the state-of-the-art ISA solvers.

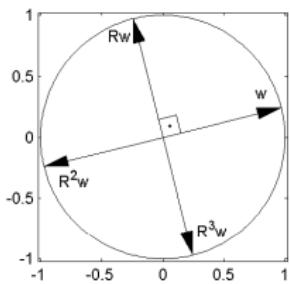
Conjecture: ISA separation theorem [Cardoso, 1998]

- ISA = ICA + permutation. $\widehat{\text{HSIC}}(\hat{s}_i, \hat{s}_j)$,



- Basis of the state-of-the-art ISA solvers.
- Sufficient conditions [Szabó et al., 2012]:
 - s^m : spherical [Fang et al., 1990].

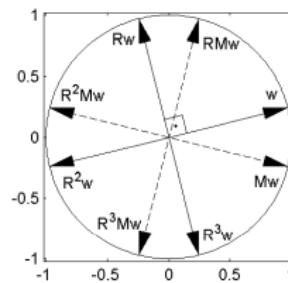
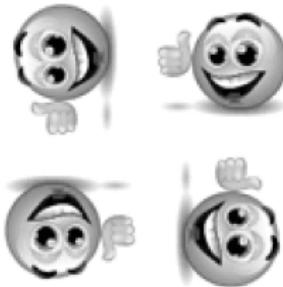
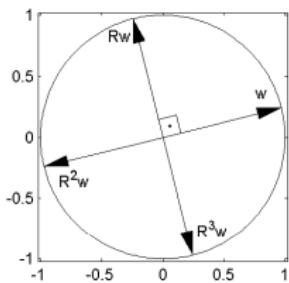
ISA separation theorem \rightarrow for $\dim(\mathbf{s}^m) = 2$ less is enough.



Invariance to

- 90° rotation: $f(u_1, u_2) = f(-u_2, u_1) = f(-u_1, -u_2) = f(u_2, -u_1)$.

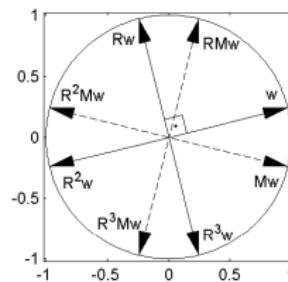
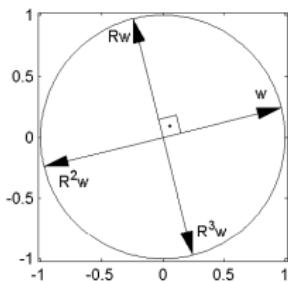
ISA separation theorem \rightarrow for $\dim(\mathbf{s}^m) = 2$ less is enough.



Invariance to

- 90° rotation: $f(u_1, u_2) = f(-u_2, u_1) = f(-u_1, -u_2) = f(u_2, -u_1)$.
- permutation and sign: $f(\pm u_1, \pm u_2) = f(\pm u_2, \pm u_1)$.

ISA separation theorem \rightarrow for $\dim(\mathbf{s}^m) = 2$ less is enough.



Invariance to

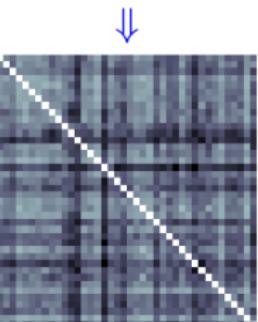
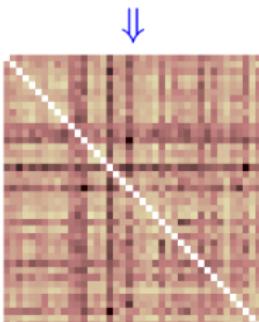
- 90° rotation: $f(u_1, u_2) = f(-u_2, u_1) = f(-u_1, -u_2) = f(u_2, -u_1)$.
- permutation and sign: $f(\pm u_1, \pm u_2) = f(\pm u_2, \pm u_1)$.
- L^p -spherical: $f(u_1, u_2) = h(\sum_i |u_i|^p)$ ($p > 0$).

Another HSIC demo: translation

- 5-line extracts.
- representation, kernel: bag-of-words, r -spectrum ($r = 5$).
- sample size: $n = 10$. repetitions: 300.

... no doubt there is great pressure on provincial and municipal governments in relation to the issue of child care, but the reality is that there have been no cuts to child care funding from the federal government to the provinces. In fact, we have increased federal investments for early childhood development...

... il est évident que les ordres de gouvernements provinciaux et municipaux subissent de fortes pressions en ce qui concerne les services de garde, mais le gouvernement n'a pas réduit le financement qu'il verse aux provinces pour les services de garde. Au contraire, nous avons augmenté le financement fédéral pour le développement des jeunes enfants...



Another HSIC demo: translation

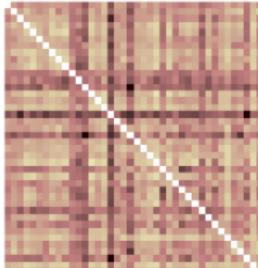
- 5-line extracts.
- representation, kernel: bag-of-words, r -spectrum ($r = 5$).
- sample size: $n = 10$. repetitions: 300.

Results:

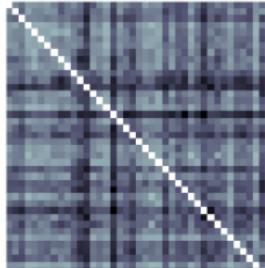
- r -spectrum: average Type-II error = 0 ($\alpha = 0.05$),
- bag-of-words: 0.18.

... no doubt there is great pressure on provincial and municipal governments in relation to the issue of child care, but the reality is that there have been no cuts to child care funding from the federal government to the provinces. In fact, we have increased federal investments for early childhood development...

... il est évident que les ordres de gouvernements provinciaux et municipaux subissent de fortes pressions en ce qui concerne les services de garde, mais le gouvernement n'a pas réduit le financement qu'il verse aux provinces pour les services de garde. Au contraire, nous avons augmenté le financement fédéral pour le développement des jeunes enfants...



⇒HSIC⇐



Recall: MMD in terms of kernel evaluations

$$\begin{aligned}\text{MMD}^2(\mathbb{P}, \mathbb{Q}) &= \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k}^2 = \\ &= \mathbb{E}_{x \sim \mathbb{P}, x' \sim \mathbb{P}} k(x, x') + \mathbb{E}_{y \sim \mathbb{Q}, y' \sim \mathbb{Q}} k(y, y') \\ &\quad - 2\mathbb{E}_{x \sim \mathbb{P}, y \sim \mathbb{Q}} k(x, y).\end{aligned}$$

Recall: MMD in terms of kernel evaluations

$$\begin{aligned}\text{MMD}^2(\mathbb{P}, \mathbb{Q}) &= \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k}^2 = \\ &= \mathbb{E}_{x \sim \mathbb{P}, x' \sim \mathbb{P}} k(x, x') + \mathbb{E}_{y \sim \mathbb{Q}, y' \sim \mathbb{Q}} k(y, y') \\ &\quad - 2\mathbb{E}_{x \sim \mathbb{P}, y \sim \mathbb{Q}} k(x, y).\end{aligned}$$

Question

Can we rewrite HSIC in terms of expected kernel values?

HSIC in terms of kernel evaluations [Gretton et al., 2005a]

$$\text{HSIC}^2(x, y) = \|C_{xy}^c\|_{HS}^2 = \|C_{xy}^u - \mu_x \otimes \mu_y\|_{HS}^2$$

HSIC in terms of kernel evaluations [Gretton et al., 2005a]

$$\begin{aligned}\text{HSIC}^2(x, y) &= \|C_{xy}^c\|_{HS}^2 = \|C_{xy}^u - \mu_x \otimes \mu_y\|_{HS}^2 \\ &= \|C_{xy}^u\|_{HS}^2 + \|\mu_x \otimes \mu_y\|_{HS}^2 - 2 \langle C_{xy}^u, \mu_x \otimes \mu_y \rangle_{HS}.\end{aligned}$$

HSIC in terms of kernel evaluations [Gretton et al., 2005a]

$$\begin{aligned}\text{HSIC}^2(x, y) &= \|C_{xy}^c\|_{HS}^2 = \|C_{xy}^u - \mu_x \otimes \mu_y\|_{HS}^2 \\ &= \|C_{xy}^u\|_{HS}^2 + \|\mu_x \otimes \mu_y\|_{HS}^2 - 2 \langle C_{xy}^u, \mu_x \otimes \mu_y \rangle_{HS}.\end{aligned}$$

First term:

$$\|C_{xy}^u\|_{HS}^2 = \langle \mathbb{E}_{xy} [\varphi(x) \otimes \psi(y)], \mathbb{E}_{x'y'} [\varphi(x') \otimes \psi(y')] \rangle_{HS}$$

HSIC in terms of kernel evaluations [Gretton et al., 2005a]

$$\begin{aligned}\text{HSIC}^2(x, y) &= \|C_{xy}^c\|_{HS}^2 = \|C_{xy}^u - \mu_x \otimes \mu_y\|_{HS}^2 \\ &= \|C_{xy}^u\|_{HS}^2 + \|\mu_x \otimes \mu_y\|_{HS}^2 - 2 \langle C_{xy}^u, \mu_x \otimes \mu_y \rangle_{HS}.\end{aligned}$$

First term:

$$\begin{aligned}\|C_{xy}^u\|_{HS}^2 &= \langle \mathbb{E}_{xy} [\varphi(x) \otimes \psi(y)], \mathbb{E}_{x'y'} [\varphi(x') \otimes \psi(y')] \rangle_{HS} \\ &= \mathbb{E}_{xy} \mathbb{E}_{x'y'} \underbrace{\langle \varphi(x) \otimes \psi(y), \varphi(x') \otimes \psi(y') \rangle_{HS}}_{\langle \varphi(x), \varphi(x') \rangle_{\mathcal{H}_k} \langle \psi(y), \psi(y') \rangle_{\mathcal{H}_\ell}}\end{aligned}$$

$$\langle e_1 \otimes f_1, e_2 \otimes f_2 \rangle_{HS(\mathcal{H}_2, \mathcal{H}_1)} = \langle e_1, e_2 \rangle_{\mathcal{H}_1} \langle f_1, f_2 \rangle_{\mathcal{H}_2}.$$

HSIC in terms of kernel evaluations [Gretton et al., 2005a]

$$\begin{aligned}\text{HSIC}^2(x, y) &= \|C_{xy}^c\|_{HS}^2 = \|C_{xy}^u - \mu_x \otimes \mu_y\|_{HS}^2 \\ &= \|C_{xy}^u\|_{HS}^2 + \|\mu_x \otimes \mu_y\|_{HS}^2 - 2 \langle C_{xy}^u, \mu_x \otimes \mu_y \rangle_{HS}.\end{aligned}$$

First term:

$$\begin{aligned}\|C_{xy}^u\|_{HS}^2 &= \langle \mathbb{E}_{xy} [\varphi(x) \otimes \psi(y)], \mathbb{E}_{x'y'} [\varphi(x') \otimes \psi(y')] \rangle_{HS} \\ &= \mathbb{E}_{xy} \mathbb{E}_{x'y'} \underbrace{\langle \varphi(x) \otimes \psi(y), \varphi(x') \otimes \psi(y') \rangle_{HS}}_{\langle \varphi(x), \varphi(x') \rangle_{\mathcal{H}_k} \langle \psi(y), \psi(y') \rangle_{\mathcal{H}_\ell}} \\ &= \mathbb{E}_{xy} \mathbb{E}_{x'y'} k(x, x') \ell(y, y').\end{aligned}$$

$$\langle e_1 \otimes f_1, e_2 \otimes f_2 \rangle_{HS(\mathcal{H}_2, \mathcal{H}_1)} = \langle e_1, e_2 \rangle_{\mathcal{H}_1} \langle f_1, f_2 \rangle_{\mathcal{H}_2}.$$

HSIC: second term

$$\|\mu_x \otimes \mu_y\|_{HS}^2 = \langle \mu_x \otimes \mu_y, \mu_x \otimes \mu_y \rangle_{HS}$$

HSIC: second term

$$\begin{aligned}\|\mu_x \otimes \mu_y\|_{HS}^2 &= \langle \mu_x \otimes \mu_y, \mu_x \otimes \mu_y \rangle_{HS} \\ &= \langle \mu_x, \mu_x \rangle_{\mathcal{H}_k} \langle \mu_y, \mu_y \rangle_{\mathcal{H}_\ell}\end{aligned}$$

HSIC: second term

$$\begin{aligned}\|\mu_x \otimes \mu_y\|_{HS}^2 &= \langle \mu_x \otimes \mu_y, \mu_x \otimes \mu_y \rangle_{HS} \\ &= \langle \mu_x, \mu_x \rangle_{\mathcal{H}_k} \langle \mu_y, \mu_y \rangle_{\mathcal{H}_\ell} \\ &= \mathbb{E}_{xx'} k(x, x') \mathbb{E}_{yy'} \ell(y, y').\end{aligned}$$

$$\langle \mathcal{C}_{xy}^u, \mu_x \otimes \mu_y \rangle_{HS} = \langle \mathbb{E}_{xy} [\varphi(x) \otimes \psi(y)], \mu_x \otimes \mu_y \rangle_{HS}$$

$$\begin{aligned}\langle C_{xy}^u, \mu_x \otimes \mu_y \rangle_{HS} &= \langle \mathbb{E}_{xy} [\varphi(x) \otimes \psi(y)], \mu_x \otimes \mu_y \rangle_{HS} \\ &= \mathbb{E}_{xy} \underbrace{\langle \varphi(x) \otimes \psi(y), \mu_x \otimes \mu_y \rangle_{HS}}_{\underbrace{\langle \varphi(x), \mu_x \rangle_{\mathcal{H}_k} \langle \psi(y), \mu_y \rangle_{\mathcal{H}_\ell}}_{\mathbb{E}_{x'} k(x, x')} \mathbb{E}_{y'} \ell(y, y')}}\end{aligned}$$

$$\begin{aligned}\langle \mathcal{C}_{xy}^u, \mu_x \otimes \mu_y \rangle_{HS} &= \langle \mathbb{E}_{xy} [\varphi(x) \otimes \psi(y)], \mu_x \otimes \mu_y \rangle_{HS} \\ &= \mathbb{E}_{xy} \underbrace{\langle \varphi(x) \otimes \psi(y), \mu_x \otimes \mu_y \rangle_{HS}}_{\underbrace{\langle \varphi(x), \mu_x \rangle_{\mathcal{H}_k} \langle \psi(y), \mu_y \rangle_{\mathcal{H}_\ell}}_{\mathbb{E}_{x'} k(x, x') \quad \mathbb{E}_{y'} \ell(y, y')}} \\ &= \mathbb{E}_{xy} [\mathbb{E}_{x'} k(x, x') \mathbb{E}_{y'} \ell(y, y')].\end{aligned}$$

HSIC: after gathering the terms

$$\begin{aligned}\text{HSIC}^2(x, y) &= \mathbb{E}_{xy} \mathbb{E}_{x'y'} k(x, x') \ell(y, y') + \mathbb{E}_{xx'} k(x, x') \mathbb{E}_{yy'} \ell(y, y') \\ &\quad - 2\mathbb{E}_{xy} [\mathbb{E}_{x'} k(x, x') \mathbb{E}_{y'} \ell(y, y')] . \\ &=: a + b - 2c.\end{aligned}$$

HSIC: after gathering the terms

$$\begin{aligned}\text{HSIC}^2(x, y) &= \mathbb{E}_{xy} \mathbb{E}_{x'y'} k(x, x') \ell(y, y') + \mathbb{E}_{xx'} k(x, x') \mathbb{E}_{yy'} \ell(y, y') \\ &\quad - 2\mathbb{E}_{xy} [\mathbb{E}_{x'} k(x, x') \mathbb{E}_{y'} \ell(y, y')] . \\ &=: a + b - 2c.\end{aligned}$$

Idea: given $\{(x_i, y_i)\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}_{xy}$,

- Let us estimate C_{xy}^u , μ_x , μ_y empirically.

HSIC: after gathering the terms

$$\begin{aligned}\text{HSIC}^2(x, y) &= \mathbb{E}_{xy} \mathbb{E}_{x'y'} k(x, x') \ell(y, y') + \mathbb{E}_{xx'} k(x, x') \mathbb{E}_{yy'} \ell(y, y') \\ &\quad - 2\mathbb{E}_{xy} [\mathbb{E}_{x'} k(x, x') \mathbb{E}_{y'} \ell(y, y')] . \\ &=: a + b - 2c.\end{aligned}$$

Idea: given $\{(x_i, y_i)\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}_{xy}$,

- Let us estimate C_{xy}^u , μ_x , μ_y empirically.

Result

$$\widehat{\text{HSIC}}_b^2(x, y) = \frac{1}{n^2} \left\langle \tilde{\mathbf{G}}_x, \tilde{\mathbf{G}}_y \right\rangle_F : \text{see the intuition. The details...}$$

HSIC estimation: from $\widehat{C}_{xy}^u, \hat{\mu}_x, \hat{\mu}_y$

First term:

$$\textcolor{blue}{a} = \|C_{xy}^u\|_{HS}^2 = \mathbb{E}_{xy}\mathbb{E}_{x'y'} k(x, x')\ell(y, y'),$$

$$\hat{a} = \|\widehat{C}_{xy}^u\|_{HS}^2 =$$

HSIC estimation: from $\widehat{C}_{xy}^u, \hat{\mu}_x, \hat{\mu}_y$

First term:

$$\textcolor{blue}{a} = \|C_{xy}^u\|_{HS}^2 = \mathbb{E}_{xy}\mathbb{E}_{x'y'} k(x, x')\ell(y, y'),$$

$$\hat{a} = \|\widehat{C}_{xy}^u\|_{HS}^2 = \left\langle \frac{1}{n} \sum_{i=1}^n \varphi(x_i) \otimes \psi(y_i), \frac{1}{n} \sum_{j=1}^n \varphi(x_j) \otimes \psi(y_j) \right\rangle_{HS}$$

HSIC estimation: from $\widehat{C}_{xy}^u, \hat{\mu}_x, \hat{\mu}_y$

First term:

$$\textcolor{blue}{a} = \|C_{xy}^u\|_{HS}^2 = \mathbb{E}_{xy}\mathbb{E}_{x'y'} k(x, x')\ell(y, y'),$$

$$\hat{a} = \|\widehat{C}_{xy}^u\|_{HS}^2 = \left\langle \frac{1}{n} \sum_{i=1}^n \varphi(x_i) \otimes \psi(y_i), \frac{1}{n} \sum_{j=1}^n \varphi(x_j) \otimes \psi(y_j) \right\rangle_{HS}$$

$$= \frac{1}{n^2} \sum_{i,j=1}^n (\mathbf{G}_x)_{ij} (\mathbf{G}_y)_{ij}$$

HSIC estimation: from \widehat{C}_{xy}^u , $\hat{\mu}_x$, $\hat{\mu}_y$

First term:

$$\textcolor{blue}{a} = \|C_{xy}^u\|_{HS}^2 = \mathbb{E}_{xy}\mathbb{E}_{x'y'} k(x, x')\ell(y, y'),$$

$$\hat{a} = \|\widehat{C}_{xy}^u\|_{HS}^2 = \left\langle \frac{1}{n} \sum_{i=1}^n \varphi(x_i) \otimes \psi(y_i), \frac{1}{n} \sum_{j=1}^n \varphi(x_j) \otimes \psi(y_j) \right\rangle_{HS}$$

$$= \frac{1}{n^2} \sum_{i,j=1}^n (\mathbf{G}_x)_{ij} (\mathbf{G}_y)_{ij} = \frac{1}{n^2} \langle \mathbf{G}_x, \mathbf{G}_y \rangle_F = \frac{1}{n^2} \text{tr}(\mathbf{G}_x \mathbf{G}_y).$$

HSIC estimation: 2nd term

$$\textcolor{blue}{b} = \|\mu_x \otimes \mu_y\|_{HS}^2 = \mathbb{E}_{xx'} k(x, x') \mathbb{E}_{yy'} \ell(y, y').$$

$$\hat{b} = \|\hat{\mu}_x \otimes \hat{\mu}_y\|_{HS}^2$$

HSIC estimation: 2nd term

$$b = \|\mu_x \otimes \mu_y\|_{HS}^2 = \mathbb{E}_{xx'} k(x, x') \mathbb{E}_{yy'} \ell(y, y').$$

$$\hat{b} = \|\hat{\mu}_x \otimes \hat{\mu}_y\|_{HS}^2 = \langle \hat{\mu}_x \otimes \hat{\mu}_y, \hat{\mu}_x \otimes \hat{\mu}_y \rangle_{HS}$$

HSIC estimation: 2nd term

$$\color{blue}{b} = \|\mu_x \otimes \mu_y\|_{HS}^2 = \mathbb{E}_{xx'} k(x, x') \mathbb{E}_{yy'} \ell(y, y').$$

$$\begin{aligned}\hat{b} &= \|\hat{\mu}_x \otimes \hat{\mu}_y\|_{HS}^2 = \langle \hat{\mu}_x \otimes \hat{\mu}_y, \hat{\mu}_x \otimes \hat{\mu}_y \rangle_{HS} \\ &= \left\langle \left[\frac{1}{n} \sum_{i=1}^n \varphi(x_i) \right] \otimes \left[\frac{1}{n} \sum_{j=1}^n \psi(y_j) \right], \left[\frac{1}{n} \sum_{i=1}^n \varphi(x_i) \right] \otimes \left[\frac{1}{n} \sum_{j=1}^n \psi(y_j) \right] \right\rangle_{HS}\end{aligned}$$

HSIC estimation: 2nd term

$$\color{blue}{b} = \|\mu_x \otimes \mu_y\|_{HS}^2 = \mathbb{E}_{xx'} k(x, x') \mathbb{E}_{yy'} \ell(y, y').$$

$$\begin{aligned}\hat{b} &= \|\hat{\mu}_x \otimes \hat{\mu}_y\|_{HS}^2 = \langle \hat{\mu}_x \otimes \hat{\mu}_y, \hat{\mu}_x \otimes \hat{\mu}_y \rangle_{HS} \\ &= \left\langle \left[\frac{1}{n} \sum_{i=1}^n \varphi(x_i) \right] \otimes \left[\frac{1}{n} \sum_{j=1}^n \psi(y_j) \right], \left[\frac{1}{n} \sum_{i=1}^n \varphi(x_i) \right] \otimes \left[\frac{1}{n} \sum_{j=1}^n \psi(y_j) \right] \right\rangle_{HS} \\ &= \left[\frac{1}{n^2} \sum_{i,j=1}^n k(x_i, x_j) \right] \left[\frac{1}{n^2} \sum_{i,j=1}^n \ell(x_i, x_j) \right]\end{aligned}$$

HSIC estimation: 2nd term

$$\begin{aligned}\textcolor{blue}{b} &= \|\mu_x \otimes \mu_y\|_{HS}^2 = \mathbb{E}_{xx'} k(x, x') \mathbb{E}_{yy'} \ell(y, y'). \\ \hat{b} &= \|\hat{\mu}_x \otimes \hat{\mu}_y\|_{HS}^2 = \langle \hat{\mu}_x \otimes \hat{\mu}_y, \hat{\mu}_x \otimes \hat{\mu}_y \rangle_{HS} \\ &= \left\langle \left[\frac{1}{n} \sum_{i=1}^n \varphi(x_i) \right] \otimes \left[\frac{1}{n} \sum_{j=1}^n \psi(y_j) \right], \left[\frac{1}{n} \sum_{i=1}^n \varphi(x_i) \right] \otimes \left[\frac{1}{n} \sum_{j=1}^n \psi(y_j) \right] \right\rangle_{HS} \\ &= \left[\frac{1}{n^2} \sum_{i,j=1}^n k(x_i, x_j) \right] \left[\frac{1}{n^2} \sum_{i,j=1}^n \ell(x_i, x_j) \right] = \frac{1}{n^4} (\mathbf{1}^\top \mathbf{G}_x \mathbf{1}) (\mathbf{1}^\top \mathbf{G}_y \mathbf{1}).\end{aligned}$$

HSIC estimation: 3rd term (without '-2')

$$c = \langle C_{xy}^u, \mu_x \otimes \mu_y \rangle_{HS},$$

$$\hat{c} = \left\langle \widehat{C}_{xy}^u, \hat{\mu}_x \otimes \hat{\mu}_y \right\rangle_{HS}$$

HSIC estimation: 3rd term (without '-2')

$$c = \langle C_{xy}^u, \mu_x \otimes \mu_y \rangle_{HS},$$

$$\hat{c} = \left\langle \widehat{C}_{xy}^u, \hat{\mu}_x \otimes \hat{\mu}_y \right\rangle_{HS}$$

$$= \left\langle \frac{1}{n} \sum_{i=1}^n \varphi(x_i) \otimes \psi(y_i), \left[\frac{1}{n} \sum_{a=1}^n \varphi(x_a) \right] \otimes \left[\frac{1}{n} \sum_{b=1}^n \psi(y_b) \right] \right\rangle_{HS}$$

HSIC estimation: 3rd term (without '-2')

$$c = \langle C_{xy}^u, \mu_x \otimes \mu_y \rangle_{HS},$$

$$\hat{c} = \left\langle \widehat{C}_{xy}^u, \hat{\mu}_x \otimes \hat{\mu}_y \right\rangle_{HS}$$

$$= \left\langle \frac{1}{n} \sum_{i=1}^n \varphi(x_i) \otimes \psi(y_i), \left[\frac{1}{n} \sum_{a=1}^n \varphi(x_a) \right] \otimes \left[\frac{1}{n} \sum_{b=1}^n \psi(y_b) \right] \right\rangle_{HS}$$

$$= \frac{1}{n} \sum_{i=1}^n \left\langle \varphi(x_i) \otimes \psi(y_i), \left[\frac{1}{n} \sum_{a=1}^n \varphi(x_a) \right] \otimes \left[\frac{1}{n} \sum_{b=1}^n \psi(y_b) \right] \right\rangle_{HS}$$

HSIC estimation: 3rd term (without '-2')

$$c = \langle C_{xy}^u, \mu_x \otimes \mu_y \rangle_{HS},$$

$$\hat{c} = \left\langle \widehat{C}_{xy}^u, \widehat{\mu}_x \otimes \widehat{\mu}_y \right\rangle_{HS}$$

$$= \left\langle \frac{1}{n} \sum_{i=1}^n \varphi(x_i) \otimes \psi(y_i), \left[\frac{1}{n} \sum_{a=1}^n \varphi(x_a) \right] \otimes \left[\frac{1}{n} \sum_{b=1}^n \psi(y_b) \right] \right\rangle_{HS}$$

$$= \frac{1}{n} \sum_{i=1}^n \left\langle \varphi(x_i) \otimes \psi(y_i), \left[\frac{1}{n} \sum_{a=1}^n \varphi(x_a) \right] \otimes \left[\frac{1}{n} \sum_{b=1}^n \psi(y_b) \right] \right\rangle_{HS}$$

$$= \underbrace{\frac{1}{n} \sum_{i=1}^n \left\langle \varphi(x_i), \frac{1}{n} \sum_{a=1}^n \varphi(x_a) \right\rangle_{\mathcal{H}_k}}_{\frac{1}{n} \sum_{a=1}^n k(x_i, x_a)} \underbrace{\left\langle \psi(y_i), \frac{1}{n} \sum_{b=1}^n \psi(y_b) \right\rangle_{\mathcal{H}_\ell}}_{\frac{1}{n} \sum_{b=1}^n \ell(y_i, y_b)}$$

HSIC estimation: 3rd term (without '-2')

$$c = \langle C_{xy}^u, \mu_x \otimes \mu_y \rangle_{HS},$$

$$\hat{c} = \left\langle \widehat{C}_{xy}^u, \hat{\mu}_x \otimes \hat{\mu}_y \right\rangle_{HS}$$

$$= \left\langle \frac{1}{n} \sum_{i=1}^n \varphi(x_i) \otimes \psi(y_i), \left[\frac{1}{n} \sum_{a=1}^n \varphi(x_a) \right] \otimes \left[\frac{1}{n} \sum_{b=1}^n \psi(y_b) \right] \right\rangle_{HS}$$

$$= \frac{1}{n} \sum_{i=1}^n \left\langle \varphi(x_i) \otimes \psi(y_i), \left[\frac{1}{n} \sum_{a=1}^n \varphi(x_a) \right] \otimes \left[\frac{1}{n} \sum_{b=1}^n \psi(y_b) \right] \right\rangle_{HS}$$

$$= \frac{1}{n} \sum_{i=1}^n \underbrace{\left\langle \varphi(x_i), \frac{1}{n} \sum_{a=1}^n \varphi(x_a) \right\rangle_{\mathcal{H}_k}}_{\frac{1}{n} \sum_{a=1}^n k(x_i, x_a)} \underbrace{\left\langle \psi(y_i), \frac{1}{n} \sum_{b=1}^n \psi(y_b) \right\rangle_{\mathcal{H}_\ell}}_{\frac{1}{n} \sum_{b=1}^n \ell(y_i, y_b)}$$

$$= \frac{1}{n^3} \sum_{a,b=1}^n \underbrace{\left[\sum_{i=1}^n k(x_i, x_a) \ell(y_i, y_b) \right]}_{(\mathbf{G}_x \mathbf{G}_y)_{a,b}}$$

HSIC estimation: 3rd term (without '-2')

$$c = \langle C_{xy}^u, \mu_x \otimes \mu_y \rangle_{HS},$$

$$\hat{c} = \left\langle \widehat{C}_{xy}^u, \hat{\mu}_x \otimes \hat{\mu}_y \right\rangle_{HS}$$

$$= \left\langle \frac{1}{n} \sum_{i=1}^n \varphi(x_i) \otimes \psi(y_i), \left[\frac{1}{n} \sum_{a=1}^n \varphi(x_a) \right] \otimes \left[\frac{1}{n} \sum_{b=1}^n \psi(y_b) \right] \right\rangle_{HS}$$

$$= \frac{1}{n} \sum_{i=1}^n \left\langle \varphi(x_i) \otimes \psi(y_i), \left[\frac{1}{n} \sum_{a=1}^n \varphi(x_a) \right] \otimes \left[\frac{1}{n} \sum_{b=1}^n \psi(y_b) \right] \right\rangle_{HS}$$

$$= \frac{1}{n} \sum_{i=1}^n \underbrace{\left\langle \varphi(x_i), \frac{1}{n} \sum_{a=1}^n \varphi(x_a) \right\rangle_{\mathcal{H}_k}}_{\frac{1}{n} \sum_{a=1}^n k(x_i, x_a)} \underbrace{\left\langle \psi(y_i), \frac{1}{n} \sum_{b=1}^n \psi(y_b) \right\rangle_{\mathcal{H}_\ell}}_{\frac{1}{n} \sum_{b=1}^n \ell(y_i, y_b)}$$

$$= \frac{1}{n^3} \sum_{a,b=1}^n \underbrace{\left[\sum_{i=1}^n k(x_i, x_a) \ell(y_i, y_b) \right]}_{(\mathbf{G}_x \mathbf{G}_y)_{a,b}} = \frac{1}{n^3} \mathbf{1}^T \mathbf{G}_x \mathbf{G}_y \mathbf{1}.$$

HSIC estimation: putting together

$$\widehat{\text{HSIC}}_b^2(x, y) =: \hat{a} + \hat{b} - 2\hat{c}$$

HSIC estimation: putting together

$$\begin{aligned}\widehat{\text{HSIC}}_b^2(x, y) &=: \hat{a} + \hat{b} - 2\hat{c} \\ &= \frac{1}{n^2} \text{tr}(\mathbf{G}_x \mathbf{G}_y) + \frac{1}{n^4} (\mathbf{1}^T \mathbf{G}_x \mathbf{1}) (\mathbf{1}^T \mathbf{G}_y \mathbf{1}) - \frac{2}{n^3} \mathbf{1}^T \mathbf{G}_x \mathbf{G}_y \mathbf{1}\end{aligned}$$

HSIC estimation: putting together

$$\begin{aligned}\widehat{\text{HSIC}}_b^2(x, y) &=: \hat{a} + \hat{b} - 2\hat{c} \\&= \frac{1}{n^2} \text{tr}(\mathbf{G}_x \mathbf{G}_y) + \frac{1}{n^4} (\mathbf{1}^T \mathbf{G}_x \mathbf{1}) (\mathbf{1}^T \mathbf{G}_y \mathbf{1}) - \frac{2}{n^3} \mathbf{1}^T \mathbf{G}_x \mathbf{G}_y \mathbf{1} \\&= \frac{1}{n^2} \text{tr} \left(\underbrace{\mathbf{G}_x \mathbf{G}_y - \frac{1}{n} \mathbf{1} \mathbf{1}^T \mathbf{G}_x \mathbf{G}_y - \frac{1}{n} \mathbf{1} \mathbf{1}^T \mathbf{G}_x \mathbf{G}_y + \frac{1}{n^2} \mathbf{1} \mathbf{1}^T \mathbf{G}_x \mathbf{1} \mathbf{1}^T \mathbf{G}_y}_{\left(\mathbf{I}_n - \frac{\mathbf{E}_n}{n} \right) \mathbf{G}_x \left(\mathbf{I}_n - \frac{\mathbf{E}_n}{n} \right) \mathbf{G}_y} \right)\end{aligned}$$

HSIC estimation: putting together

$$\begin{aligned}\widehat{\text{HSIC}}_b^2(x, y) &=: \hat{a} + \hat{b} - 2\hat{c} \\&= \frac{1}{n^2} \text{tr}(\mathbf{G}_x \mathbf{G}_y) + \frac{1}{n^4} (\mathbf{1}^T \mathbf{G}_x \mathbf{1}) (\mathbf{1}^T \mathbf{G}_y \mathbf{1}) - \frac{2}{n^3} \mathbf{1}^T \mathbf{G}_x \mathbf{G}_y \mathbf{1} \\&= \frac{1}{n^2} \text{tr} \left(\underbrace{\mathbf{G}_x \mathbf{G}_y - \frac{1}{n} \mathbf{1} \mathbf{1}^T \mathbf{G}_x \mathbf{G}_y - \frac{1}{n} \mathbf{1} \mathbf{1}^T \mathbf{G}_x \mathbf{G}_y + \frac{1}{n^2} \mathbf{1} \mathbf{1}^T \mathbf{G}_x \mathbf{1} \mathbf{1}^T \mathbf{G}_y}_{\left(\mathbf{I}_n - \frac{\mathbf{1} \mathbf{1}^T}{n} \right) \mathbf{G}_x \left(\mathbf{I}_n - \frac{\mathbf{1} \mathbf{1}^T}{n} \right) \mathbf{G}_y} \right) \\&= \frac{1}{n^2} \text{tr} (\mathbf{H} \mathbf{G}_x \mathbf{H} \mathbf{G}_y)\end{aligned}$$

HSIC estimation: putting together

$$\begin{aligned}\widehat{\text{HSIC}}_b^2(x, y) &=: \hat{a} + \hat{b} - 2\hat{c} \\&= \frac{1}{n^2} \text{tr}(\mathbf{G}_x \mathbf{G}_y) + \frac{1}{n^4} (\mathbf{1}^T \mathbf{G}_x \mathbf{1}) (\mathbf{1}^T \mathbf{G}_y \mathbf{1}) - \frac{2}{n^3} \mathbf{1}^T \mathbf{G}_x \mathbf{G}_y \mathbf{1} \\&= \frac{1}{n^2} \text{tr} \left(\underbrace{\mathbf{G}_x \mathbf{G}_y - \frac{1}{n} \mathbf{1} \mathbf{1}^T \mathbf{G}_x \mathbf{G}_y - \frac{1}{n} \mathbf{1} \mathbf{1}^T \mathbf{G}_x \mathbf{G}_y + \frac{1}{n^2} \mathbf{1} \mathbf{1}^T \mathbf{G}_x \mathbf{1} \mathbf{1}^T \mathbf{G}_y}_{\left(\mathbf{I}_n - \frac{\mathbf{E}_n}{n} \right) \mathbf{G}_x \left(\mathbf{I}_n - \frac{\mathbf{E}_n}{n} \right) \mathbf{G}_y} \right) \\&= \frac{1}{n^2} \text{tr} (\mathbf{H} \mathbf{G}_x \mathbf{H} \mathbf{G}_y) = \frac{1}{n^2} \text{tr} \left(\underbrace{\mathbf{H} \mathbf{G}_x \mathbf{H}}_{\tilde{\mathbf{G}}_x} \underbrace{\mathbf{H} \mathbf{G}_y \mathbf{H}}_{\tilde{\mathbf{G}}_y} \right)\end{aligned}$$

HSIC estimation: putting together

$$\begin{aligned}\widehat{\text{HSIC}}_b^2(x, y) &=: \hat{a} + \hat{b} - 2\hat{c} \\&= \frac{1}{n^2} \text{tr}(\mathbf{G}_x \mathbf{G}_y) + \frac{1}{n^4} (\mathbf{1}^T \mathbf{G}_x \mathbf{1}) (\mathbf{1}^T \mathbf{G}_y \mathbf{1}) - \frac{2}{n^3} \mathbf{1}^T \mathbf{G}_x \mathbf{G}_y \mathbf{1} \\&= \frac{1}{n^2} \text{tr} \left(\underbrace{\mathbf{G}_x \mathbf{G}_y - \frac{1}{n} \mathbf{1} \mathbf{1}^T \mathbf{G}_x \mathbf{G}_y - \frac{1}{n} \mathbf{1} \mathbf{1}^T \mathbf{G}_x \mathbf{G}_y + \frac{1}{n^2} \mathbf{1} \mathbf{1}^T \mathbf{G}_x \mathbf{1} \mathbf{1}^T \mathbf{G}_y}_{\left(\mathbf{I}_n - \frac{\mathbf{E}_n}{n} \right) \mathbf{G}_x \left(\mathbf{I}_n - \frac{\mathbf{E}_n}{n} \right) \mathbf{G}_y} \right) \\&= \frac{1}{n^2} \text{tr} (\mathbf{H} \mathbf{G}_x \mathbf{H} \mathbf{G}_y) = \frac{1}{n^2} \text{tr} \left(\underbrace{\mathbf{H} \mathbf{G}_x \mathbf{H}}_{\tilde{\mathbf{G}}_x} \underbrace{\mathbf{H} \mathbf{G}_y \mathbf{H}}_{\tilde{\mathbf{G}}_y} \right) = \frac{1}{n^2} \langle \tilde{\mathbf{G}}_x, \tilde{\mathbf{G}}_y \rangle_F.\end{aligned}$$

HSIC estimation: putting together

$$\begin{aligned}\widehat{\text{HSIC}}_b^2(x, y) &=: \hat{a} + \hat{b} - 2\hat{c} \\ &= \frac{1}{n^2} \text{tr}(\mathbf{G}_x \mathbf{G}_y) + \frac{1}{n^4} (\mathbf{1}^\top \mathbf{G}_x \mathbf{1}) (\mathbf{1}^\top \mathbf{G}_y \mathbf{1}) - \frac{2}{n^3} \mathbf{1}^\top \mathbf{G}_x \mathbf{G}_y \mathbf{1} \\ &= \frac{1}{n^2} \text{tr} \left(\underbrace{\mathbf{G}_x \mathbf{G}_y - \frac{1}{n} \mathbf{1} \mathbf{1}^\top \mathbf{G}_x \mathbf{G}_y - \frac{1}{n} \mathbf{1} \mathbf{1}^\top \mathbf{G}_x \mathbf{G}_y + \frac{1}{n^2} \mathbf{1} \mathbf{1}^\top \mathbf{G}_x \mathbf{1} \mathbf{1}^\top \mathbf{G}_y}_{\left(\mathbf{I}_n - \frac{\mathbf{E}_n}{n} \right) \mathbf{G}_x \left(\mathbf{I}_n - \frac{\mathbf{E}_n}{n} \right) \mathbf{G}_y} \right) \\ &= \frac{1}{n^2} \text{tr} (\mathbf{H} \mathbf{G}_x \mathbf{H} \mathbf{G}_y) = \frac{1}{n^2} \text{tr} \left(\underbrace{\mathbf{H} \mathbf{G}_x \mathbf{H}}_{\tilde{\mathbf{G}}_x} \underbrace{\mathbf{H} \mathbf{G}_y \mathbf{H}}_{\tilde{\mathbf{G}}_y} \right) = \frac{1}{n^2} \langle \tilde{\mathbf{G}}_x, \tilde{\mathbf{G}}_y \rangle_F.\end{aligned}$$

Bias of $\widehat{\text{HSIC}}_b$: $\mathcal{O}\left(\frac{1}{n}\right)$.

Reminder: MMD^2 , $\widehat{\text{MMD}}_b^2$, $\widehat{\text{MMD}}_u^2$

$$\text{MMD}^2(\mathbb{P}, \mathbb{Q}) := \mathbb{E}_{\mathbf{x}\mathbf{x}'} k(x, x') + \mathbb{E}_{\mathbf{y}\mathbf{y}'} k(y, y') - 2\mathbb{E}_{\mathbf{x}\mathbf{y}} k(x, y),$$

$$\begin{aligned} \widehat{\text{MMD}}_b^2(\mathbb{P}, \mathbb{Q}) &= \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m k(x_i, x_j) + \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n k(y_i, y_j) \\ &\quad - \frac{2}{mn} \sum_{i=1}^m \sum_{j=1}^n k(x_i, y_j), \end{aligned}$$

$$\begin{aligned} \widehat{\text{MMD}}_u^2(\mathbb{P}, \mathbb{Q}) &= \frac{1}{m(m-1)} \sum_{i=1}^m \sum_{j \neq i}^m k(x_i, x_j) + \frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j \neq i}^n k(y_i, y_j) \\ &\quad - \frac{2}{mn} \sum_{i=1}^m \sum_{j=1}^n k(x_i, y_j). \end{aligned}$$

$\widehat{\text{HSIC}}_b^2$ until now

$$\begin{aligned}\text{HSIC}^2(x, y) &= \mathbb{E}_{xy} \mathbb{E}_{x'y'} k(\textcolor{blue}{x}, \textcolor{blue}{x}') \ell(y, y') + \mathbb{E}_{xx'} k(x, x') \mathbb{E}_{yy'} \ell(y, y') \\ &\quad - 2\mathbb{E}_{xy} [\mathbb{E}_{x'} k(x, x') \mathbb{E}_{y'} \ell(y, y')],\end{aligned}$$

$$\widehat{\text{HSIC}}_b^2(x, y) = \frac{1}{n^2} \sum_{i,j=1}^n k(\textcolor{blue}{x}_i, \textcolor{blue}{x}_j) \ell(y_i, y_j) + \dots$$

$\widehat{\text{HSIC}}_b^2$ until now

$$\begin{aligned}\text{HSIC}^2(x, y) &= \mathbb{E}_{xy} \mathbb{E}_{x'y'} k(\textcolor{blue}{x}, \textcolor{blue}{x}') \ell(y, y') + \mathbb{E}_{xx'} k(x, x') \mathbb{E}_{yy'} \ell(y, y') \\ &\quad - 2\mathbb{E}_{xy} [\mathbb{E}_{x'} k(x, x') \mathbb{E}_{y'} \ell(y, y')],\end{aligned}$$

$$\widehat{\text{HSIC}}_b^2(x, y) = \frac{1}{n^2} \sum_{i,j=1}^n k(\textcolor{blue}{x}_i, \textcolor{blue}{x}_j) \ell(y_i, y_j) + \dots$$

- $\textcolor{blue}{x}, \textcolor{blue}{x}'$ should be independent, but
- with plug-in: $i = j$, it introduces **bias**.

HSIC: unbiased estimator

Idea: get rid of the $i = j$ -type terms. Let $k_{ij} := k(x_i, x_j)$, $\ell_{ij} := \ell(y_i, y_j)$.

$$\hat{a}_{\mathbf{b}} = \frac{1}{n^2} \sum_{i,j=1}^n k_{ij} \ell_{ij},$$

HSIC: unbiased estimator

Idea: get rid of the $i = j$ -type terms. Let $k_{ij} := k(x_i, x_j)$, $\ell_{ij} := \ell(y_i, y_j)$.

$$\hat{a}_b = \frac{1}{n^2} \sum_{i,j=1}^n k_{ij} \ell_{ij}, \quad \hat{a}_u = \frac{1}{n(n-1)} \sum_{i \neq j} k_{ij} \ell_{ij}$$

HSIC: unbiased estimator

Idea: get rid of the $i = j$ -type terms. Let $k_{ij} := k(x_i, x_j)$, $\ell_{ij} := \ell(y_i, y_j)$.

$$\hat{a}_b = \frac{1}{n^2} \sum_{i,j=1}^n k_{ij} \ell_{ij}, \quad \hat{a}_u = \underbrace{\frac{1}{n(n-1)} \sum_{i \neq j} k_{ij} \ell_{ij}}_{\frac{1}{(n)_2} \sum_{(i,j) \in I_2^n} k_{ij} \ell_{ij}}$$

$$I_p^n = \{(i_1, \dots, i_p) : i_j \in \{1, \dots, n\} \text{ without replacement}\}, \quad (n)_p = |I_p^n|.$$

HSIC: unbiased estimator

Idea: get rid of the $i = j$ -type terms. Let $k_{ij} := k(x_i, x_j)$, $\ell_{ij} := \ell(y_i, y_j)$.

$$\hat{a}_b = \frac{1}{n^2} \sum_{i,j=1}^n k_{ij} \ell_{ij}, \quad \hat{a}_u = \underbrace{\frac{1}{n(n-1)} \sum_{i \neq j} k_{ij} \ell_{ij}}_{\frac{1}{(n)_2} \sum_{(i,j) \in I_2^n} k_{ij} \ell_{ij}}$$

$$\hat{c}_b = \frac{1}{n^3} \sum_{i,q,r=1}^n k_{iq} \ell_{ir},$$

$$I_p^n = \{(i_1, \dots, i_p) : i_j \in \{1, \dots, n\} \text{ without replacement}\}, \quad (n)_p = |I_p^n|.$$

HSIC: unbiased estimator

Idea: get rid of the $i = j$ -type terms. Let $k_{ij} := k(x_i, x_j)$, $\ell_{ij} := \ell(y_i, y_j)$.

$$\hat{a}_b = \frac{1}{n^2} \sum_{i,j=1}^n k_{ij} \ell_{ij}, \quad \hat{a}_u = \underbrace{\frac{1}{n(n-1)} \sum_{i \neq j} k_{ij} \ell_{ij}}_{\frac{1}{(n)_2} \sum_{(i,j) \in I_2^n} k_{ij} \ell_{ij}}$$

$$\hat{c}_b = \frac{1}{n^3} \sum_{i,q,r=1}^n k_{iq} \ell_{ir}, \quad \hat{c}_u = \frac{1}{(n)_3} \sum_{(i,q,r) \in I_3^n}^n k_{iq} \ell_{ir},$$

$$I_p^n = \{(i_1, \dots, i_p) : i_j \in \{1, \dots, n\} \text{ without replacement}\}, \quad (n)_p = |I_p^n|.$$

HSIC: unbiased estimator

Idea: get rid of the $i = j$ -type terms. Let $k_{ij} := k(x_i, x_j)$, $\ell_{ij} := \ell(y_i, y_j)$.

$$\hat{a}_b = \frac{1}{n^2} \sum_{i,j=1}^n k_{ij} \ell_{ij}, \quad \hat{a}_u = \underbrace{\frac{1}{n(n-1)} \sum_{i \neq j} k_{ij} \ell_{ij}}_{\frac{1}{(n)_2} \sum_{(i,j) \in I_2^n} k_{ij} \ell_{ij}}$$

$$\hat{c}_b = \frac{1}{n^3} \sum_{i,q,r=1}^n k_{iq} \ell_{ir}, \quad \hat{c}_u = \frac{1}{(n)_3} \sum_{(i,q,r) \in I_3^n}^n k_{iq} \ell_{ir},$$

$$\hat{b}_b = \frac{1}{n^4} \sum_{i,j,q,r=1}^n k_{ij} \ell_{qr},$$

$$I_p^n = \{(i_1, \dots, i_p) : i_j \in \{1, \dots, n\} \text{ without replacement}\}, \quad (n)_p = |I_p^n|.$$

HSIC: unbiased estimator

Idea: get rid of the $i = j$ -type terms. Let $k_{ij} := k(x_i, x_j)$, $\ell_{ij} := \ell(y_i, y_j)$.

$$\hat{a}_b = \frac{1}{n^2} \sum_{i,j=1}^n k_{ij} \ell_{ij},$$

$$\hat{a}_u = \underbrace{\frac{1}{n(n-1)} \sum_{i \neq j} k_{ij} \ell_{ij}}_{\frac{1}{(n)_2} \sum_{(i,j) \in I_2^n} k_{ij} \ell_{ij}},$$

$$\hat{c}_b = \frac{1}{n^3} \sum_{i,q,r=1}^n k_{iq} \ell_{ir},$$

$$\hat{c}_u = \frac{1}{(n)_3} \sum_{(i,q,r) \in I_3^n} k_{iq} \ell_{ir},$$

$$\hat{b}_b = \frac{1}{n^4} \sum_{i,j,q,r=1}^n k_{ij} \ell_{qr},$$

$$\hat{b}_u = \frac{1}{(n)_4} \sum_{(i,j,q,r) \in I_4^n} k_{ij} \ell_{qr}.$$

$$I_p^n = \{(i_1, \dots, i_p) : i_j \in \{1, \dots, n\} \text{ without replacement}\}, (n)_p = |I_p^n|.$$

HSIC: resulting unbiased estimator

After some linear algebra [Gretton et al., 2005a], $(M)_{++} := \sum_{i,j} M_{ij}$,

$$\widehat{\text{HSIC}}_b^2(x, y) = \frac{1}{n^2} \left\langle \tilde{\mathbf{G}}_x, \tilde{\mathbf{G}}_y \right\rangle_F,$$

$$\begin{aligned} \widehat{\text{HSIC}}_u^2(x, y) &= \frac{1}{n(n-3)} \left[\left\langle \tilde{\mathbf{G}}_x, \tilde{\mathbf{G}}_y \right\rangle_F - \frac{2}{n-2} (\tilde{\mathbf{G}}_x \tilde{\mathbf{G}}_y)_{++} \right. \\ &\quad \left. + \frac{1}{(n-1)(n-2)} (\tilde{\mathbf{G}}_x)_{++} (\tilde{\mathbf{G}}_y)_{++} \right]. \end{aligned}$$

Estimation in practice: few ITE examples

(<https://bitbucket.org/szzoli/ite/>)

(<https://bitbucket.org/szzoli/ite-in-python/>)

KCCA estimation: Matlab

Goal: estimate KCCA,

```
>ds = [2;3;4]; Y = rand(sum(ds),5000);  
>mult = 1;  
>co = IKCCA_initialization(mult);  
>KCCA = IKCCA_estimation(Y,ds,co);
```

KCCA estimation: Matlab

Goal: estimate KCCA,

```
>ds = [2;3;4]; Y = rand(sum(ds),5000);  
>mult = 1;  
>co = IKCCA_initialization(mult);  
>KCCA = IKCCA_estimation(Y,ds,co);
```

Alternative initialization:

```
>co = IKCCA_initialization(mult,{’kappa’,0.01,’eta’,0.001});  
where  $\kappa$ : regularization constant,  $\eta$ : low-rank approximation.
```

KCCA & HSIC estimation: Matlab

Goal: estimate KCCA,

```
>ds = [2;3;4]; Y = rand(sum(ds),5000);  
>mult = 1;  
>co = IKCCA_initialization(mult);  
>KCCA = IKCCA_estimation(Y,ds,co);
```

Alternative initialization:

```
>co = IKCCA_initialization(mult,{'kappa',0.01,'eta',0.001});  
where  $\kappa$ : regularization constant,  $\eta$ : low-rank approximation.
```

Note: HSIC similarly.

MMD estimation: Matlab

Using for example U-statistic:

```
>X1 = randn(3,2000); X2 = randn(3,3000);
>mult = 1;
>co = DMMD_Ustat_initialization(mult);
>MMD = DMMD_Ustat_estimation(X1,X2,co);
```

MMD estimation: Matlab

Using for example U-statistic:

```
>X1 = randn(3,2000); X2 = randn(3,3000);
>mult = 1;
>co = DMMD_Ustat_initialization(mult);
>MMD = DMMD_Ustat_estimation(X1,X2,co);
```

With low-rank approximation, and setting some parameters:

```
co2 = DMMD_Ustat_iChol_initialization(mult)
co3 = DMMD_Ustat_iChol_initialization(mult,{'sigma',0.2,
'eta',0.01})
```

HSIC estimation: Python

Import ITE (1x), generate observations:

```
>>> import ite
>>> from numpy.random import randn
>>> from numpy import array
>>> ds = array([2, 3, 4])
>>> t = 1000
>>> y = randn(t, sum(ds))
```

HSIC estimation: Python

Import ITE (1x), generate observations:

```
>>> import ite
>>> from numpy.random import randn
>>> from numpy import array
>>> ds = array([2, 3, 4])
>>> t = 1000
>>> y = randn(t, sum(ds))
```

Estimate HSIC:

```
>>> co = ite.cost.BIHSIC_IChol()
>>> hsic = co.estimation(y, ds)
```

HSIC estimation: Python

Alternative initialization-1:

```
>>> co2 = ite.cost.BIHSIC_IChol(eta=1e-3)
>>> hsic2 = co2.estimation(y, ds)
```

HSIC estimation: Python

Alternative initialization-1:

```
>>> co2 = ite.cost.BIHSIC_IChol(eta=1e-3)
>>> hsic2 = co2.estimation(y, ds)
```

Alternative-2:

```
>>> from ite.cost.x_kernel import Kernel
>>> k = Kernel({'name': 'RBF', 'sigma': 1})
>>> co3 = ite.cost.BIHSIC_IChol(kernel=k, eta=1e-3)
>>> hsic3 = co3.estimation(y, ds)
```

HSIC & KCCA estimation: Python

Alternative initialization-1:

```
>>> co2 = ite.cost.BIHSIC_IChol(eta=1e-3)
>>> hsic2 = co2.estimation(y, ds)
```

Alternative-2:

```
>>> from ite.cost.x_kernel import Kernel
>>> k = Kernel({'name': 'RBF', 'sigma': 1})
>>> co3 = ite.cost.BIHSIC_IChol(kernel=k, eta=1e-3)
>>> hsic3 = co3.estimation(y, ds)
```

Note: KCCA similarly.

MMD estimation: Python

Import ITE, generate observations:

```
>>> import ite
>>> from numpy.random import randn
>>> dim = 3
>>> t1, t2 = 2000, 3000
>>> y1 = randn(t1, dim)
>>> y2 = randn(t2, dim)
```

MMD estimation: Python

Import ITE, generate observations:

```
>>> import ite
>>> from numpy.random import randn
>>> dim = 3
>>> t1, t2 = 2000, 3000
>>> y1 = randn(t1, dim)
>>> y2 = randn(t2, dim)
```

Estimate MMD:

```
>>> co = ite.cost.BDMMD_UStat_IChol()
>>> mmd = co.estimation(y1, y2)
```

MMD estimation: Python

Alternative initialization-1:

```
>>> co2 = ite.cost.BDMMD_UStat_IChol(eta=1e-2)
>>> mmd2 = co2.estimation(y1, y2)
```

MMD estimation: Python

Alternative initialization-1:

```
>>> co2 = ite.cost.BDMMD_UStat_IChol(eta=1e-2)
>>> mmd2 = co2.estimation(y1, y2)
```

Alternative-2:

```
>>> k = Kernel('name': 'RBF', 'sigma': 1)
>>> co3 = ite.cost.BDMMD_UStat_IChol(kernel=k, eta=1e-2)
>>> mmd3 = co3.estimation(y1, y2)
```

Towards unbiased estimators

- MMD, HSIC: $\mathbb{E}_{x,x'} k(x, x')$ -type quantities.

Towards unbiased estimators

- MMD, HSIC: $\mathbb{E}_{x,x'} k(x, x')$ -type quantities.
- x, x' : independence.

Towards unbiased estimators

- MMD, HSIC: $\mathbb{E}_{x,x'} k(x, x')$ -type quantities.
- x, x' : independence.
- Plugin methods: $i = j$, biased.

Towards unbiased estimators

- MMD, HSIC: $\mathbb{E}_{x,x'} k(x, x')$ -type quantities.
- x, x' : independence.
- Plugin methods: $i = j$, biased.
- If we restrict to $i \neq j$, we got **unbiased** estimators.

Towards unbiased estimators

- MMD, HSIC: $\mathbb{E}_{x,x'} k(x, x')$ -type quantities.
- x, x' : independence.
- Plugin methods: $i = j$, biased.
- If we restrict to $i \neq j$, we got **unbiased** estimators.

Question

What is happening here? Concentration of the estimators?

→ hypothesis testing: our statistics := these estimators

Unbiased estimators for $\mathbb{E}_{x,x'} k(x, x')$ -type quantities – extensions of **average**

Task

- Goal: estimate

$$\theta(\mathbb{P}) := \mathbb{E}_{\mathbb{P}} h(X_1, \dots, X_m).$$

Task

- Goal: estimate

$$\theta(\mathbb{P}) := \mathbb{E}_{\mathbb{P}} h(X_1, \dots, X_m).$$

- Given: $x_1, \dots, x_n \stackrel{i.i.d.}{\sim} \mathbb{P}$, $n \geq m$.

Task

- Goal: estimate

$$\theta(\mathbb{P}) := \mathbb{E}_{\mathbb{P}} h(X_1, \dots, X_m).$$

- Given: $x_1, \dots, x_n \stackrel{i.i.d.}{\sim} \mathbb{P}$, $n \geq m$.
- Assume (w.l.o.g.): h is **symmetric**,

$$h(x_1, \dots, x_m) = h(x_{\pi(1)}, \dots, x_{\pi(m)}) \quad \forall \pi \text{ permutation.}$$

Example: $k(x, x') = k(x', x)$.

Task

- Goal: estimate

$$\theta(\mathbb{P}) := \mathbb{E}_{\mathbb{P}} h(X_1, \dots, X_m).$$

- Given: $x_1, \dots, x_n \stackrel{i.i.d.}{\sim} \mathbb{P}$, $n \geq m$.
- Assume (w.l.o.g.): h is **symmetric**,

$$h(x_1, \dots, x_m) = h(x_{\pi(1)}, \dots, x_{\pi(m)}) \quad \forall \pi \text{ permutation.}$$

Example: $k(x, x') = k(x', x)$.

- Otherwise: $h \leftarrow \frac{1}{m!} \sum_{\pi} h(x_{\pi(1)}, \dots, x_{\pi(m)})$.

- Estimator for $\mathbb{E}_{\mathbb{P}} h(X_1, \dots, X_m)$:

$$U_n = U(x_1, \dots, x_n) = \frac{1}{\binom{n}{m}} \sum_c h(x_{i_1}, \dots, x_{i_m}),$$

\sum_c : m -tuples without replacement.

- Estimator for $\mathbb{E}_{\mathbb{P}} h(X_1, \dots, X_m)$:

$$U_n = U(x_1, \dots, x_n) = \frac{1}{\binom{n}{m}} \sum_c h(x_{i_1}, \dots, x_{i_m}),$$

\sum_c : m -tuples without replacement.

- U_n : unbiased, i.e. $\mathbb{E}_{\mathbb{P}}(U_n) = \theta$.

- Estimator for $\mathbb{E}_{\mathbb{P}} h(X_1, \dots, X_m)$:

$$V_n = V(x_1, \dots, x_n) = \frac{1}{n^m} \sum_{i_1=1}^n \dots \sum_{i_m=1}^n h(x_{i_1}, \dots, x_{i_m}).$$

- Estimator for $\mathbb{E}_{\mathbb{P}} h(X_1, \dots, X_m)$:

$$V_n = V(x_1, \dots, x_n) = \frac{1}{n^m} \sum_{i_1=1}^n \dots \sum_{i_m=1}^n h(x_{i_1}, \dots, x_{i_m}).$$

- Samples with replacement.

U-statistic: examples

- $\theta(\mathbb{P}) = \mathbb{E}_{X \sim \mathbb{P}} X$. Sample average:

$$h(x) = x, \quad U(x_1, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n x_i.$$

U-statistic: examples

- $\theta(\mathbb{P}) = \mathbb{E}_{X \sim \mathbb{P}} X$. Sample average:

$$h(x) = x, \quad U(x_1, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n x_i.$$

- $\theta(\mathbb{P}) = \mathbb{E}_{X \sim \mathbb{P}} X^k$. Sample k^{th} moment:

$$h(x) = x^k, \quad U(x_1, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n x_i^k.$$

U-statistic: examples

- $\theta(\mathbb{P}) = \mathbb{E}_{X \sim \mathbb{P}} X$. Sample average:

$$h(x) = x, \quad U(x_1, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n x_i.$$

- $\theta(\mathbb{P}) = \mathbb{E}_{X \sim \mathbb{P}} X^k$. Sample k^{th} moment:

$$h(x) = x^k, \quad U(x_1, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n x_i^k.$$

- $\theta(\mathbb{P}) = \sigma^2(\mathbb{P}) = \int (x - \mu)^2 d\mathbb{P}(x)$, $\mu = \mathbb{E}_{X \sim \mathbb{P}} X$. Sample variance:

$$\sigma^2(\mathbb{P}) = \mathbb{E} X^2 - \mathbb{E}^2 X$$

U-statistic: examples

- $\theta(\mathbb{P}) = \mathbb{E}_{X \sim \mathbb{P}} X$. Sample average:

$$h(x) = x, \quad U(x_1, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n x_i.$$

- $\theta(\mathbb{P}) = \mathbb{E}_{X \sim \mathbb{P}} X^k$. Sample k^{th} moment:

$$h(x) = x^k, \quad U(x_1, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n x_i^k.$$

- $\theta(\mathbb{P}) = \sigma^2(\mathbb{P}) = \int (x - \mu)^2 d\mathbb{P}(x)$, $\mu = \mathbb{E}_{X \sim \mathbb{P}} X$. Sample variance:

$$\sigma^2(\mathbb{P}) = \mathbb{E}X^2 - \mathbb{E}^2 X = \frac{\mathbb{E}X_1^2 + \mathbb{E}X_2^2}{2} - \mathbb{E}X_1 \mathbb{E}X_2$$

U-statistic: examples

- $\theta(\mathbb{P}) = \mathbb{E}_{X \sim \mathbb{P}} X$. Sample average:

$$h(x) = x, \quad U(x_1, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n x_i.$$

- $\theta(\mathbb{P}) = \mathbb{E}_{X \sim \mathbb{P}} X^k$. Sample k^{th} moment:

$$h(x) = x^k, \quad U(x_1, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n x_i^k.$$

- $\theta(\mathbb{P}) = \sigma^2(\mathbb{P}) = \int (x - \mu)^2 d\mathbb{P}(x)$, $\mu = \mathbb{E}_{X \sim \mathbb{P}} X$. Sample variance:

$$\sigma^2(\mathbb{P}) = \mathbb{E}X^2 - \mathbb{E}^2 X = \frac{\mathbb{E}X_1^2 + \mathbb{E}X_2^2}{2} - \mathbb{E}X_1 \mathbb{E}X_2 = \mathbb{E}h(X_1, X_2),$$

$$h(x_1, x_2) = \frac{x_1^2 + x_2^2 - 2x_1 x_2}{2}$$

U-statistic: examples

- $\theta(\mathbb{P}) = \mathbb{E}_{X \sim \mathbb{P}} X$. Sample average:

$$h(x) = x, \quad U(x_1, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n x_i.$$

- $\theta(\mathbb{P}) = \mathbb{E}_{X \sim \mathbb{P}} X^k$. Sample k^{th} moment:

$$h(x) = x^k, \quad U(x_1, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n x_i^k.$$

- $\theta(\mathbb{P}) = \sigma^2(\mathbb{P}) = \int (x - \mu)^2 d\mathbb{P}(x)$, $\mu = \mathbb{E}_{X \sim \mathbb{P}} X$. Sample variance:

$$\sigma^2(\mathbb{P}) = \mathbb{E}X^2 - \mathbb{E}^2 X = \frac{\mathbb{E}X_1^2 + \mathbb{E}X_2^2}{2} - \mathbb{E}X_1 \mathbb{E}X_2 = \mathbb{E}h(X_1, X_2),$$

$$h(x_1, x_2) = \frac{x_1^2 + x_2^2 - 2x_1 x_2}{2} = \frac{(x_1 - x_2)^2}{2},$$

U-statistic: examples

- $\theta(\mathbb{P}) = \mathbb{E}_{X \sim \mathbb{P}} X$. Sample average:

$$h(x) = x, \quad U(x_1, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n x_i.$$

- $\theta(\mathbb{P}) = \mathbb{E}_{X \sim \mathbb{P}} X^k$. Sample k^{th} moment:

$$h(x) = x^k, \quad U(x_1, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n x_i^k.$$

- $\theta(\mathbb{P}) = \sigma^2(\mathbb{P}) = \int (x - \mu)^2 d\mathbb{P}(x)$, $\mu = \mathbb{E}_{X \sim \mathbb{P}} X$. Sample variance:

$$\sigma^2(\mathbb{P}) = \mathbb{E}X^2 - \mathbb{E}^2 X = \frac{\mathbb{E}X_1^2 + \mathbb{E}X_2^2}{2} - \mathbb{E}X_1 \mathbb{E}X_2 = \mathbb{E}h(X_1, X_2),$$

$$h(x_1, x_2) = \frac{x_1^2 + x_2^2 - 2x_1 x_2}{2} = \frac{(x_1 - x_2)^2}{2},$$

$$U(x_1, \dots, x_n) = \frac{2}{n(n-1)} \sum_{1 \leq i < j \leq n} h(x_i, x_j)$$

U-statistic: examples

- $\theta(\mathbb{P}) = \mathbb{E}_{X \sim \mathbb{P}} X$. Sample average:

$$h(x) = x, \quad U(x_1, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n x_i.$$

- $\theta(\mathbb{P}) = \mathbb{E}_{X \sim \mathbb{P}} X^k$. Sample k^{th} moment:

$$h(x) = x^k, \quad U(x_1, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n x_i^k.$$

- $\theta(\mathbb{P}) = \sigma^2(\mathbb{P}) = \int (x - \mu)^2 d\mathbb{P}(x)$, $\mu = \mathbb{E}_{X \sim \mathbb{P}} X$. Sample variance:

$$\sigma^2(\mathbb{P}) = \mathbb{E}X^2 - \mathbb{E}^2 X = \frac{\mathbb{E}X_1^2 + \mathbb{E}X_2^2}{2} - \mathbb{E}X_1 \mathbb{E}X_2 = \mathbb{E}h(X_1, X_2),$$

$$h(x_1, x_2) = \frac{x_1^2 + x_2^2 - 2x_1 x_2}{2} = \frac{(x_1 - x_2)^2}{2},$$

$$U(x_1, \dots, x_n) = \frac{2}{n(n-1)} \sum_{1 \leq i < j \leq n} h(x_i, x_j) = \frac{1}{n(n-1)} \sum_{i \neq j} h(x_i, x_j)$$

U-statistic: examples

- $\theta(\mathbb{P}) = \mathbb{E}_{X \sim \mathbb{P}} X$. Sample average:

$$h(x) = x, \quad U(x_1, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n x_i.$$

- $\theta(\mathbb{P}) = \mathbb{E}_{X \sim \mathbb{P}} X^k$. Sample k^{th} moment:

$$h(x) = x^k, \quad U(x_1, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n x_i^k.$$

- $\theta(\mathbb{P}) = \sigma^2(\mathbb{P}) = \int (x - \mu)^2 d\mathbb{P}(x)$, $\mu = \mathbb{E}_{X \sim \mathbb{P}} X$. Sample variance:

$$\sigma^2(\mathbb{P}) = \mathbb{E}X^2 - \mathbb{E}^2 X = \frac{\mathbb{E}X_1^2 + \mathbb{E}X_2^2}{2} - \mathbb{E}X_1 \mathbb{E}X_2 = \mathbb{E}h(X_1, X_2),$$

$$h(x_1, x_2) = \frac{x_1^2 + x_2^2 - 2x_1 x_2}{2} = \frac{(x_1 - x_2)^2}{2},$$

$$U(x_1, \dots, x_n) = \frac{2}{n(n-1)} \sum_{1 \leq i < j \leq n} h(x_i, x_j) = \frac{1}{n(n-1)} \sum_{i \neq j} h(x_i, x_j) = s_n^2.$$

U-statistic: examples+

$$\theta(\mathbb{P}) = F_{\mathbb{P}}(t_0) = \mathbb{P}(X \leq t_0)$$

U-statistic: examples+

$$\theta(\mathbb{P}) = F_{\mathbb{P}}(t_0) = \mathbb{P}(X \leq t_0) = \mathbb{E}_{\mathbb{P}} \chi_{(-\infty, t_0]},$$

U-statistic: examples+

$$\theta(\mathbb{P}) = F_{\mathbb{P}}(t_0) = \mathbb{P}(X \leq t_0) = \mathbb{E}_{\mathbb{P}} \chi_{(-\infty, t_0]},$$

$$h(x) = \chi_{\{x \leq t_0\}},$$

$$U(x_1, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n \chi_{\{x_i \leq t_0\}}$$

U-statistic: examples+

$$\theta(\mathbb{P}) = F_{\mathbb{P}}(t_0) = \mathbb{P}(X \leq t_0) = \mathbb{E}_{\mathbb{P}} \chi_{(-\infty, t_0]},$$

$$h(x) = \chi_{\{x \leq t_0\}},$$

$$U(x_1, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n \chi_{\{x_i \leq t_0\}} = F_n(t_0),$$

F_n : empirical cdf.

Extension: if we have L independent samples \rightarrow MMD:
 $L = 2$

- Given: $x_1^{(j)}, \dots, x_{n_j}^{(j)} \stackrel{i.i.d.}{\sim} \mathbb{P}_j$ ($j = 1, \dots, L$), $n_j \geq m_j$.

Extension: if we have L independent samples \rightarrow MMD:
 $L = 2$

- Given: $x_1^{(j)}, \dots, x_{n_j}^{(j)} \stackrel{i.i.d.}{\sim} \mathbb{P}_j$ ($j = 1, \dots, L$), $n_j \geq m_j$.
- Goal: estimate: $\theta = \mathbb{E}h\left(\underbrace{X_1^{(1)}, \dots, X_{m_1}^{(1)}}_{1^{st} \text{ block}}, \dots, \underbrace{X_1^{(L)}, \dots, X_{m_L}^{(L)}}_{L^{th} \text{ block}}\right)$.

Extension: if we have L independent samples \rightarrow MMD:
 $L = 2$

- Given: $x_1^{(j)}, \dots, x_{n_j}^{(j)} \stackrel{i.i.d.}{\sim} \mathbb{P}_j$ ($j = 1, \dots, L$), $n_j \geq m_j$.
- Goal: estimate: $\theta = \mathbb{E}h\left(\underbrace{X_1^{(1)}, \dots, X_{m_1}^{(1)}}_{1^{st} \text{ block}}, \dots, \underbrace{X_1^{(L)}, \dots, X_{m_L}^{(L)}}_{L^{th} \text{ block}}\right)$.
- Assumption: symmetry for each block.

Extension: if we have L independent samples \rightarrow MMD:
 $L = 2$

- Given: $x_1^{(j)}, \dots, x_{n_j}^{(j)} \stackrel{i.i.d.}{\sim} \mathbb{P}_j$ ($j = 1, \dots, L$), $n_j \geq m_j$.
- Goal: estimate: $\theta = \mathbb{E} h\left(\underbrace{X_1^{(1)}, \dots, X_{m_1}^{(1)}}_{1^{st} \text{ block}}, \dots, \underbrace{X_1^{(L)}, \dots, X_{m_L}^{(L)}}_{L^{th} \text{ block}}\right)$.
- Assumption: symmetry for each block.
- L -sample U-statistic

$$U_n = \frac{1}{\prod_{j=1}^L \binom{n_j}{m_j}} \sum_c h\left(X_1^{(1)}, \dots, X_{m_1}^{(1)}, \dots, X_1^{(L)}, \dots, X_{m_L}^{(L)}\right).$$

Results [Serfling, 1980]

U_n : minimum variance unbiased estimator.

Results [Serfling, 1980]

U_n : minimum variance unbiased estimator. Notation:

- Asymptotics: depends on $\text{var} \neq 0$ condition.

$$h_c(x_1, \dots, x_c) := \mathbb{E} h(x_1, \dots, x_c, X_{c+1}, \dots, X_m),$$

Results [Serfling, 1980]

U_n : minimum variance unbiased estimator. Notation:

- Asymptotics: depends on $\text{var} \neq 0$ condition.

$$h_c(x_1, \dots, x_c) := \mathbb{E} h(x_1, \dots, x_c, X_{c+1}, \dots, X_m),$$

$$v_c := \text{var } h_c(X_1, \dots, X_c), v_0 = 0.$$

Results [Serfling, 1980]

U_n : minimum variance unbiased estimator. Notation:

- Asymptotics: depends on $\text{var} \neq 0$ condition.

$$h_c(x_1, \dots, x_c) := \mathbb{E} h(x_1, \dots, x_c, X_{c+1}, \dots, X_m),$$
$$v_c := \text{var } h_c(X_1, \dots, X_c), v_0 = 0.$$

- If $\mathbb{E} h^2(X_1, \dots, X_m) < \infty$:

$$0 = v_0 \leq v_1 \leq \dots \leq v_m = \text{var } h(X_1, \dots, X_m) < \infty.$$

Results [Serfling, 1980]

U_n : minimum variance unbiased estimator. Notation:

- Asymptotics: depends on $\text{var} \neq 0$ condition.

$$h_c(x_1, \dots, x_c) := \mathbb{E} h(x_1, \dots, x_c, X_{c+1}, \dots, X_m),$$
$$v_c := \text{var } h_c(X_1, \dots, X_c), v_0 = 0.$$

- If $\mathbb{E} h^2(X_1, \dots, X_m) < \infty$:

$$0 = v_0 \leq v_1 \leq \dots \leq v_m = \text{var } h(X_1, \dots, X_m) < \infty.$$

- c : $0 = v_1 = \dots = v_{c-1} < v_c$. $c = 1$: non-degenerate, $c \geq 2$: degenerate U-statistic.

Results [Serfling, 1980]

U_n : minimum variance unbiased estimator. Notation:

- Asymptotics: depends on $\text{var} \neq 0$ condition.

$$h_c(x_1, \dots, x_c) := \mathbb{E} h(x_1, \dots, x_c, X_{c+1}, \dots, X_m),$$
$$v_c := \text{var } h_c(X_1, \dots, X_c), v_0 = 0.$$

- If $\mathbb{E} h^2(X_1, \dots, X_m) < \infty$:

$$0 = v_0 \leq v_1 \leq \dots \leq v_m = \text{var } h(X_1, \dots, X_m) < \infty.$$

- c : $0 = v_1 = \dots = v_{c-1} < v_c$. $c = 1$: non-degenerate, $c \geq 2$: degenerate U-statistic.

In most applications

$c = 1$ or $c = 2$.

Asymptotics for $c = 1$

Assume: $\mathbb{E}_{\mathbb{P}} h^2 < \infty$, $c = 1$.

$$n^{\frac{1}{2}}(U_n - \theta) \xrightarrow{d} N(0, m^2 v_1),$$

i.e.

$$U_n \text{ is AN} \left(\theta, \frac{m^2 v_1}{n} \right),$$

AN := asymptotically normal.

Asymptotics for $c = 2$

Assume: $\mathbb{E}_{\mathbb{P}} h^2 < \infty$, $c = 2$.

$$n(U_n - \theta) \xrightarrow{d} \frac{m(m-1)}{2} Y, \quad Y = \sum_{j=1}^{\infty} \lambda_j (\chi_j^2 - 1),$$

where

Asymptotics for $c = 2$

Assume: $\mathbb{E}_{\mathbb{P}} h^2 < \infty$, $c = 2$.

$$n(U_n - \theta) \xrightarrow{d} \frac{m(m-1)}{2} Y, \quad Y = \sum_{j=1}^{\infty} \lambda_j (\chi_j^2 - 1),$$

where

- χ_j^2 : i.i.d. $N^2(0, 1)$ variables,

Asymptotics for $c = 2$

Assume: $\mathbb{E}_{\mathbb{P}} h^2 < \infty$, $c = 2$.

$$n(U_n - \theta) \xrightarrow{d} \frac{m(m-1)}{2} Y, \quad Y = \sum_{j=1}^{\infty} \lambda_j (\chi_j^2 - 1),$$

where

- χ_j^2 : i.i.d. $N^2(0, 1)$ variables,
- λ_j : \mathbb{R} -eigenvalues of $T = T(\tilde{h}_2)$, $\tilde{h}_2 = h_2 - \theta$

$$(Tg)(x) = \int \tilde{h}_2(x, y) g(y) d\mathbb{P}(y), \quad g \in L^2.$$

Theorem (Hoeffding inequality)

Let $h(x_1, \dots, x_m) \in [a, b]$. If $\sigma^2 = \text{var } h$, then for any $t > 0$

$$\mathbb{P}(U_n - \theta \geq t) \leq e^{-\frac{2[n/m]t^2}{(b-a)^2}}.$$

- Minimum variance unbiased estimator.
- $c = 1$: asymptotically normal.
- $c = 2$: asymptotically ∞ -sum of weighted χ^2 .
- For bounded h : Hoeffding inequality.

Application

Hypothesis testing!

Hypothesis testing

What is a two-sample test?

- Given:
 - $X = \{x_i\}_{i=1}^m \stackrel{i.i.d.}{\sim} \mathbb{P}$, $Y = \{y_j\}_{j=1}^n \stackrel{i.i.d.}{\sim} \mathbb{Q}$.
 - Example: $x_i = i^{th}$ happy face, $y_j = j^{th}$ sad face.

What is a two-sample test?

- Given:
 - $X = \{x_i\}_{i=1}^m \stackrel{i.i.d.}{\sim} \mathbb{P}$, $Y = \{y_j\}_{j=1}^n \stackrel{i.i.d.}{\sim} \mathbb{Q}$.
 - Example: $x_i = i^{th}$ happy face, $y_j = j^{th}$ sad face.
- Problem: using X , Y test

$$H_0 : \mathbb{P} = \mathbb{Q}, \text{ vs}$$
$$H_1 : \mathbb{P} \neq \mathbb{Q}.$$

What is a two-sample test?

- Given:
 - $X = \{x_i\}_{i=1}^m \stackrel{i.i.d.}{\sim} \mathbb{P}$, $Y = \{y_j\}_{j=1}^n \stackrel{i.i.d.}{\sim} \mathbb{Q}$.
 - Example: $x_i = i^{th}$ happy face, $y_j = j^{th}$ sad face.
- Problem: using X , Y test

$$\begin{aligned}H_0 : \mathbb{P} &= \mathbb{Q}, \text{ vs} \\H_1 : \mathbb{P} &\neq \mathbb{Q}.\end{aligned}$$

- Assumption: $x, y \in \mathcal{X}$. (\mathcal{X}, k) : kernel-endowed domain.

What is a two-sample test?

- Given:
 - $X = \{x_i\}_{i=1}^m \stackrel{i.i.d.}{\sim} \mathbb{P}$, $Y = \{y_j\}_{j=1}^n \stackrel{i.i.d.}{\sim} \mathbb{Q}$.
 - Example: $x_i = i^{th}$ happy face, $y_j = j^{th}$ sad face.
- Problem: using X , Y test

$$H_0 : \mathbb{P} = \mathbb{Q}, \text{ vs}$$
$$H_1 : \mathbb{P} \neq \mathbb{Q}.$$

- Assumption: $x, y \in \mathcal{X}$. (\mathcal{X}, k) : kernel-endowed domain.

Discrepancy measure

Example: MMD

What is an independence test?

- Given: **paired samples**
 - $Z = \{(x_i, y_i)\}_{i=1}^n \stackrel{i.i.d.}{\sim} \mathbb{P}_{xy}$.
 - Example:
 - x_i : i^{th} text in English, y_i : i^{th} text translated to French.

What is an independence test?

- Given: paired samples
 - $Z = \{(x_i, y_i)\}_{i=1}^n \stackrel{i.i.d.}{\sim} \mathbb{P}_{xy}$.
 - Example:
 - x_i : i^{th} text in English, y_i : i^{th} text translated to French.
- Problem: using data Z test

$$H_0 : \mathbb{P}_{xy} = \mathbb{P}_x \otimes \mathbb{P}_y, \quad H_1 : \mathbb{P}_{xy} \neq \mathbb{P}_x \otimes \mathbb{P}_y.$$

What is an independence test?

- Given: paired samples
 - $Z = \{(x_i, y_i)\}_{i=1}^n \stackrel{i.i.d.}{\sim} \mathbb{P}_{xy}$.
 - Example:
 - x_i : i^{th} text in English, y_i : i^{th} text translated to French.
- Problem: using data Z test

$$H_0 : \mathbb{P}_{xy} = \mathbb{P}_x \otimes \mathbb{P}_y, \quad H_1 : \mathbb{P}_{xy} \neq \mathbb{P}_x \otimes \mathbb{P}_y.$$

- Assumption: $(x, y) \in \mathcal{X} \times \mathcal{Y}$. (\mathcal{X}, k) , (\mathcal{Y}, ℓ) : with kernels.

What is an independence test?

- Given: paired samples
 - $Z = \{(x_i, y_i)\}_{i=1}^n \stackrel{i.i.d.}{\sim} \mathbb{P}_{xy}$.
 - Example:
 - x_i : i^{th} text in English, y_i : i^{th} text translated to French.
- Problem: using data Z test

$$H_0 : \mathbb{P}_{xy} = \mathbb{P}_x \otimes \mathbb{P}_y, \quad H_1 : \mathbb{P}_{xy} \neq \mathbb{P}_x \otimes \mathbb{P}_y.$$

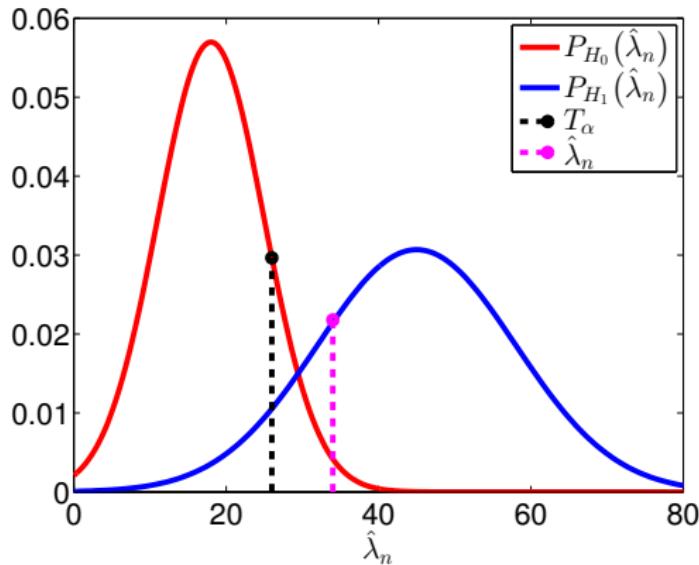
- Assumption: $(x, y) \in \mathcal{X} \times \mathcal{Y}$. (\mathcal{X}, k) , (\mathcal{Y}, ℓ) : with kernels.

Discrepancy measure

Example: HSIC

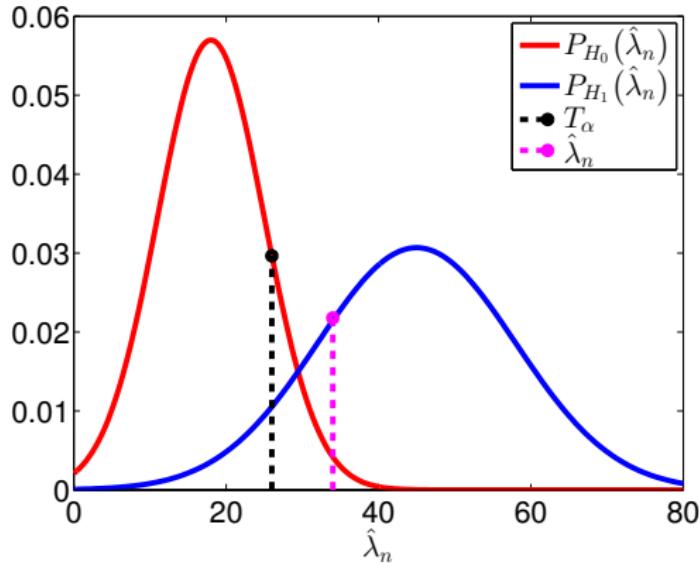
Concepts in hypothesis testing

- Test statistic: $\hat{\lambda}_n = \hat{\lambda}_n(X, Y)$, random.
- Significance level: $\alpha = 0.01$.
- Under H_0 : $P_{H_0}(\underbrace{\hat{\lambda}_n \leq T_\alpha}_{\text{correctly accepting } H_0}) = 1 - \alpha$.



Concepts in hypothesis testing

- Test statistic: $\hat{\lambda}_n = \hat{\lambda}_n(X, Y)$, random.
- Significance level: $\alpha = 0.01$.
- Under H_0 : $P_{H_0}(\underbrace{\hat{\lambda}_n \leq T_\alpha}_{\text{correctly accepting } H_0}) = 1 - \alpha$.
- Under H_1 : $P_{H_1}(T_\alpha < \hat{\lambda}_n) = P(\text{correctly rejecting } H_0) =: \text{power}$.



Two-sample testing (aka homogeneity testing) – details.

Two-sample testing with MMD

[Gretton et al., 2007, Gretton et al., 2012]

- Statistic: $\hat{\lambda}_n = \widehat{\text{MMD}}_b^2$ or $\widehat{\text{MMD}}_u^2$.

Two-sample testing with MMD

[Gretton et al., 2007, Gretton et al., 2012]

- Statistic: $\hat{\lambda}_n = \widehat{\text{MMD}}_b^2$ or $\widehat{\text{MMD}}_u^2$.
- Reject H_0 : if $\hat{\lambda}_n$ is 'large'.

Two-sample testing with MMD

[Gretton et al., 2007, Gretton et al., 2012]

- Statistic: $\hat{\lambda}_n = \widehat{\text{MMD}}_b^2$ or $\widehat{\text{MMD}}_u^2$.
- Reject H_0 : if $\hat{\lambda}_n$ is 'large'.
- We need to control $\hat{\lambda}_n$.

Two-sample testing with MMD

[Gretton et al., 2007, Gretton et al., 2012]

- Statistic: $\hat{\lambda}_n = \widehat{\text{MMD}}_b^2$ or $\widehat{\text{MMD}}_u^2$.
- Reject H_0 : if $\hat{\lambda}_n$ is 'large'.
- We need to control $\hat{\lambda}_n$.
- We will use U-statistic theory.

Finite-sample control

- Large deviation inequalities.
- $P\left(\left|\widehat{\text{MMD}}(\mathbb{P}, \mathbb{Q}) - \text{MMD}(\mathbb{P}, \mathbb{Q})\right| \geq \epsilon\right) \leq f(\epsilon, m, n) \xrightarrow{m, n \rightarrow \infty} 0.$

Finite-sample control

- Large deviation inequalities.
- $P\left(\left|\widehat{\text{MMD}}(\mathbb{P}, \mathbb{Q}) - \text{MMD}(\mathbb{P}, \mathbb{Q})\right| \geq \epsilon\right) \leq f(\epsilon, m, n) \xrightarrow{m, n \rightarrow \infty} 0.$
- \Rightarrow tests: **consistent** against fixed alternative.

- Distribution-free tests
 - quantile: \mathbb{P}, \mathbb{Q} -independent!
 - less sensitive to detecting differences.

- Distribution-free tests
 - quantile: \mathbb{P}, \mathbb{Q} -independent!
 - less sensitive to detecting differences.
- Proof idea:

- Distribution-free tests
 - quantile: \mathbb{P}, \mathbb{Q} -independent!
 - less sensitive to detecting differences.
- Proof idea:
 - $\widehat{\text{MMD}}_b^2$: bounded difference property, McDiarmid inequality.

- Distribution-free tests
 - quantile: \mathbb{P}, \mathbb{Q} -independent!
 - less sensitive to detecting differences.
- Proof idea:
 - $\widehat{\text{MMD}}_b^2$: bounded difference property, McDiarmid inequality.
 - $\widehat{\text{MMD}}_u^2$: large deviation bound of U-statistics.

Asymptotics based test

Needed: Asymptotic distribution of $\widehat{\text{MMD}}_u^2$.

$$\begin{aligned}\widehat{\text{MMD}}_u^2(\mathbb{P}, \mathbb{Q}) &= \frac{1}{m(m-1)} \sum_{i=1}^m \sum_{j \neq i}^m k(x_i, x_j) + \frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j \neq i}^n k(y_i, y_j) \\ &\quad - \frac{2}{mn} \sum_{i=1}^m \sum_{j=1}^n k(x_i, y_j).\end{aligned}$$

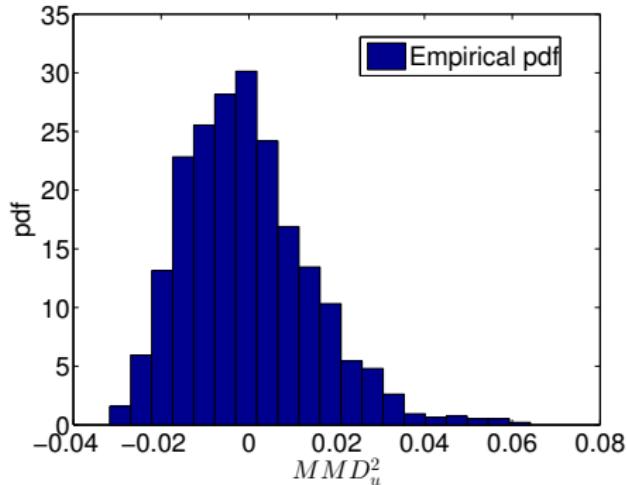
Two-sample test using MMD asymptotics: H_0 [$c = 2!$]

Under H_0 ($\mathbb{P} = \mathbb{Q}$): asymptotic distribution is

$$\widehat{n\text{MMD}_u^2}(\mathbb{P}, \mathbb{P}) \sim \sum_{i=1}^{\infty} \lambda_i(z_i^2 - 2),$$

where $z_i \sim N(0, 2)$ i.i.d.,

$$\int_{\mathcal{X}} \tilde{k}(x, x') v_i(x) d\mathbb{P}(x) = \lambda_i v_i(x'), \quad \tilde{k}(x, x') = \langle \varphi_x - \mu_{\mathbb{P}}, \varphi_{x'} - \mu_{\mathbb{P}} \rangle_{\mathcal{H}_k}.$$



Approximate the null by

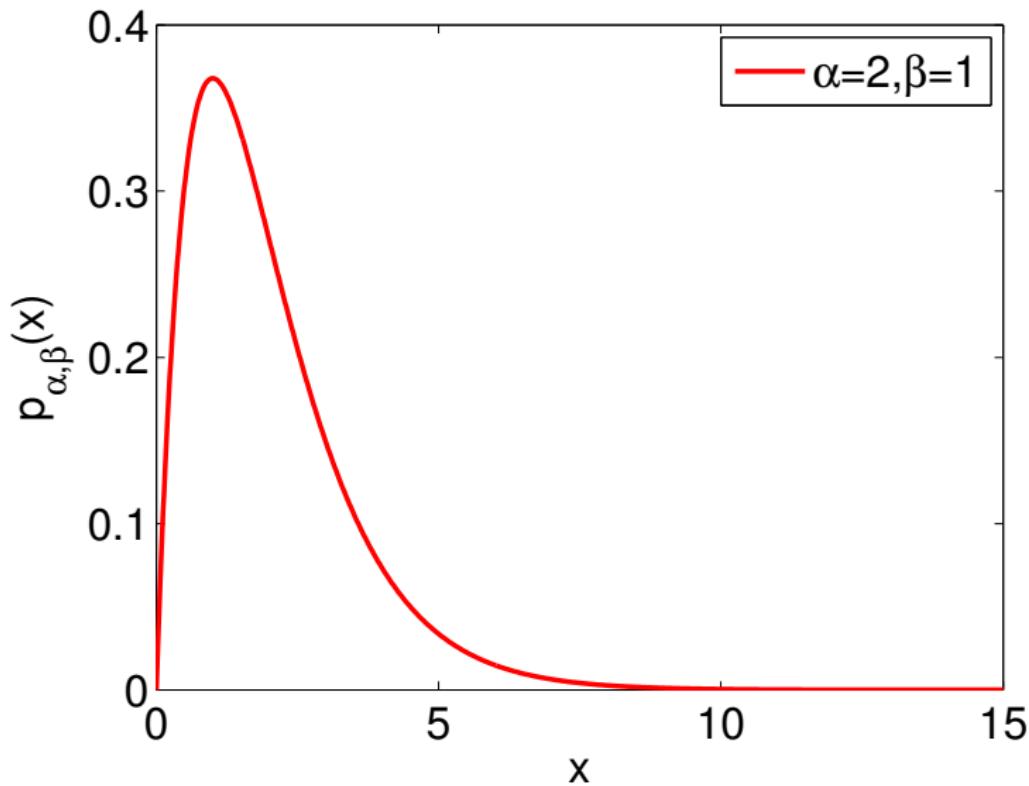
- **permutation-test**: slow.

Approximate the null by

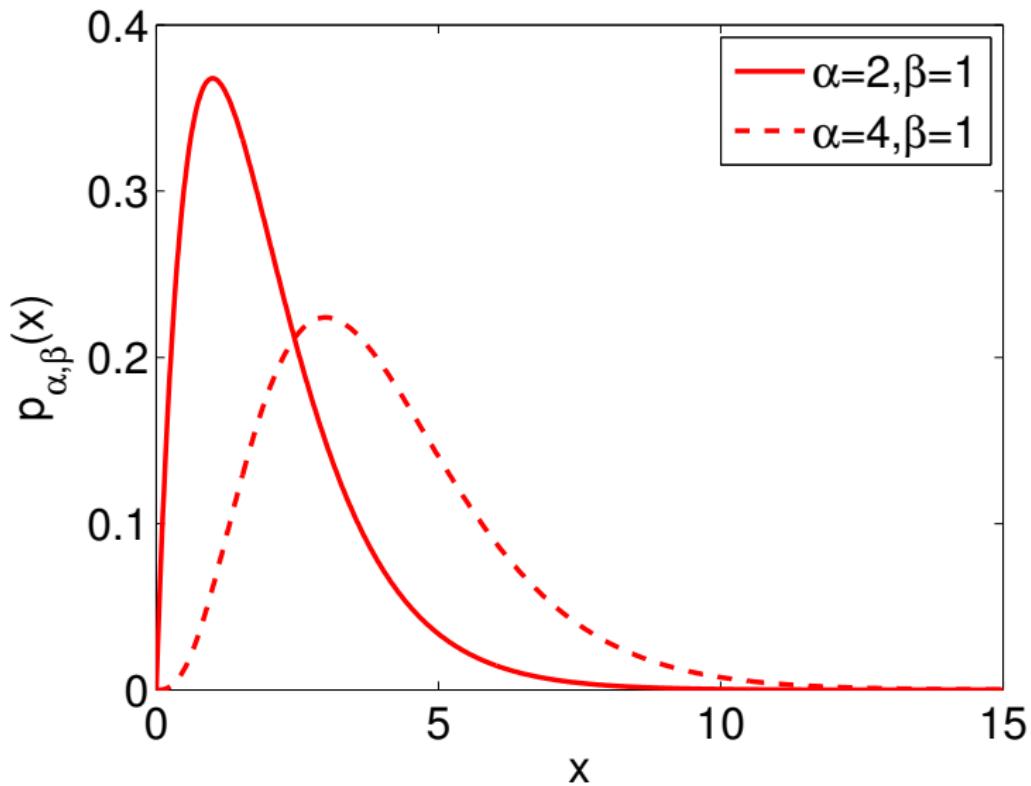
- **permutation-test**: slow.
- two-parameter **gamma distribution** [Johnson et al., 1994]:

$$p_{\alpha,\beta}(x) = \frac{x^{\alpha-1} e^{-\frac{x}{\beta}}}{\beta^\alpha \Gamma(\alpha)} \quad (x > 0, \alpha: \text{shape} > 0, \beta: \text{scale} > 0).$$

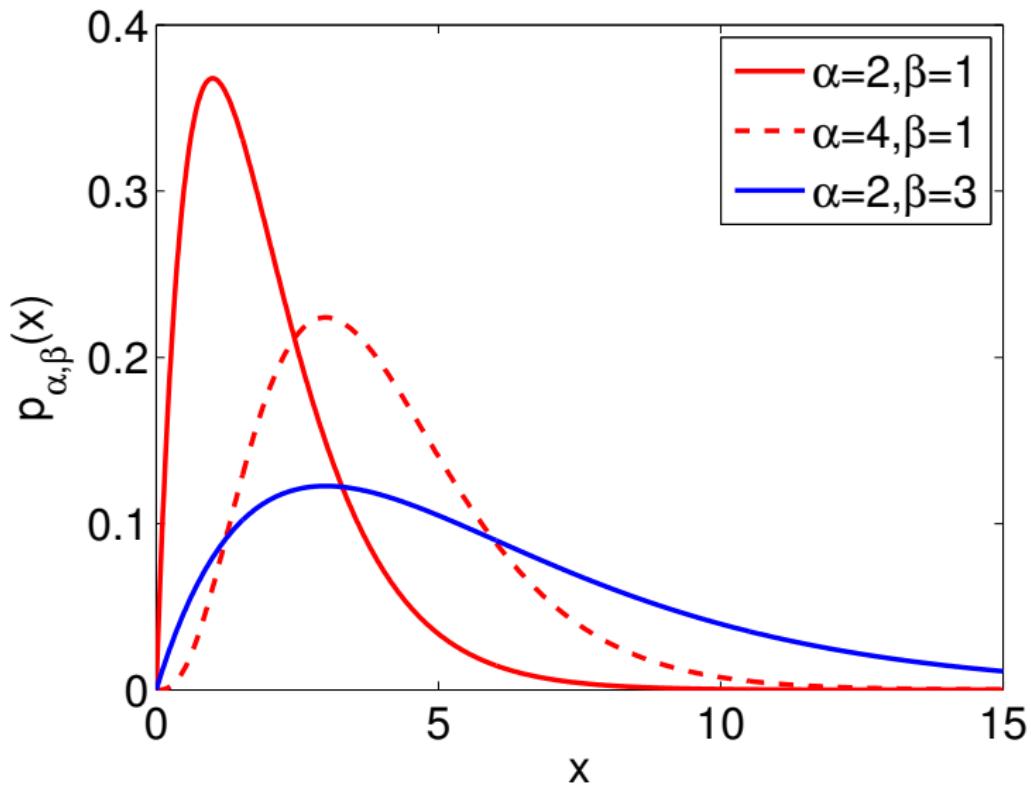
Gamma distribution: demo



Gamma distribution: demo



Gamma distribution: demo



Gamma approximation

- Assumption: statistic $T \sim p_{\alpha, \beta}$.

Gamma approximation

- Assumption: statistic $T \sim p_{\alpha, \beta}$.
- For $p_{\alpha, \beta}$ gamma distribution:

$$\mathbb{E} T = \alpha\beta, \quad \text{var}(T) = \alpha\beta^2$$

Gamma approximation

- Assumption: statistic $T \sim p_{\alpha,\beta}$.
- For $p_{\alpha,\beta}$ gamma distribution:

$$\mathbb{E} T = \alpha\beta, \quad \text{var}(T) = \alpha\beta^2 \quad \Rightarrow \quad \beta = \frac{\text{var}(T)}{\mathbb{E} T} \quad \Rightarrow \quad \alpha = \frac{\mathbb{E}^2 T}{\text{var}(T)}.$$

Gamma approximation

- Assumption: statistic $T \sim p_{\alpha,\beta}$.
- For $p_{\alpha,\beta}$ gamma distribution:

$$\mathbb{E}T = \alpha\beta, \quad \text{var}(T) = \alpha\beta^2 \quad \Rightarrow \quad \beta = \frac{\text{var}(T)}{\mathbb{E}T} \quad \Rightarrow \quad \alpha = \frac{\mathbb{E}^2 T}{\text{var}(T)}.$$

- Thus, $\widehat{\mathbb{E}T}$ and $\widehat{\text{var}(T)}$ $\rightarrow \hat{\alpha}, \hat{\beta}$.

Gamma approximation

- Assumption: statistic $T \sim p_{\alpha, \beta}$.
- For $p_{\alpha, \beta}$ gamma distribution:

$$\mathbb{E} T = \alpha\beta, \quad \text{var}(T) = \alpha\beta^2 \quad \Rightarrow \quad \beta = \frac{\text{var}(T)}{\mathbb{E} T} \quad \Rightarrow \quad \alpha = \frac{\mathbb{E}^2 T}{\text{var}(T)}.$$

- Thus, $\widehat{\mathbb{E} T}$ and $\widehat{\text{var}(T)}$ $\rightarrow \hat{\alpha}, \hat{\beta}$.
- **Consistency** of the test is **lost**.

Which null approximation to use?

Rules-of-thumb:

- Small sample size: permutation test.

Which null approximation to use?

Rules-of-thumb:

- Small sample size: permutation test.
- Medium sample size: gamma approximation, truncated expansion [Gretton et al., 2009],

Which null approximation to use?

Rules-of-thumb:

- Small sample size: permutation test.
- Medium sample size: gamma approximation, truncated expansion [Gretton et al., 2009],
- Large sample size:
 - online techniques [Gretton et al., 2012], or
 - recent linear methods (next time).

Independence testing: HSIC

Independence testing

Theorem ([Gretton et al., 2008, Pfister et al., 2017])

Under H_0

$$n \widehat{\text{HSIC}}_b^2 \xrightarrow{d} \sum_{i=1}^{\infty} \lambda_i z_i^2, \quad z_i \sim N(0, 1).$$

Independence testing

Theorem ([Gretton et al., 2008, Pfister et al., 2017])

Under H_0

$$n \widehat{\text{HSIC}}_b^2 \xrightarrow{d} \sum_{i=1}^{\infty} \lambda_i z_i^2, \quad z_i \sim N(0, 1).$$

Notes:

- For U-statistic: $\sum_i \lambda_i (z_i^2 - 1)$.

Independence testing

Theorem ([Gretton et al., 2008, Pfister et al., 2017])

Under H_0

$$n \widehat{\text{HSIC}}_b^2 \xrightarrow{d} \sum_{i=1}^{\infty} \lambda_i z_i^2, \quad z_i \sim N(0, 1).$$

Notes:

- For U-statistic: $\sum_i \lambda_i (z_i^2 - 1)$.
- In practice: permutation-test/gamma-approximation.

Related work

Two-sample problem: truncated expansion

[Gretton et al., 2009]: $n = m$, $z_i = (x_i, y_i)$. Estimator:

$$\widehat{\text{MMD}}_{u'}^2(\mathbb{P}, \mathbb{Q}) = \frac{1}{n(n-1)} \sum_{i \neq j} h(z_i, z_j),$$

$$h(z, z') = k(x, x') + k(y, y') - k(x, y') - k(x', y).$$

Two-sample problem: truncated expansion

[Gretton et al., 2009]: $n = m$, $z_i = (x_i, y_i)$. Estimator:

$$\widehat{\text{MMD}}_{u'}^2(\mathbb{P}, \mathbb{Q}) = \frac{1}{n(n-1)} \sum_{i \neq j} h(z_i, z_j),$$

$$h(z, z') = k(x, x') + k(y, y') - k(x, y') - k(x', y).$$

$\widehat{\text{MMD}}_{u'}^2$: unbiased.

Theorem

Assuming $\sum_{i=1}^{\infty} \lambda_i^{\frac{1}{2}} < \infty$, the empirical null converges as $n \rightarrow \infty$

$$T_n := \sum_{i=1}^n \hat{\lambda}_{i,n} (a_i^2 - 2) \xrightarrow{d} \sum_{i=1}^{\infty} \lambda_i (a_i^2 - 2), \quad a_i \sim N(0, 2).$$

Theorem

Assuming $\sum_{i=1}^{\infty} \lambda_i^{\frac{1}{2}} < \infty$, the empirical null converges as $n \rightarrow \infty$

$$T_n := \sum_{i=1}^n \hat{\lambda}_{i,n} (a_i^2 - 2) \xrightarrow{d} \sum_{i=1}^{\infty} \lambda_i (a_i^2 - 2), \quad a_i \sim N(0, 2).$$

Note:

$$\hat{\lambda}_{i,n} := \frac{\lambda_i(\tilde{\mathbf{G}}_x)}{n} \quad (i = 1, \dots, n), \quad \tilde{\mathbf{G}}_x \in \mathbb{R}^{n \times n}.$$

Online variant [Gretton et al., 2012]

$$\widehat{\text{MMD}}_{u'}^2(\mathbb{P}, \mathbb{Q}) = \frac{1}{n(n-1)} \sum_{i \neq j} h(z_i, z_j),$$

has a natural online approximation, $n_2 := \lceil n/2 \rceil$

$$\widehat{\text{MMD}}_I^2(\mathbb{P}, \mathbb{Q}) = \frac{1}{n_2} \sum_{i=1}^{n_2} h((x_{2i-1}, y_{2i-1}), (x_{2i}, y_{2i})).$$

Online variant [Gretton et al., 2012]

$$\widehat{\text{MMD}}_{u'}^2(\mathbb{P}, \mathbb{Q}) = \frac{1}{n(n-1)} \sum_{i \neq j} h(z_i, z_j),$$

has a natural online approximation, $n_2 := [n/2]$

$$\widehat{\text{MMD}}_I^2(\mathbb{P}, \mathbb{Q}) = \frac{1}{n_2} \sum_{i=1}^{n_2} h((x_{2i-1}, y_{2i-1}), (x_{2i}, y_{2i})).$$

- Unbiased.

Online variant [Gretton et al., 2012]

$$\widehat{\text{MMD}}_{u'}^2(\mathbb{P}, \mathbb{Q}) = \frac{1}{n(n-1)} \sum_{i \neq j} h(z_i, z_j),$$

has a natural online approximation, $n_2 := \lceil n/2 \rceil$

$$\widehat{\text{MMD}}_I^2(\mathbb{P}, \mathbb{Q}) = \frac{1}{n_2} \sum_{i=1}^{n_2} h((x_{2i-1}, y_{2i-1}), (x_{2i}, y_{2i})).$$

- Unbiased.
- Linear-time: streaming data.

Online variant [Gretton et al., 2012]

$$\widehat{\text{MMD}}_{u'}^2(\mathbb{P}, \mathbb{Q}) = \frac{1}{n(n-1)} \sum_{i \neq j} h(z_i, z_j),$$

has a natural online approximation, $n_2 := [n/2]$

$$\widehat{\text{MMD}}_I^2(\mathbb{P}, \mathbb{Q}) = \frac{1}{n_2} \sum_{i=1}^{n_2} h((x_{2i-1}, y_{2i-1}), (x_{2i}, y_{2i})).$$

- Unbiased.
- Linear-time: streaming data.
- In practice: **high** variance.

By the **average** the CLT kicks in:

Theorem

Assuming $\mathbb{E} h^2 \in (0, \infty)$, $\widehat{\text{MMD}}_I^2$ is asymptotically normal

$$\sqrt{n} \left[\widehat{\text{MMD}}_I^2(\mathbb{P}, \mathbb{Q}) - \text{MMD}^2(\mathbb{P}, \mathbb{Q}) \right] \xrightarrow{d} \mathcal{N}(0, \sigma^2),$$

where $\sigma^2 = 2 \left[\mathbb{E}_{z,z'} h^2(z, z') - \mathbb{E}_{z,z'}^2 h(z, z') \right]$.

Idea:

- partition the data to blocks of size B ,
- on each block: compute $\widehat{\text{MMD}}_I^2$,
- average the results.

Properties:

- Statistic: asymptotically normal (H_0, H_1).
- For consistency: increase B_m s.t. $\frac{m}{B_m} \rightarrow \infty$.
- **Reduced variance.**

Three-variable interaction test

- Goal (interaction):

$$([x_1; x_2] \perp x_3) \vee ([x_1; x_3] \perp x_2) \vee ([x_2; x_3] \perp x_1).$$

Example: $\mathbb{P} = \mathbb{P}_{12} \otimes \mathbb{P}_3$.

Three-variable interaction test

- Goal (interaction):

$$([x_1; x_2] \perp x_3) \vee ([x_1; x_3] \perp x_2) \vee ([x_2; x_3] \perp x_1).$$

Example: $\mathbb{P} = \mathbb{P}_{12} \otimes \mathbb{P}_3$.

- Applications:

- structure learning of graphical models,
- discovering V-structures.

Analogy

Independence $\Leftrightarrow \mathbb{P} = \mathbb{P}_1 \otimes \mathbb{P}_2 \Leftrightarrow \mathbb{P} - \mathbb{P}_1 \otimes \mathbb{P}_2 = 0.$

Three-variable interaction test – continued

Analogy

$$\text{Independence} \Leftrightarrow \mathbb{P} = \mathbb{P}_1 \otimes \mathbb{P}_2 \Leftrightarrow \mathbb{P} - \mathbb{P}_1 \otimes \mathbb{P}_2 = 0.$$

- Lancaster 3-variable interaction [Lancaster, 1969]:

$$L(\mathbb{P}) = \mathbb{P} - \mathbb{P}_{1,2} \otimes \mathbb{P}_3 - \mathbb{P}_{2,3} \otimes \mathbb{P}_1 - \mathbb{P}_{1,3} \otimes \mathbb{P}_2 + 2\mathbb{P}_1 \otimes \mathbb{P}_2 \otimes \mathbb{P}_3.$$

is a signed measure,

$$\text{interaction} \Rightarrow L(\mathbb{P}) = 0.$$

Three-variable interaction test – continued

Analogy

$$\text{Independence} \Leftrightarrow \mathbb{P} = \mathbb{P}_1 \otimes \mathbb{P}_2 \Leftrightarrow \mathbb{P} - \mathbb{P}_1 \otimes \mathbb{P}_2 = 0.$$

- Lancaster 3-variable interaction [Lancaster, 1969]:

$$L(\mathbb{P}) = \mathbb{P} - \mathbb{P}_{1,2} \otimes \mathbb{P}_3 - \mathbb{P}_{2,3} \otimes \mathbb{P}_1 - \mathbb{P}_{1,3} \otimes \mathbb{P}_2 + 2\mathbb{P}_1 \otimes \mathbb{P}_2 \otimes \mathbb{P}_3.$$

is a signed measure,

$$\text{interaction} \Rightarrow L(\mathbb{P}) = 0.$$

- $x_i \in (\mathcal{X}_i, k_i)$ are kernel endowed domains.

Three-variable interaction test – continued

- Interaction index [Sejdinovic et al., 2013a]:

$$I = \left\| \mu_{L(\mathbb{P})} \right\|_{\mathcal{H}_{k_1} \otimes \mathcal{H}_{k_2} \otimes \mathcal{H}_{k_3}}^2.$$

Three-variable interaction test – continued

- Interaction index [Sejdinovic et al., 2013a]:

$$I = \left\| \mu_{L(\mathbb{P})} \right\|_{\mathcal{H}_{k_1} \otimes \mathcal{H}_{k_2} \otimes \mathcal{H}_{k_3}}^2.$$

- Empirical estimate:

$$\hat{I} = \frac{\left(\tilde{\mathbf{G}}_{x_1} \circ \tilde{\mathbf{G}}_{x_3} \circ \tilde{\mathbf{G}}_{x_3} \right)_{++}}{n^2}.$$

Three-variable interaction test – continued

- Interaction index [Sejdinovic et al., 2013a]:

$$I = \left\| \mu_{L(\mathbb{P})} \right\|_{\mathcal{H}_{k_1} \otimes \mathcal{H}_{k_2} \otimes \mathcal{H}_{k_3}}^2.$$

- Empirical estimate:

$$\hat{I} = \frac{\left(\tilde{\mathbf{G}}_{x_1} \circ \tilde{\mathbf{G}}_{x_3} \circ \tilde{\mathbf{G}}_{x_3} \right)_{++}}{n^2}.$$

- Null approximation: permutation-test.

Time-series tests: independence

- Goal: test independence of **stationary** processes.
- Independence tests:
 - Statistic: HSIC.

Time-series tests: independence

- Goal: test independence of **stationary** processes.
- Independence tests:
 - Statistic: HSIC.
 - i.i.d. **permutation** technique: would **fail**.

- Goal: test independence of **stationary** processes.
- Independence tests:
 - Statistic: HSIC.
 - i.i.d. **permutation** technique: would **fail**.
 - Idea: **shift**-approach = preserves 'time structure'
[Chwialkowski and Gretton, 2014].

- Permutation approach (i.i.d): ± 1 .

- Permutation approach (i.i.d): ± 1 .
- Idea: mask according to the memory of the processes.

- Permutation approach (i.i.d): ± 1 .
- Idea: mask according to the memory of the processes.
- Implementation [Chwialkowski et al., 2014]: based on wild bootstrap [Leucht and H.Neumann, 2013].

- Permutation approach (i.i.d): ± 1 .
- Idea: mask according to the memory of the processes.
- Implementation [Chwialkowski et al., 2014]: based on wild bootstrap [Leucht and H.Neumann, 2013].

3-variable interaction:

- Lancaster interaction + wild bootstrap [Rubenstein et al., 2016].

Goodness-of-fit test

- Given:
 - $\{x_i\}_{i=1}^n \stackrel{i.i.d.}{\sim} q$,
 - p : target distribution.

Goodness-of-fit test

- Given:
 - $\{x_i\}_{i=1}^n \stackrel{i.i.d.}{\sim} q$,
 - p : target distribution.
- p, q live on $\mathcal{X} \subset \mathbb{R}^d$ (differentiability), kernel k on \mathcal{X} .

Goodness-of-fit test

- Given:
 - $\{x_i\}_{i=1}^n \stackrel{i.i.d.}{\sim} q$,
 - p : target distribution.
- p, q live on $\mathcal{X} \subset \mathbb{R}^d$ (differentiability), kernel k on \mathcal{X} .
- Goal:

$$H_0 : p = q,$$

$$H_1 : p \neq q.$$

Goodness-of-fit test: continued

- Idea [Chwialkowski et al., 2016, Liu et al., 2016]: Stein operator

$$(\mathcal{S}_p f)(x) = \sum_{i=1}^d \left[\frac{\partial \log p(x)}{\partial x_i} f_i(x) + \frac{\partial f_i(x)}{\partial x_i} \right], \quad f \in \mathcal{H} := \otimes_{i=1}^d \mathcal{H}_k,$$

Goodness-of-fit test: continued

- Idea [Chwialkowski et al., 2016, Liu et al., 2016]: Stein operator

$$(\mathcal{S}_p f)(x) = \sum_{i=1}^d \left[\frac{\partial \log p(x)}{\partial x_i} f_i(x) + \frac{\partial f_i(x)}{\partial x_i} \right], \quad f \in \mathcal{H} := \otimes_{i=1}^d \mathcal{H}_k,$$
$$\mathbb{E}_{x \sim q} (\mathcal{S}_p f)(x) = 0 \iff p = q \ (\forall f).$$

Goodness-of-fit test: continued

- Idea [Chwialkowski et al., 2016, Liu et al., 2016]: Stein operator

$$(\mathcal{S}_p f)(x) = \sum_{i=1}^d \left[\frac{\partial \log p(x)}{\partial x_i} f_i(x) + \frac{\partial f_i(x)}{\partial x_i} \right], \quad f \in \mathcal{H} := \otimes_{i=1}^d \mathcal{H}_k,$$

$$T_p(q) = \sup_{\|f\|_{\mathcal{H}} \leq 1} \mathbb{E}_{x \sim q} (\mathcal{S}_p f)(x).$$

Goodness-of-fit test: continued

- Idea [Chwialkowski et al., 2016, Liu et al., 2016]: Stein operator

$$(\mathcal{S}_p f)(x) = \sum_{i=1}^d \left[\frac{\partial \log p(x)}{\partial x_i} f_i(x) + \frac{\partial f_i(x)}{\partial x_i} \right], \quad f \in \mathcal{H} := \otimes_{i=1}^d \mathcal{H}_k,$$

$$T_p(q) = \sup_{\|f\|_{\mathcal{H}} \leq 1} \mathbb{E}_{x \sim q} (\mathcal{S}_p f)(x).$$

- For c_0 -universal k : $T_p(q) = 0 \Leftrightarrow p = q$.

Goodness-of-fit test: continued

- Idea [Chwialkowski et al., 2016, Liu et al., 2016]: Stein operator

$$(\mathcal{S}_p f)(x) = \sum_{i=1}^d \left[\frac{\partial \log p(x)}{\partial x_i} f_i(x) + \frac{\partial f_i(x)}{\partial x_i} \right], \quad f \in \mathcal{H} := \otimes_{i=1}^d \mathcal{H}_k,$$

$$T_p(q) = \sup_{\|f\|_{\mathcal{H}} \leq 1} \mathbb{E}_{x \sim q} (\mathcal{S}_p f)(x).$$

- For c_0 -universal k : $T_p(q) = 0 \Leftrightarrow p = q$.
- Enough: p up to multiplicative constant $(\nabla \log p)$.

Goodness-of-fit test: continued

- Idea [Chwialkowski et al., 2016, Liu et al., 2016]: Stein operator

$$(\mathcal{S}_p f)(x) = \sum_{i=1}^d \left[\frac{\partial \log p(x)}{\partial x_i} f_i(x) + \frac{\partial f_i(x)}{\partial x_i} \right], \quad f \in \mathcal{H} := \otimes_{i=1}^d \mathcal{H}_k,$$

$$T_p(q) = \sup_{\|f\|_{\mathcal{H}} \leq 1} \mathbb{E}_{x \sim q} (\mathcal{S}_p f)(x).$$

- For c_0 -universal k : $T_p(q) = 0 \Leftrightarrow p = q$.
- Enough: p up to multiplicative constant ($\nabla \log p$).
- Null approximation: wild bootstrap (including non-i.i.d.).

Quadratic-time methods

- Two-sample, independence, interaction, goodness-of-fit test.

Quadratic-time methods

- Two-sample, independence, interaction, goodness-of-fit test.
- Kernel endowed domain (goodness-of-fit: \mathbb{R}^d).

Quadratic-time methods

- Two-sample, independence, interaction, goodness-of-fit test.
- Kernel endowed domain (goodness-of-fit: \mathbb{R}^d).
- Typically: null can be ugly, techniques **do not scale** well.

Quadratic-time methods

- Two-sample, independence, interaction, goodness-of-fit test.
- Kernel endowed domain (goodness-of-fit: \mathbb{R}^d).
- Typically: null can be ugly, techniques **do not scale** well.

Next step

Linear-time tests, with high-power!

Questions

- Lancaster-interaction measure: reason of the last term?
- Stein operator: why does it work?
- Stein operator: how to estimate it?

Interaction measure:

$$L(\mathbb{P}) = \mathbb{P} - \mathbb{P}_{1,2} \otimes \mathbb{P}_3 - \mathbb{P}_{2,3} \otimes \mathbb{P}_1 - \mathbb{P}_{1,3} \otimes \mathbb{P}_2 + 2\mathbb{P}_1 \otimes \mathbb{P}_2 \otimes \mathbb{P}_3.$$

Assume for example:

$$\begin{aligned}\mathbb{P} &= \mathbb{P}_1 \otimes \mathbb{P}_{2,3} \quad \Rightarrow \quad \mathbb{P}_{1,2} = \mathbb{P}_1 \otimes \mathbb{P}_2, \quad \mathbb{P}_{1,3} = \mathbb{P}_1 \otimes \mathbb{P}_3, \\ x_1 &\perp [x_2; x_3], \quad \quad \quad x_1 \perp x_2, \quad \quad \quad x_1 \perp x_3,\end{aligned}$$

Lancaster interaction

Interaction measure:

$$L(\mathbb{P}) = \mathbb{P} - \mathbb{P}_{1,2} \otimes \mathbb{P}_3 - \mathbb{P}_{2,3} \otimes \mathbb{P}_1 - \mathbb{P}_{1,3} \otimes \mathbb{P}_2 + 2\mathbb{P}_1 \otimes \mathbb{P}_2 \otimes \mathbb{P}_3.$$

Assume for example:

$$\mathbb{P} = \mathbb{P}_1 \otimes \mathbb{P}_{2,3} \quad \Rightarrow \quad \mathbb{P}_{1,2} = \mathbb{P}_1 \otimes \mathbb{P}_2, \quad \mathbb{P}_{1,3} = \mathbb{P}_1 \otimes \mathbb{P}_3,$$
$$x_1 \perp [x_2; x_3], \quad x_1 \perp x_2, \quad x_1 \perp x_3,$$

and L simplifies to

$$L(\mathbb{P}) = \mathbb{P} - \underbrace{\mathbb{P}_{1,2} \otimes \mathbb{P}_3}_{\mathbb{P}_1 \otimes \mathbb{P}_2 \otimes \mathbb{P}_3} - \underbrace{\mathbb{P}_{2,3} \otimes \mathbb{P}_1}_{\mathbb{P}_1 \otimes \mathbb{P}_2 \otimes \mathbb{P}_3} - \underbrace{\mathbb{P}_{1,3} \otimes \mathbb{P}_2}_{\mathbb{P}_1 \otimes \mathbb{P}_2 \otimes \mathbb{P}_3} + 2\mathbb{P}_1 \otimes \mathbb{P}_2 \otimes \mathbb{P}_3 = 0.$$

Stein operator ($d = 1$ for simplicity): why?

Let $f \in \mathcal{H}_k$.

$$(S_p f)(x) = [\log p(x)]' f(x) + f'(x)$$

Stein operator ($d = 1$ for simplicity): why?

Let $f \in \mathcal{H}_k$.

$$\begin{aligned}(S_p f)(x) &= [\log p(x)]' f(x) + f'(x) \\ &= \frac{p'(x)}{p(x)} f(x) + \frac{f'(x)}{p(x)} p(x)\end{aligned}$$

Stein operator ($d = 1$ for simplicity): why?

Let $f \in \mathcal{H}_k$.

$$\begin{aligned}(S_p f)(x) &= [\log p(x)]' f(x) + f'(x) \\&= \frac{p'(x)}{p(x)} f(x) + \frac{f'(x)}{p(x)} p(x) = \frac{[p(x)f(x)]'}{p(x)}.\end{aligned}$$

Stein operator ($d = 1$ for simplicity): why?

Let $f \in \mathcal{H}_k$.

$$\begin{aligned}(S_p f)(x) &= [\log p(x)]' f(x) + f'(x) \\&= \frac{p'(x)}{p(x)} f(x) + \frac{f'(x)}{p(x)} p(x) = \frac{[p(x)f(x)]'}{p(x)}.\end{aligned}$$

$p = q$ implies: for any f

$$\mathbb{E}_{x \sim q} (S_p f)(x) = \int_{\mathbb{R}} \frac{[p(x)f(x)]'}{p(x)} p(x) dx$$

Stein operator ($d = 1$ for simplicity): why?

Let $f \in \mathcal{H}_k$.

$$\begin{aligned}(S_p f)(x) &= [\log p(x)]' f(x) + f'(x) \\&= \frac{p'(x)}{p(x)} f(x) + \frac{f'(x)}{p(x)} p(x) = \frac{[p(x)f(x)]'}{p(x)}.\end{aligned}$$

$p = q$ implies: for any f

$$\begin{aligned}\mathbb{E}_{x \sim q} (S_p f)(x) &= \int_{\mathbb{R}} \frac{[p(x)f(x)]'}{p(x)} p(x) dx \\&= \int_{\mathbb{R}} [p(x)f(x)]' dx\end{aligned}$$

Stein operator ($d = 1$ for simplicity): why?

Let $f \in \mathcal{H}_k$.

$$\begin{aligned}(S_p f)(x) &= [\log p(x)]' f(x) + f'(x) \\&= \frac{p'(x)}{p(x)} f(x) + \frac{f'(x)}{p(x)} p(x) = \frac{[p(x)f(x)]'}{p(x)}.\end{aligned}$$

$p = q$ implies: for any f

$$\begin{aligned}\mathbb{E}_{x \sim q} (S_p f)(x) &= \int_{\mathbb{R}} \frac{[p(x)f(x)]'}{p(x)} p(x) dx \\&= \int_{\mathbb{R}} [p(x)f(x)]' dx = [p(x)f(x)]_{x=-\infty}^{x=\infty} = 0.\end{aligned}$$

Assumption: $\lim_{|x| \rightarrow \infty} p(x)f(x) = 0$.

Stein operator: computation

Test statistics:

$$T_p(q) = \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \mathbb{E}_{x \sim q}(S_p f)(x).$$

We rewrite $(S_p f)(x)$ by the reproducing property:

$$(S_p f)(x) = [\log p(x)]' \color{red}{f(x)} + \color{blue}{f'(x)}$$

Stein operator: computation

Test statistics:

$$T_p(q) = \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \mathbb{E}_{x \sim q}(S_p f)(x).$$

We rewrite $(S_p f)(x)$ by the reproducing property:

$$\begin{aligned} (S_p f)(x) &= [\log p(x)]' \color{red}{f(x)} + \color{blue}{f'(x)} \\ &= [\log p(x)]' \langle \color{red}{f}, k(\cdot, x) \rangle_{\mathcal{H}_k} + \langle \color{blue}{f}, k'(\cdot, x) \rangle_{\mathcal{H}_k} \end{aligned}$$

Stein operator: computation

Test statistics:

$$T_p(q) = \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \mathbb{E}_{x \sim q}(S_p f)(x).$$

We rewrite $(S_p f)(x)$ by the reproducing property:

$$\begin{aligned} (S_p f)(x) &= [\log p(x)]' \color{red}{f(x)} + \color{blue}{f'(x)} \\ &= [\log p(x)]' \langle \color{red}{f}, k(\cdot, x) \rangle_{\mathcal{H}_k} + \langle \color{blue}{f}, k'(\cdot, x) \rangle_{\mathcal{H}_k} \\ &= \underbrace{\langle f, [\log p(x)]' \color{red}{k(\cdot, x)} + \color{blue}{k'(\cdot, x)} \rangle_{\mathcal{H}_k}}_{=: \xi_p(\cdot, x)}. \end{aligned}$$

Stein operator: computation

Test statistics:

$$T_p(q) = \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \mathbb{E}_{x \sim q} (S_p f)(x).$$

We rewrite $(S_p f)(x)$ by the reproducing property:

$$\begin{aligned}(S_p f)(x) &= [\log p(x)]' \color{red}{f(x)} + \color{blue}{f'(x)} \\&= [\log p(x)]' \langle \color{red}{f}, k(\cdot, x) \rangle_{\mathcal{H}_k} + \langle \color{blue}{f}, k'(\cdot, x) \rangle_{\mathcal{H}_k} \\&= \langle f, \underbrace{[\log p(x)]' k(\cdot, x) + k'(\cdot, x)}_{=: \xi_p(\cdot, x)} \rangle_{\mathcal{H}_k}.\end{aligned}$$

Thus,

$$T_p(q) = \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \mathbb{E}_{x \sim q} \langle f, \xi_p(\cdot, x) \rangle_{\mathcal{H}_k}$$

Stein operator: computation

Test statistics:

$$T_p(q) = \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \mathbb{E}_{x \sim q}(S_p f)(x).$$

We rewrite $(S_p f)(x)$ by the reproducing property:

$$\begin{aligned} (S_p f)(x) &= [\log p(x)]' \color{red}{f(x)} + \color{blue}{f'(x)} \\ &= [\log p(x)]' \langle \color{red}{f}, k(\cdot, x) \rangle_{\mathcal{H}_k} + \langle \color{blue}{f}, k'(\cdot, x) \rangle_{\mathcal{H}_k} \\ &= \langle f, \underbrace{[\log p(x)]' k(\cdot, x) + k'(\cdot, x)}_{=: \xi_p(\cdot, x)} \rangle_{\mathcal{H}_k}. \end{aligned}$$

Thus,

$$T_p(q) = \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \mathbb{E}_{x \sim q} \langle f, \xi_p(\cdot, x) \rangle_{\mathcal{H}_k} = \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \langle f, \mathbb{E}_{x \sim q} \xi_p(\cdot, x) \rangle_{\mathcal{H}_k}$$

Stein operator: computation

Test statistics:

$$T_p(q) = \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \mathbb{E}_{x \sim q}(S_p f)(x).$$

We rewrite $(S_p f)(x)$ by the reproducing property:

$$\begin{aligned} (S_p f)(x) &= [\log p(x)]' \color{red}{f(x)} + \color{blue}{f'(x)} \\ &= [\log p(x)]' \langle \color{red}{f}, k(\cdot, x) \rangle_{\mathcal{H}_k} + \langle \color{blue}{f}, k'(\cdot, x) \rangle_{\mathcal{H}_k} \\ &= \underbrace{\langle f, [\log p(x)]' \color{red}{k(\cdot, x)} + \color{blue}{k'(\cdot, x)} \rangle_{\mathcal{H}_k}}_{=: \xi_p(\cdot, x)}. \end{aligned}$$

Thus,

$$\begin{aligned} T_p(q) &= \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \mathbb{E}_{x \sim q} \langle f, \xi_p(\cdot, x) \rangle_{\mathcal{H}_k} = \sup_{\|f\|_{\mathcal{H}_k} \leq 1} \langle f, \mathbb{E}_{x \sim q} \xi_p(\cdot, x) \rangle_{\mathcal{H}_k} \\ &= \|\mathbb{E}_{x \sim q} \xi_p(\cdot, x)\|_{\mathcal{H}_k} =: \|g\|_{\mathcal{H}_k}, g : \text{Stein witness function}. \end{aligned}$$

Stein operator: computation finished

Until now: with $\mathbf{g} = \mathbb{E}_{x \sim q} \xi_p(\cdot, x)$, $\xi_p(\cdot, x) = [\log p(x)]' k(\cdot, x) + k'(\cdot, x)$

$$[T_p(q)]^2 = \|\mathbf{g}\|_{\mathcal{H}_k}^2 = \langle \mathbb{E}_{x \sim q} \xi_p(\cdot, x), \mathbb{E}_{x' \sim q} \xi_p(\cdot, x') \rangle_{\mathcal{H}_k}$$

Stein operator: computation finished

Until now: with $\mathbf{g} = \mathbb{E}_{x \sim q} \xi_p(\cdot, x)$, $\xi_p(\cdot, x) = [\log p(x)]' k(\cdot, x) + k'(\cdot, x)$

$$\begin{aligned}[T_p(q)]^2 &= \|\mathbf{g}\|_{\mathcal{H}_k}^2 = \langle \mathbb{E}_{x \sim q} \xi_p(\cdot, x), \mathbb{E}_{x' \sim q} \xi_p(\cdot, x') \rangle_{\mathcal{H}_k} \\ &= \mathbb{E}_{x \sim q} \mathbb{E}_{x' \sim q} \langle \xi_p(\cdot, x), \xi_p(\cdot, x') \rangle_{\mathcal{H}_k},\end{aligned}$$

Stein operator: computation finished

Until now: with $\mathbf{g} = \mathbb{E}_{x \sim q} \xi_p(\cdot, x)$, $\xi_p(\cdot, x) = [\log p(x)]' k(\cdot, x) + k'(\cdot, x)$

$$\begin{aligned}[T_p(q)]^2 &= \|\mathbf{g}\|_{\mathcal{H}_k}^2 = \langle \mathbb{E}_{x \sim q} \xi_p(\cdot, x), \mathbb{E}_{x' \sim q} \xi_p(\cdot, x') \rangle_{\mathcal{H}_k} \\ &= \mathbb{E}_{x \sim q} \mathbb{E}_{x' \sim q} \langle \xi_p(\cdot, x), \xi_p(\cdot, x') \rangle_{\mathcal{H}_k},\end{aligned}$$

$$\begin{aligned}\langle \xi_p(\cdot, x), \xi_p(\cdot, x') \rangle_{\mathcal{H}_k} &= \langle [\log p(x)]' k(\cdot, x) + k'(\cdot, x), \\ &\quad [\log p(x')]' k(\cdot, x') + k'(\cdot, x') \rangle_{\mathcal{H}_k} =: h_p(x, x'),\end{aligned}$$

Stein operator: computation finished

Until now: with $\mathbf{g} = \mathbb{E}_{x \sim q} \xi_p(\cdot, x)$, $\xi_p(\cdot, x) = [\log p(x)]' k(\cdot, x) + k'(\cdot, x)$

$$\begin{aligned}[T_p(q)]^2 &= \|\mathbf{g}\|_{\mathcal{H}_k}^2 = \langle \mathbb{E}_{x \sim q} \xi_p(\cdot, x), \mathbb{E}_{x' \sim q} \xi_p(\cdot, x') \rangle_{\mathcal{H}_k} \\ &= \mathbb{E}_{x \sim q} \mathbb{E}_{x' \sim q} \langle \xi_p(\cdot, x), \xi_p(\cdot, x') \rangle_{\mathcal{H}_k},\end{aligned}$$

$$\begin{aligned}\langle \xi_p(\cdot, x), \xi_p(\cdot, x') \rangle_{\mathcal{H}_k} &= \langle [\log p(x)]' k(\cdot, x) + k'(\cdot, x), \\ &\quad [\log p(x')]' k(\cdot, x') + k'(\cdot, x') \rangle_{\mathcal{H}_k} =: h_p(x, x'), \\ h_p(x, y) &= [\log p(x)]' [\log p(y)]' k(x, y) + [\log p(x)]' k'_y(x, y) + \\ &\quad [\log p(y)]' k'_x(x, y) + k''_{xy}(x, y).\end{aligned}$$

Stein operator: computation finished

Until now: with $\mathbf{g} = \mathbb{E}_{x \sim q} \xi_p(\cdot, x)$, $\xi_p(\cdot, x) = [\log p(x)]' k(\cdot, x) + k'(\cdot, x)$

$$\begin{aligned}[T_p(q)]^2 &= \|\mathbf{g}\|_{\mathcal{H}_k}^2 = \langle \mathbb{E}_{x \sim q} \xi_p(\cdot, x), \mathbb{E}_{x' \sim q} \xi_p(\cdot, x') \rangle_{\mathcal{H}_k} \\ &= \mathbb{E}_{x \sim q} \mathbb{E}_{x' \sim q} \langle \xi_p(\cdot, x), \xi_p(\cdot, x') \rangle_{\mathcal{H}_k},\end{aligned}$$

$$\begin{aligned}\langle \xi_p(\cdot, x), \xi_p(\cdot, x') \rangle_{\mathcal{H}_k} &= \langle [\log p(x)]' k(\cdot, x) + k'(\cdot, x), \\ &\quad [\log p(x')]' k(\cdot, x') + k'(\cdot, x') \rangle_{\mathcal{H}_k} =: h_p(x, x'), \\ h_p(x, y) &= [\log p(x)]' [\log p(y)]' k(x, y) + [\log p(x)]' k'_y(x, y) + \\ &\quad [\log p(y)]' k'_x(x, y) + k''_{xy}(x, y).\end{aligned}$$

⇒ Quadratic-time estimator (U-statistic):

$$\widehat{[T_p(q)]^2} = \frac{1}{n(n-1)} \sum_{i \neq j} h_p(x_i, x_j).$$

Hypothesis testing: linear-time methods

- Nyström method, random Fourier features.
- **Analytic representations** → linear-time two-sample testing.
- **High-power** linear-time techniques:
 - two-sample testing,
 - independence testing.
 - goodness-of-fit testing.

Three schemes

Exemplified in independence testing [Zhang et al., 2017]:

- **block-HSIC**: analog of block-MMD.

Three schemes

Exemplified in independence testing [Zhang et al., 2017]:

- **block-HSIC**: analog of block-MMD.
- 2 low-rank schemes:

Three schemes

Exemplified in independence testing [Zhang et al., 2017]:

- **block-HSIC**: analog of block-MMD.
- 2 low-rank schemes:
 - **Nyström method**
[Williams and Seeger, 2001, Drineas and Mahoney, 2005].

Three schemes

Exemplified in independence testing [Zhang et al., 2017]:

- **block-HSIC**: analog of block-MMD.
- 2 low-rank schemes:
 - **Nyström method**
[Williams and Seeger, 2001, Drineas and Mahoney, 2005].
 - **random Fourier features**: [Rahimi and Recht, 2007,
Sutherland and Schneider, 2015, Sriperumbudur and Szabó, 2015].

$$\begin{aligned}\mathcal{C}_{xy}^c &= \mathbb{E}_{xy} \left[(\varphi(x) - \mu_x) \otimes (\psi(y) - \mu_y) \right] \\ &= \mathbb{E}_{xy} [\varphi(x) \otimes \psi(y)] - \mu_x \otimes \mu_y, \\ \text{HSIC}(x, y) &= \|\mathcal{C}_{xy}^c\|_{HS}.\end{aligned}$$

Nyström method

Idea

Approximate $\mathbf{G} \in \mathbb{R}^{n \times n}$ with a (random) subset of size $r \ll n$.

Nyström method

Idea

Approximate $\mathbf{G} \in \mathbb{R}^{n \times n}$ with a (random) subset of size $r \ll n$.

- Without centering:

$$\mathbb{R}^{n \times n} \ni \hat{\mathbf{G}} \approx \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-1} \mathbf{G}_{r,n}$$

Nyström method

Idea

Approximate $\mathbf{G} \in \mathbb{R}^{n \times n}$ with a (random) subset of size $r \ll n$.

- Without centering:

$$\mathbb{R}^{n \times n} \ni \hat{\mathbf{G}} \approx \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-1} \mathbf{G}_{r,n} = \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} \mathbf{G}_{r,r}^{-\frac{1}{2}} \mathbf{G}_{n,r}^T$$

Nyström method

Idea

Approximate $\mathbf{G} \in \mathbb{R}^{n \times n}$ with a (random) subset of size $r \ll n$.

- Without centering:

$$\begin{aligned}\mathbb{R}^{n \times n} &\ni \hat{\mathbf{G}} \approx \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-1} \mathbf{G}_{r,n} = \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} \mathbf{G}_{r,r}^{-\frac{1}{2}} \mathbf{G}_{n,r}^T \\ &= \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} \left[\mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} \right]^T\end{aligned}$$

Nyström method

Idea

Approximate $\mathbf{G} \in \mathbb{R}^{n \times n}$ with a (random) subset of size $r \ll n$.

- Without centering:

$$\begin{aligned}\mathbb{R}^{n \times n} &\ni \hat{\mathbf{G}} \approx \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-1} \mathbf{G}_{r,n} = \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} \mathbf{G}_{r,r}^{-\frac{1}{2}} \mathbf{G}_{n,r}^T \\ &= \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} \left[\mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} \right]^T = \Phi^u (\Phi^u)^T, \quad \Phi^u \in \mathbb{R}^{n \times r},\end{aligned}$$

Nyström method

Idea

Approximate $\mathbf{G} \in \mathbb{R}^{n \times n}$ with a (random) subset of size $r \ll n$.

- Without centering:

$$\begin{aligned}\mathbb{R}^{n \times n} &\ni \hat{\mathbf{G}} \approx \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-1} \mathbf{G}_{r,n} = \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} \mathbf{G}_{r,r}^{-\frac{1}{2}} \mathbf{G}_{n,r}^T \\ &= \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} \left[\mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} \right]^T = \Phi^u (\Phi^u)^T, \quad \Phi^u \in \mathbb{R}^{n \times r}, \\ \mathbb{R}^{r \times r} &\ni \mathbf{C}^u = \left[\mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} \right]^T \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}}\end{aligned}$$

Nyström method

Idea

Approximate $\mathbf{G} \in \mathbb{R}^{n \times n}$ with a (random) subset of size $r \ll n$.

- Without centering:

$$\begin{aligned}\mathbb{R}^{n \times n} &\ni \hat{\mathbf{G}} \approx \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-1} \mathbf{G}_{r,n} = \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} \mathbf{G}_{r,r}^{-\frac{1}{2}} \mathbf{G}_{n,r}^T \\ &= \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} \left[\mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} \right]^T = \Phi^u (\Phi^u)^T, \quad \Phi^u \in \mathbb{R}^{n \times r}, \\ \mathbb{R}^{r \times r} &\ni \mathbf{C}^u = \left[\mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} \right]^T \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} =: (\Phi^u)^T \Phi^u.\end{aligned}$$

Nyström method

Idea

Approximate $\mathbf{G} \in \mathbb{R}^{n \times n}$ with a (random) subset of size $r \ll n$.

- Without centering:

$$\begin{aligned}\mathbb{R}^{n \times n} &\ni \hat{\mathbf{G}} \approx \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-1} \mathbf{G}_{r,n} = \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} \mathbf{G}_{r,r}^{-\frac{1}{2}} \mathbf{G}_{n,r}^T \\ &= \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} \left[\mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} \right]^T = \Phi^u (\Phi^u)^T, \quad \Phi^u \in \mathbb{R}^{n \times r}, \\ \mathbb{R}^{r \times r} &\ni \mathbf{C}^u = \left[\mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} \right]^T \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} =: (\Phi^u)^T \Phi^u.\end{aligned}$$

- With centering:

$$\mathbb{R}^{n \times n} \ni \tilde{\hat{\mathbf{G}}} = \mathbf{H}_n \hat{\mathbf{G}} \mathbf{H}_n$$

Nyström method

Idea

Approximate $\mathbf{G} \in \mathbb{R}^{n \times n}$ with a (random) subset of size $r \ll n$.

- Without centering:

$$\begin{aligned}\mathbb{R}^{n \times n} &\ni \hat{\mathbf{G}} \approx \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-1} \mathbf{G}_{r,n} = \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} \mathbf{G}_{r,r}^{-\frac{1}{2}} \mathbf{G}_{n,r}^T \\ &= \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} \left[\mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} \right]^T = \Phi^u (\Phi^u)^T, \quad \Phi^u \in \mathbb{R}^{n \times r}, \\ \mathbb{R}^{r \times r} &\ni \mathbf{C}^u = \left[\mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} \right]^T \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} =: (\Phi^u)^T \Phi^u.\end{aligned}$$

- With centering:

$$\mathbb{R}^{n \times n} \ni \tilde{\hat{\mathbf{G}}} = \mathbf{H}_n \hat{\mathbf{G}} \mathbf{H}_n = \mathbf{H}_n \Phi^u (\Phi^u)^T \mathbf{H}_n,$$

Nyström method

Idea

Approximate $\mathbf{G} \in \mathbb{R}^{n \times n}$ with a (random) subset of size $r \ll n$.

- Without centering:

$$\begin{aligned}\mathbb{R}^{n \times n} &\ni \hat{\mathbf{G}} \approx \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-1} \mathbf{G}_{r,n} = \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} \mathbf{G}_{r,r}^{-\frac{1}{2}} \mathbf{G}_{n,r}^T \\ &= \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} \left[\mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} \right]^T = \Phi^u (\Phi^u)^T, \quad \Phi^u \in \mathbb{R}^{n \times r}, \\ \mathbb{R}^{r \times r} &\ni \mathbf{C}^u = \left[\mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} \right]^T \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} =: (\Phi^u)^T \Phi^u.\end{aligned}$$

- With centering:

$$\begin{aligned}\mathbb{R}^{n \times n} &\ni \tilde{\hat{\mathbf{G}}} = \mathbf{H}_n \hat{\mathbf{G}} \mathbf{H}_n = \mathbf{H}_n \Phi^u (\Phi^u)^T \mathbf{H}_n, \\ \mathbb{R}^{r \times r} &\ni \mathbf{C}^c = (\Phi^u)^T \mathbf{H}_n \mathbf{H}_n \Phi^u\end{aligned}$$

Nyström method

Idea

Approximate $\mathbf{G} \in \mathbb{R}^{n \times n}$ with a (random) subset of size $r \ll n$.

- Without centering:

$$\begin{aligned}\mathbb{R}^{n \times n} &\ni \hat{\mathbf{G}} \approx \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-1} \mathbf{G}_{r,n} = \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} \mathbf{G}_{r,r}^{-\frac{1}{2}} \mathbf{G}_{n,r}^T \\ &= \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} \left[\mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} \right]^T = \Phi^u (\Phi^u)^T, \quad \Phi^u \in \mathbb{R}^{n \times r}, \\ \mathbb{R}^{r \times r} &\ni \mathbf{C}^u = \left[\mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} \right]^T \mathbf{G}_{n,r} \mathbf{G}_{r,r}^{-\frac{1}{2}} =: (\Phi^u)^T \Phi^u.\end{aligned}$$

- With centering:

$$\begin{aligned}\mathbb{R}^{n \times n} &\ni \tilde{\hat{\mathbf{G}}} = \mathbf{H}_n \hat{\mathbf{G}} \mathbf{H}_n = \mathbf{H}_n \Phi^u (\Phi^u)^T \mathbf{H}_n, \\ \mathbb{R}^{r \times r} &\ni \mathbf{C}^c = (\Phi^u)^T \mathbf{H}_n \mathbf{H}_n \Phi^u =: (\Phi^c)^T \Phi^c.\end{aligned}$$

Implementation for x and y , separately

On x :

$$\hat{\mathbf{G}}_x \approx \Phi_x^u (\Phi_x^u)^T$$

Implementation for x and y , separately

On x :

$$\hat{\mathbf{G}}_x \approx \Phi_x^u (\Phi_x^u)^T \Rightarrow \mathbf{C}_x^u = (\Phi_x^u)^T \Phi_x^u$$

Implementation for x and y , separately

On x :

$$\hat{\mathbf{G}}_x \approx \Phi_x^u (\Phi_x^u)^T \Rightarrow \mathbf{C}_x^u = (\Phi_x^u)^T \Phi_x^u, \quad \Phi_x^u = \left[(\Phi_{x,1}^u)^T; \dots; (\Phi_{x,n}^u)^T \right] \in \mathbb{R}^{n \times r_x},$$

Implementation for x and y , separately

On x :

$$\hat{\mathbf{G}}_x \approx \Phi_x^u (\Phi_x^u)^T \Rightarrow \mathbf{C}_x^u = (\Phi_x^u)^T \Phi_x^u, \quad \Phi_x^u = \left[(\Phi_{x,1}^u)^T; \dots; (\Phi_{x,n}^u)^T \right] \in \mathbb{R}^{n \times r_x},$$

$$\tilde{\hat{\mathbf{G}}}_x = \mathbf{H}_n \hat{\mathbf{G}}_x \mathbf{H}_n$$

Implementation for x and y , separately

On x :

$$\hat{\mathbf{G}}_x \approx \Phi_x^u (\Phi_x^u)^T \Rightarrow \mathbf{C}_x^u = (\Phi_x^u)^T \Phi_x^u, \quad \Phi_x^u = \left[(\Phi_{x,1}^u)^T; \dots; (\Phi_{x,n}^u)^T \right] \in \mathbb{R}^{n \times r_x},$$

$$\tilde{\hat{\mathbf{G}}}_x = \mathbf{H}_n \hat{\mathbf{G}}_x \mathbf{H}_n = \mathbf{H}_n \Phi_x^u (\Phi_x^u)^T \mathbf{H}_n$$

Implementation for x and y , separately

On x :

$$\hat{\mathbf{G}}_x \approx \Phi_x^u (\Phi_x^u)^T \Rightarrow \mathbf{C}_x^u = (\Phi_x^u)^T \Phi_x^u, \quad \Phi_x^u = \left[(\Phi_{x,1}^u)^T; \dots; (\Phi_{x,n}^u)^T \right] \in \mathbb{R}^{n \times r_x},$$

$$\tilde{\hat{\mathbf{G}}}_x = \mathbf{H}_n \hat{\mathbf{G}}_x \mathbf{H}_n = \mathbf{H}_n \Phi_x^u (\Phi_x^u)^T \mathbf{H}_n, \quad \mathbf{C}_x^c = (\Phi_x^u)^T \mathbf{H}_n \mathbf{H}_n \Phi_x^u$$

Implementation for x and y , separately

On x :

$$\hat{\mathbf{G}}_x \approx \Phi_x^u (\Phi_x^u)^T \Rightarrow \mathbf{C}_x^u = (\Phi_x^u)^T \Phi_x^u, \quad \Phi_x^u = \left[(\Phi_{x,1}^u)^T; \dots; (\Phi_{x,n}^u)^T \right] \in \mathbb{R}^{n \times r_x},$$

$$\tilde{\hat{\mathbf{G}}}_x = \mathbf{H}_n \hat{\mathbf{G}}_x \mathbf{H}_n = \mathbf{H}_n \Phi_x^u (\Phi_x^u)^T \mathbf{H}_n, \quad \mathbf{C}_x^c = (\Phi_x^u)^T \mathbf{H}_n \mathbf{H}_n \Phi_x^u =: (\Phi_x^c)^T \Phi_x^c.$$

Implementation for x and y , separately

On x :

$$\hat{\mathbf{G}}_x \approx \Phi_x^u (\Phi_x^u)^T \Rightarrow \mathbf{C}_x^u = (\Phi_x^u)^T \Phi_x^u, \quad \Phi_x^u = \left[(\Phi_{x,1}^u)^T; \dots; (\Phi_{x,n}^u)^T \right] \in \mathbb{R}^{n \times r_x},$$

$$\tilde{\hat{\mathbf{G}}}_x = \mathbf{H}_n \hat{\mathbf{G}}_x \mathbf{H}_n = \mathbf{H}_n \Phi_x^u (\Phi_x^u)^T \mathbf{H}_n, \quad \mathbf{C}_x^c = (\Phi_x^u)^T \mathbf{H}_n \mathbf{H}_n \Phi_x^u =: (\Phi_x^c)^T \Phi_x^c.$$

On y :

$$\hat{\mathbf{G}}_y \approx \Phi_y^u (\Phi_y^u)^T \Rightarrow \mathbf{C}_y^u = (\Phi_y^u)^T \Phi_y^u, \quad \Phi_y^u = \left[(\Phi_{y,1}^u)^T; \dots; (\Phi_{y,n}^u)^T \right] \in \mathbb{R}^{n \times r_y},$$

Implementation for x and y , separately

On x :

$$\hat{\mathbf{G}}_x \approx \Phi_x^u (\Phi_x^u)^T \Rightarrow \mathbf{C}_x^u = (\Phi_x^u)^T \Phi_x^u, \quad \Phi_x^u = \left[(\Phi_{x,1}^u)^T; \dots; (\Phi_{x,n}^u)^T \right] \in \mathbb{R}^{n \times r_x},$$

$$\tilde{\hat{\mathbf{G}}}_x = \mathbf{H}_n \hat{\mathbf{G}}_x \mathbf{H}_n = \mathbf{H}_n \Phi_x^u (\Phi_x^u)^T \mathbf{H}_n, \quad \mathbf{C}_x^c = (\Phi_x^u)^T \mathbf{H}_n \mathbf{H}_n \Phi_x^u =: (\Phi_x^c)^T \Phi_x^c.$$

On y :

$$\hat{\mathbf{G}}_y \approx \Phi_y^u (\Phi_y^u)^T \Rightarrow \mathbf{C}_y^u = (\Phi_y^u)^T \Phi_y^u, \quad \Phi_y^u = \left[(\Phi_{y,1}^u)^T; \dots; (\Phi_{y,n}^u)^T \right] \in \mathbb{R}^{n \times r_y},$$

$$\tilde{\hat{\mathbf{G}}}_y \approx \mathbf{H}_n \mathbf{G}_y \mathbf{H}_n = \mathbf{H}_n \Phi_y^u (\Phi_y^u)^T \mathbf{H}_n, \quad \mathbf{C}_y^c = (\Phi_y^u)^T \mathbf{H}_n \mathbf{H}_n \Phi_y^u =: (\Phi_y^c)^T \Phi_y^c.$$

Nyström-based HSIC estimator

Population quantity:

$$\begin{aligned}\text{HSIC}^2(x, y) &= \|\mathbb{E}_{xy} [\varphi(x) \otimes \psi(y)] - \mu_x \otimes \mu_y\|_{HS}^2 \\ &= \left\| \mathbb{E}_{xy} \left[(\varphi(x) - \mu_x) \otimes (\psi(y) - \mu_y) \right] \right\|_{HS}^2.\end{aligned}$$

Estimator:

$$\widehat{\text{HSIC}}_{b,N}^2(x, y) = \left\| \frac{1}{n} \sum_{i=1}^n \Phi_{x,i}^u (\Phi_{y,i}^u)^T - \left(\frac{1}{n} \sum_{i=1}^n \Phi_{x,i}^u \right) \left(\frac{1}{n} \sum_{i=1}^n \Phi_{y,i}^u \right)^T \right\|_F^2$$

Nyström-based HSIC estimator

Population quantity:

$$\begin{aligned}\text{HSIC}^2(x, y) &= \|\mathbb{E}_{xy} [\varphi(x) \otimes \psi(y)] - \mu_x \otimes \mu_y\|_{HS}^2 \\ &= \left\| \mathbb{E}_{xy} \left[(\varphi(x) - \mu_x) \otimes (\psi(y) - \mu_y) \right] \right\|_{HS}^2.\end{aligned}$$

Estimator:

$$\begin{aligned}\widehat{\text{HSIC}}_{b,N}^2(x, y) &= \left\| \frac{1}{n} \sum_{i=1}^n \Phi_{x,i}^u (\Phi_{y,i}^u)^T - \left(\frac{1}{n} \sum_{i=1}^n \Phi_{x,i}^u \right) \left(\frac{1}{n} \sum_{i=1}^n \Phi_{y,i}^u \right)^T \right\|_F^2 \\ &= \left\| \frac{1}{n} (\Phi_x^u)^T \Phi_y^u - \frac{1}{n^2} (\Phi_x^u)^T \mathbf{1}_n \mathbf{1}_n^T \Phi_y^u \right\|_F^2\end{aligned}$$

Nyström-based HSIC estimator

Population quantity:

$$\begin{aligned}\text{HSIC}^2(x, y) &= \|\mathbb{E}_{xy} [\varphi(x) \otimes \psi(y)] - \mu_x \otimes \mu_y\|_{HS}^2 \\ &= \left\| \mathbb{E}_{xy} \left[(\varphi(x) - \mu_x) \otimes (\psi(y) - \mu_y) \right] \right\|_{HS}^2.\end{aligned}$$

Estimator:

$$\begin{aligned}\widehat{\text{HSIC}}_{b,N}^2(x, y) &= \left\| \frac{1}{n} \sum_{i=1}^n \Phi_{x,i}^u (\Phi_{y,i}^u)^T - \left(\frac{1}{n} \sum_{i=1}^n \Phi_{x,i}^u \right) \left(\frac{1}{n} \sum_{i=1}^n \Phi_{y,i}^u \right)^T \right\|_F^2 \\ &= \left\| \frac{1}{n} (\Phi_x^u)^T \Phi_y^u - \frac{1}{n^2} (\Phi_x^u)^T \mathbf{1}_n \mathbf{1}_n^T \Phi_y^u \right\|_F^2 \\ &= \left\| \frac{1}{n} (\Phi_x^u)^T \underbrace{\left(\mathbf{I}_n - \frac{\mathbf{1}_n \mathbf{1}_n^T}{n} \right)}_{\mathbf{H}_n = \mathbf{H}_n^T \mathbf{H}_n} \Phi_y^u \right\|_F^2\end{aligned}$$

Nyström-based HSIC estimator

Population quantity:

$$\begin{aligned}\text{HSIC}^2(x, y) &= \|\mathbb{E}_{xy} [\varphi(x) \otimes \psi(y)] - \mu_x \otimes \mu_y\|_{HS}^2 \\ &= \left\| \mathbb{E}_{xy} \left[(\varphi(x) - \mu_x) \otimes (\psi(y) - \mu_y) \right] \right\|_{HS}^2.\end{aligned}$$

Estimator:

$$\begin{aligned}\widehat{\text{HSIC}}_{b,N}^2(x, y) &= \left\| \frac{1}{n} \sum_{i=1}^n \Phi_{x,i}^u (\Phi_{y,i}^u)^T - \left(\frac{1}{n} \sum_{i=1}^n \Phi_{x,i}^u \right) \left(\frac{1}{n} \sum_{i=1}^n \Phi_{y,i}^u \right)^T \right\|_F^2 \\ &= \left\| \frac{1}{n} (\Phi_x^u)^T \Phi_y^u - \frac{1}{n^2} (\Phi_x^u)^T \mathbf{1}_n \mathbf{1}_n^T \Phi_y^u \right\|_F^2 \\ &= \left\| \frac{1}{n} (\Phi_x^u)^T \underbrace{\left(\mathbf{I}_n - \frac{\mathbf{1}_n \mathbf{1}_n^T}{n} \right)}_{\mathbf{H}_n = \mathbf{H}_n^T \mathbf{H}_n} \Phi_y^u \right\|_F^2 = \left\| \frac{1}{n} (\Phi_x^c)^T \Phi_y^c \right\|_F^2.\end{aligned}$$

Nyström-based HSIC estimator – conclusion

$$\text{HSIC}^2(x, y) = \|C_{xy}^c\|_{HS}^2,$$
$$\widehat{\text{HSIC}}_{b,N}^2(x, y) = \left\| \frac{1}{n} (\Phi_x^c)^T \Phi_y^c \right\|_F^2.$$

Nyström-based HSIC estimator – conclusion

$$\text{HSIC}^2(x, y) = \|C_{xy}^c\|_{HS}^2,$$
$$\widehat{\text{HSIC}}_{b,N}^2(x, y) = \left\| \frac{1}{n} (\Phi_x^c)^T \Phi_y^c \right\|_F^2.$$

In short

C_{xy}^c changed to $\frac{1}{n} (\Phi_x^c)^T \Phi_y^c$, with Frobenius norm.

Nyström technique: notes

- Use $\widehat{\text{HSIC}}_{b,N}$ in
 - permutation test, or spectral approach.

Nyström technique: notes

- Use $\widehat{\text{HSIC}}_{b,N}$ in
 - permutation test, or spectral approach.
- Computational complexity (null approximation):

$$\mathcal{O}(n^3) \rightarrow \mathcal{O}(r_x^3 + r_y^3 + (r_x^2 + r_y^2)n + r_x r_y n).$$

In practice: $r_x, r_y \ll n$.

Nyström technique: notes

- Use $\widehat{\text{HSIC}}_{b,N}$ in
 - permutation test, or spectral approach.
- Computational complexity (null approximation):

$$\mathcal{O}(n^3) \rightarrow \mathcal{O}\left(r_x^3 + r_y^3 + (r_x^2 + r_y^2)n + r_x r_y n\right).$$

In practice: $r_x, r_y \ll n$.

- GP [Snelson and Ghahramani, 2006, Titsias, 2009]:
 - subset → optimized subset of size r ,
 - inducing points.

Random Fourier features

Characteristic functions: quick summary [Sasvári, 2013]

$\mathbb{P} \mapsto \phi_{\mathbb{P}}$:

$$\phi_{\mathbb{P}}(\mathbf{t}) := \mathbb{E}_{\mathbf{x} \sim \mathbb{P}} \left[e^{i \langle \mathbf{t}, \mathbf{x} \rangle} \right] = \int_{\mathbb{R}^d} e^{i \langle \mathbf{t}, \mathbf{x} \rangle} d\mathbb{P}(\mathbf{x}), \quad \mathbf{t} \in \mathbb{R}^d.$$

$\mathbb{P} \mapsto \phi_{\mathbb{P}}$:

$$\phi_{\mathbb{P}}(\mathbf{t}) := \mathbb{E}_{\mathbf{x} \sim \mathbb{P}} \left[e^{i\langle \mathbf{t}, \mathbf{x} \rangle} \right] = \int_{\mathbb{R}^d} e^{i\langle \mathbf{t}, \mathbf{x} \rangle} d\mathbb{P}(\mathbf{x}), \quad \mathbf{t} \in \mathbb{R}^d.$$

Properties:

- $\exists, \mathbb{P} \xleftrightarrow{1:1} \phi_{\mathbb{P}},$

Characteristic functions: quick summary [Sasvári, 2013]

$\mathbb{P} \mapsto \phi_{\mathbb{P}}$:

$$\phi_{\mathbb{P}}(\mathbf{t}) := \mathbb{E}_{\mathbf{x} \sim \mathbb{P}} \left[e^{i \langle \mathbf{t}, \mathbf{x} \rangle} \right] = \int_{\mathbb{R}^d} e^{i \langle \mathbf{t}, \mathbf{x} \rangle} d\mathbb{P}(\mathbf{x}), \quad \mathbf{t} \in \mathbb{R}^d.$$

Properties:

- $\exists, \mathbb{P} \xleftrightarrow{1:1} \phi_{\mathbb{P}}$,
- $|\phi_{\mathbb{P}}(\mathbf{t})| \leq 1, \phi_{\mathbb{P}}(-\mathbf{t}) = \overline{\phi_{\mathbb{P}}(\mathbf{t})} \quad \forall \mathbf{t} \in \mathbb{R}^d$.

$\mathbb{P} \mapsto \phi_{\mathbb{P}}$:

$$\phi_{\mathbb{P}}(\mathbf{t}) := \mathbb{E}_{\mathbf{x} \sim \mathbb{P}} \left[e^{i\langle \mathbf{t}, \mathbf{x} \rangle} \right] = \int_{\mathbb{R}^d} e^{i\langle \mathbf{t}, \mathbf{x} \rangle} d\mathbb{P}(\mathbf{x}), \quad \mathbf{t} \in \mathbb{R}^d.$$

Properties:

- $\exists, \mathbb{P} \xleftrightarrow{1:1} \phi_{\mathbb{P}}$,
- $|\phi_{\mathbb{P}}(\mathbf{t})| \leq 1, \phi_{\mathbb{P}}(-\mathbf{t}) = \overline{\phi_{\mathbb{P}}(\mathbf{t})} \quad \forall \mathbf{t} \in \mathbb{R}^d$.
- $\phi_{\mathbb{P}}$: uniformly continuous on \mathbb{R}^d .

$\mathbb{P} \mapsto \phi_{\mathbb{P}}$:

$$\phi_{\mathbb{P}}(\mathbf{t}) := \mathbb{E}_{\mathbf{x} \sim \mathbb{P}} \left[e^{i\langle \mathbf{t}, \mathbf{x} \rangle} \right] = \int_{\mathbb{R}^d} e^{i\langle \mathbf{t}, \mathbf{x} \rangle} d\mathbb{P}(\mathbf{x}), \quad \mathbf{t} \in \mathbb{R}^d.$$

Properties:

- $\exists, \mathbb{P} \xleftrightarrow{1:1} \phi_{\mathbb{P}}$,
- $|\phi_{\mathbb{P}}(\mathbf{t})| \leq 1, \phi_{\mathbb{P}}(-\mathbf{t}) = \overline{\phi_{\mathbb{P}}(\mathbf{t})} \quad \forall \mathbf{t} \in \mathbb{R}^d$.
- $\phi_{\mathbb{P}}$: uniformly continuous on \mathbb{R}^d .
- pd: $\sum_{i,j=1}^n \phi_{\mathbb{P}}(\mathbf{t}_i - \mathbf{t}_j) c_i \bar{c}_j \geq 0$, for $\forall n \in \mathbb{Z}^+, \mathbf{t}_i \in \mathbb{R}^d, c_i \in \mathbb{C}$.

$\mathbb{P} \mapsto \phi_{\mathbb{P}}:$

$$\phi_{\mathbb{P}}(\mathbf{t}) := \mathbb{E}_{\mathbf{x} \sim \mathbb{P}} \left[e^{i \langle \mathbf{t}, \mathbf{x} \rangle} \right] = \int_{\mathbb{R}^d} e^{i \langle \mathbf{t}, \mathbf{x} \rangle} d\mathbb{P}(\mathbf{x}), \quad \mathbf{t} \in \mathbb{R}^d.$$

Properties:

- $\exists, \mathbb{P} \xrightarrow{1:1} \phi_{\mathbb{P}},$
- $|\phi_{\mathbb{P}}(\mathbf{t})| \leq 1, \phi_{\mathbb{P}}(-\mathbf{t}) = \overline{\phi_{\mathbb{P}}(\mathbf{t})} \quad \forall \mathbf{t} \in \mathbb{R}^d.$
- $\phi_{\mathbb{P}}$: uniformly continuous on \mathbb{R}^d .
- pd: $\sum_{i,j=1}^n \phi_{\mathbb{P}}(\mathbf{t}_i - \mathbf{t}_j) c_i \bar{c}_j \geq 0$, for $\forall n \in \mathbb{Z}^+, \mathbf{t}_i \in \mathbb{R}^d, c_i \in \mathbb{C}$.

Recall

Bochner's theorem & $\mathbf{G} \geq 0$ definition of kernels!

Characteristic functions: continued

Operations, closedness:

- Sum of independent variables:

$$\phi_{\sum_{i=1}^n \mathbf{x}_i}(\mathbf{t}) = \prod_{i=1}^n \phi_{\mathbf{x}_i}(\mathbf{t}), \quad \forall \mathbf{t} \in \mathbb{R}^d.$$

Characteristic functions: continued

Operations, closedness:

- Sum of independent variables:

$$\phi_{\sum_{i=1}^n \mathbf{x}_i}(\mathbf{t}) = \prod_{i=1}^n \phi_{\mathbf{x}_i}(\mathbf{t}), \quad \forall \mathbf{t} \in \mathbb{R}^d.$$

- Affine transformation ($\mathbf{A} \in \mathbb{R}^{d_1 \times d_2}$):

$$\phi_{\mathbf{Ax} + \mathbf{b}}(\mathbf{t}) = e^{i\langle \mathbf{t}, \mathbf{b} \rangle} \phi_{\mathbf{x}}(\mathbf{A}^T \mathbf{t}), \quad \forall \mathbf{t} \in \mathbb{R}^d.$$

Characteristic functions: continued

Operations, closedness:

- Sum of independent variables:

$$\phi_{\sum_{i=1}^n \mathbf{x}_i}(\mathbf{t}) = \prod_{i=1}^n \phi_{\mathbf{x}_i}(\mathbf{t}), \quad \forall \mathbf{t} \in \mathbb{R}^d.$$

- Affine transformation ($\mathbf{A} \in \mathbb{R}^{d_1 \times d_2}$):

$$\phi_{\mathbf{Ax} + \mathbf{b}}(\mathbf{t}) = e^{i\langle \mathbf{t}, \mathbf{b} \rangle} \phi_{\mathbf{x}}(\mathbf{A}^T \mathbf{t}), \quad \forall \mathbf{t} \in \mathbb{R}^d.$$

- Concatenation of independent variables: $\mathbf{x} = [\mathbf{x}_1; \dots; \mathbf{x}_n]$

$$\phi_{\mathbf{x}}(\mathbf{t}) = \prod_{i=1}^n \phi_{\mathbf{x}_i}(\mathbf{t}_i), \quad \mathbf{t} = [\mathbf{t}_1; \dots; \mathbf{t}_n] \in \mathbb{R}^d.$$

Characteristic functions: continued

Operations, closedness:

- Sum of independent variables:

$$\phi_{\sum_{i=1}^n \mathbf{x}_i}(\mathbf{t}) = \prod_{i=1}^n \phi_{\mathbf{x}_i}(\mathbf{t}), \quad \forall \mathbf{t} \in \mathbb{R}^d.$$

- Affine transformation ($\mathbf{A} \in \mathbb{R}^{d_1 \times d_2}$):

$$\phi_{\mathbf{Ax} + \mathbf{b}}(\mathbf{t}) = e^{i\langle \mathbf{t}, \mathbf{b} \rangle} \phi_{\mathbf{x}}(\mathbf{A}^T \mathbf{t}), \quad \forall \mathbf{t} \in \mathbb{R}^d.$$

- Concatenation of independent variables: $\mathbf{x} = [\mathbf{x}_1; \dots; \mathbf{x}_n]$

$$\phi_{\mathbf{x}}(\mathbf{t}) = \prod_{i=1}^n \phi_{\mathbf{x}_i}(\mathbf{t}_i), \quad \mathbf{t} = [\mathbf{t}_1; \dots; \mathbf{t}_n] \in \mathbb{R}^d.$$

Recall

Distance covariance!

Characteristic functions: continued

Moment condition on $\mathbb{P} \Rightarrow$ differentiability of $\phi_{\mathbb{P}}$.

Assume that exists:

$$M_{\mathbf{a}} = \mathbb{E}_{\mathbf{x} \sim \mathbb{P}} [\mathbf{x}^{\mathbf{a}}] \quad \mathbf{a} \in \mathbb{N}^d, \quad \left(\mathbf{x}^{\mathbf{a}} := \prod_{i=1}^d x_i^{a_i} \right).$$

Then $\exists \partial^{\mathbf{a}} \phi_{\mathbb{P}}$ and

$$\partial^{\mathbf{a}} \phi_{\mathbb{P}}(\mathbf{t}) = i^{|\mathbf{a}|} \int_{\mathbb{R}^d} \mathbf{x}^{\mathbf{a}} e^{i\langle \mathbf{t}, \mathbf{x} \rangle} d\mathbb{P}(x), \quad \forall \mathbf{t} \in \mathbb{R}^d,$$

$$\partial^{\mathbf{a}} \phi_{\mathbb{P}}(\mathbf{0}) = i^{|\mathbf{a}|} M_{\mathbf{a}}, \quad |\mathbf{a}| = \sum_{i=1}^d a_i,$$

and $\partial^{\mathbf{a}} \phi_{\mathbb{P}}$ is uniformly continuous.

RFF idea

- k : continuous bounded & shift-invariant on \mathbb{R}^d [$k(\mathbf{x}, \mathbf{y}) = k_0(\mathbf{x} - \mathbf{y})$].
By Bochner:

$$k(\mathbf{x}, \mathbf{y}) = \int_{\mathbb{R}^d} \underbrace{e^{i\omega^T(\mathbf{x}-\mathbf{y})}}_{\cos(\omega^T(\mathbf{x}-\mathbf{y})) + i \sin(\omega^T(\mathbf{x}-\mathbf{y}))} d\Lambda(\omega)$$

RFF idea

- k : continuous bounded & shift-invariant on \mathbb{R}^d [$k(\mathbf{x}, \mathbf{y}) = k_0(\mathbf{x} - \mathbf{y})$].
By Bochner:

$$\begin{aligned} k(\mathbf{x}, \mathbf{y}) &= \int_{\mathbb{R}^d} \underbrace{e^{i\omega^T(\mathbf{x}-\mathbf{y})}}_{\cos(\omega^T(\mathbf{x}-\mathbf{y}))+i\sin(\omega^T(\mathbf{x}-\mathbf{y}))} d\Lambda(\omega) \\ &= \int_{\mathbb{R}^d} \cos(\omega^T(\mathbf{x} - \mathbf{y})) d\Lambda(\omega). \end{aligned}$$

RFF idea

- k : continuous bounded & shift-invariant on \mathbb{R}^d [$k(\mathbf{x}, \mathbf{y}) = k_0(\mathbf{x} - \mathbf{y})$].
By Bochner:

$$\begin{aligned} k(\mathbf{x}, \mathbf{y}) &= \int_{\mathbb{R}^d} \underbrace{e^{i\omega^T(\mathbf{x}-\mathbf{y})}}_{\cos(\omega^T(\mathbf{x}-\mathbf{y}))+i\sin(\omega^T(\mathbf{x}-\mathbf{y}))} d\Lambda(\omega) \\ &= \int_{\mathbb{R}^d} \cos(\omega^T(\mathbf{x} - \mathbf{y})) d\Lambda(\omega). \end{aligned}$$

- RFF trick [Rahimi and Recht, 2007] (MC): $\boldsymbol{\omega}_{1:m} := (\omega_j)_{j=1}^m \stackrel{i.i.d.}{\sim} \Lambda$,

$$\hat{k}(\mathbf{x}, \mathbf{y}) = \frac{1}{m} \sum_{j=1}^m \cos(\omega_j^T(\mathbf{x} - \mathbf{y}))$$

RFF idea

- k : continuous bounded & shift-invariant on \mathbb{R}^d [$k(\mathbf{x}, \mathbf{y}) = k_0(\mathbf{x} - \mathbf{y})$].
By Bochner:

$$\begin{aligned} k(\mathbf{x}, \mathbf{y}) &= \int_{\mathbb{R}^d} \underbrace{e^{i\omega^T(\mathbf{x}-\mathbf{y})}}_{\cos(\omega^T(\mathbf{x}-\mathbf{y}))+i\sin(\omega^T(\mathbf{x}-\mathbf{y}))} d\Lambda(\omega) \\ &= \int_{\mathbb{R}^d} \cos(\omega^T(\mathbf{x} - \mathbf{y})) d\Lambda(\omega). \end{aligned}$$

- RFF trick [Rahimi and Recht, 2007] (MC): $\omega_{1:m} := (\omega_j)_{j=1}^m \stackrel{i.i.d.}{\sim} \Lambda$,

$$\hat{k}(\mathbf{x}, \mathbf{y}) = \frac{1}{m} \sum_{j=1}^m \cos(\omega_j^T(\mathbf{x} - \mathbf{y})) = \int_{\mathbb{R}^d} \cos(\omega^T(\mathbf{x} - \mathbf{y})) d\Lambda_m(\omega).$$

RFF idea

- k : continuous bounded & shift-invariant on \mathbb{R}^d [$k(\mathbf{x}, \mathbf{y}) = k_0(\mathbf{x} - \mathbf{y})$].
By Bochner:

$$\begin{aligned} k(\mathbf{x}, \mathbf{y}) &= \int_{\mathbb{R}^d} \underbrace{e^{i\omega^T(\mathbf{x}-\mathbf{y})}}_{\cos(\omega^T(\mathbf{x}-\mathbf{y}))+i\sin(\omega^T(\mathbf{x}-\mathbf{y}))} d\Lambda(\omega) \\ &= \int_{\mathbb{R}^d} \cos(\omega^T(\mathbf{x} - \mathbf{y})) d\Lambda(\omega). \end{aligned}$$

- RFF trick [Rahimi and Recht, 2007] (MC): $\omega_{1:m} := (\omega_j)_{j=1}^m \stackrel{i.i.d.}{\sim} \Lambda$,

$$\hat{k}(\mathbf{x}, \mathbf{y}) = \frac{1}{m} \sum_{j=1}^m \cos(\omega_j^T(\mathbf{x} - \mathbf{y})) = \int_{\mathbb{R}^d} \cos(\omega^T(\mathbf{x} - \mathbf{y})) d\Lambda_m(\omega).$$

Recall (characteristic kernels)

We saw many $k \rightarrow \Lambda$ examples!

Questions

- Why is RFF useful?
- Does it converge ($k - \hat{k} \approx 0$)? Rates?

Why is RFF useful?

Kernel approximation:

$$\hat{k}(\mathbf{x}, \mathbf{y}) = \frac{1}{m} \sum_{j=1}^m \cos(\omega_j^T (\mathbf{x} - \mathbf{y})).$$

Why is RFF useful?

Kernel approximation:

$$\hat{k}(\mathbf{x}, \mathbf{y}) = \frac{1}{m} \sum_{j=1}^m \cos(\omega_j^T (\mathbf{x} - \mathbf{y})).$$

By the trigonometric identity:

$$\cos(a - b) = \cos(a)\cos(b) + \sin(a)\sin(b),$$

Why is RFF useful?

Kernel approximation:

$$\hat{k}(\mathbf{x}, \mathbf{y}) = \frac{1}{m} \sum_{j=1}^m \cos(\omega_j^T (\mathbf{x} - \mathbf{y})).$$

By the trigonometric identity:

$$\cos(a - b) = \cos(a)\cos(b) + \sin(a)\sin(b),$$

$$\hat{k}(\mathbf{x}, \mathbf{y}) = \left\langle \hat{\phi}(\mathbf{x}), \hat{\phi}(\mathbf{y}) \right\rangle_{\mathbb{R}^{2m}},$$

$$\begin{aligned}\hat{\phi}(\mathbf{x}) &= \frac{1}{\sqrt{m}} \left[\cos(\omega_1^T \mathbf{x}); \dots; \cos(\omega_m^T \mathbf{x}); \right. \\ &\quad \left. \sin(\omega_1^T \mathbf{x}); \dots; \sin(\omega_m^T \mathbf{x}) \right] \in \mathbb{R}^{2m}.\end{aligned}$$

Why is RFF useful?

Kernel approximation:

$$\hat{k}(\mathbf{x}, \mathbf{y}) = \frac{1}{m} \sum_{j=1}^m \cos(\omega_j^T (\mathbf{x} - \mathbf{y})).$$

By the trigonometric identity:

$$\cos(a - b) = \cos(a)\cos(b) + \sin(a)\sin(b),$$

$$\hat{k}(\mathbf{x}, \mathbf{y}) = \left\langle \hat{\phi}(\mathbf{x}), \hat{\phi}(\mathbf{y}) \right\rangle_{\mathbb{R}^{2m}},$$

$$\begin{aligned}\hat{\phi}(\mathbf{x}) &= \frac{1}{\sqrt{m}} \left[\cos(\omega_1^T \mathbf{x}); \dots; \cos(\omega_m^T \mathbf{x}); \right. \\ &\quad \left. \sin(\omega_1^T \mathbf{x}); \dots; \sin(\omega_m^T \mathbf{x}) \right] \in \mathbb{R}^{2m}.\end{aligned}$$

Key

We got (random) explicit feature maps!

RFF application in independence testing

Previous slide ⇒

$$(\Phi_x^u)^T := \left[\hat{\phi}(x_1); \dots; \hat{\phi}(x_n) \right], (\Phi_y^u)^T := \left[\hat{\phi}(y_1); \dots; \hat{\phi}(y_n) \right],$$

$$\mathbf{G}_x \approx \Phi_x^u (\Phi_x^u)^T, \quad \mathbf{G}_y \approx \Phi_y^u (\Phi_y^u)^T,$$

RFF application in independence testing

Previous slide \Rightarrow

$$(\Phi_x^u)^T := \left[\hat{\phi}(x_1); \dots; \hat{\phi}(x_n) \right], (\Phi_y^u)^T := \left[\hat{\phi}(y_1); \dots; \hat{\phi}(y_n) \right],$$

$$\mathbf{G}_x \approx \Phi_x^u (\Phi_x^u)^T, \quad \mathbf{G}_y \approx \Phi_y^u (\Phi_y^u)^T,$$

and hence

$$\widehat{\text{HSIC}}_{b,\text{RFF}}^2(x, y) = \left\| \frac{1}{n} \sum_{i=1}^n \Phi_{x,i}^u (\Phi_{y,i}^u)^T - \left(\frac{1}{n} \sum_{i=1}^n \Phi_{x,i}^u \right) \left(\frac{1}{n} \sum_{i=1}^n \Phi_{y,i}^u \right)^T \right\|_F^2$$

RFF application in independence testing

Previous slide \Rightarrow

$$(\Phi_x^u)^T := \left[\hat{\phi}(x_1); \dots; \hat{\phi}(x_n) \right], (\Phi_y^u)^T := \left[\hat{\phi}(y_1); \dots; \hat{\phi}(y_n) \right],$$
$$\mathbf{G}_x \approx \Phi_x^u (\Phi_x^u)^T, \quad \mathbf{G}_y \approx \Phi_y^u (\Phi_y^u)^T,$$

and hence

$$\widehat{\text{HSIC}}_{b,\text{RFF}}^2(x, y) = \left\| \frac{1}{n} \sum_{i=1}^n \Phi_{x,i}^u (\Phi_{y,i}^u)^T - \left(\frac{1}{n} \sum_{i=1}^n \Phi_{x,i}^u \right) \left(\frac{1}{n} \sum_{i=1}^n \Phi_{y,i}^u \right)^T \right\|_F^2$$
$$= \dots = \left\| \frac{1}{n} (\Phi_x^c)^T \Phi_y^c \right\|_F^2.$$

Briefly

We simply '**overloaded**' the features with the RFF ones.

Some further RFF-accelerated measures

- **KCCA** [Lopez-Paz et al., 2014].
- **MMD** [Sutherland and Schneider, 2015,
Zhao and Meng, 2015, Lopez-Paz, 2016].

RFF: in kernel ridge regression

- Given: $\{(x_i, y_i)\}_{i=1}^\ell$.
- Task: find $f \in \mathcal{H}_k$ s.t. $f(x_i) \approx y_i$,

$$J(f) = \frac{1}{\ell} \sum_{i=1}^{\ell} [f(x_i) - y_i]^2 + \lambda \|f\|_{\mathcal{H}_k}^2 \rightarrow \min_{f \in \mathcal{H}_k} \quad (\lambda > 0).$$

RFF: in kernel ridge regression

- Given: $\{(x_i, y_i)\}_{i=1}^\ell$.
- Task: find $f \in \mathcal{H}_k$ s.t. $f(x_i) \approx y_i$,

$$J(f) = \frac{1}{\ell} \sum_{i=1}^{\ell} [f(x_i) - y_i]^2 + \lambda \|f\|_{\mathcal{H}_k}^2 \rightarrow \min_{f \in \mathcal{H}_k} \quad (\lambda > 0).$$

- Analytical solution, $\mathcal{O}(\ell^3)$ – **expensive**:

$$f(x) = [k(x_1, x), \dots, k(x_\ell, x)](\mathbf{G} + \lambda \ell I)^{-1}[y_1; \dots; y_\ell],$$

$$\mathbf{G} = [k(x_i, x_j)]_{i,j=1}^\ell.$$

RFF: in kernel ridge regression

- Given: $\{(x_i, y_i)\}_{i=1}^\ell$.
- Task: find $f \in \mathcal{H}_k$ s.t. $f(x_i) \approx y_i$,

$$J(f) = \frac{1}{\ell} \sum_{i=1}^{\ell} [f(x_i) - y_i]^2 + \lambda \|f\|_{\mathcal{H}_k}^2 \rightarrow \min_{f \in \mathcal{H}_k} \quad (\lambda > 0).$$

- Analytical solution, $\mathcal{O}(\ell^3)$ – expensive:

$$f(x) = [k(x_1, x), \dots, k(x_\ell, x)](\mathbf{G} + \lambda \ell I)^{-1}[y_1; \dots; y_\ell],$$

$$\mathbf{G} = [k(x_i, x_j)]_{i,j=1}^\ell.$$

- Idea: $\hat{\mathbf{G}}$, matrix-inversion lemma, fast primal solvers \rightarrow RFF.

Approximation quality

- Hoeffding inequality + union bound
[Rahimi and Recht, 2007, Sutherland and Schneider, 2015]:

$$\|k - \hat{k}\|_{L^\infty(\mathcal{S})} = \mathcal{O}_p \left(|\mathcal{S}| \frac{\sqrt{\log(m)}}{\sqrt{m}} \right).$$

Approximation quality

- Hoeffding inequality + union bound
[Rahimi and Recht, 2007, Sutherland and Schneider, 2015]:

$$\|k - \hat{k}\|_{L^\infty(\mathcal{S})} = \mathcal{O}_p \left(|\mathcal{S}| \frac{\sqrt{\log(m)}}{\sqrt{m}} \right).$$

- ECFs [Csörgő and Totik, 1983]: $|\mathcal{S}_m| = e^{o(m)}$ – optimal rate, asymptotic!

Approximation quality

- Hoeffding inequality + union bound
[Rahimi and Recht, 2007, Sutherland and Schneider, 2015]:

$$\|k - \hat{k}\|_{L^\infty(\mathcal{S})} = \mathcal{O}_p \left(|\mathcal{S}| \frac{\sqrt{\log(m)}}{\sqrt{m}} \right).$$

- ECFs [Csörgő and Totik, 1983]: $|\mathcal{S}_m| = e^{o(m)}$ – optimal rate, asymptotic!
- Finite-sample L^∞ -bound [Sriperumbudur and Szabó, 2015] $\xrightarrow{\text{spec.}}$

$$\|k - \hat{k}\|_{L^\infty(\mathcal{S})} = \mathcal{O}_{a.s.} \left(\frac{\sqrt{\log |\mathcal{S}|}}{\sqrt{m}} \right).$$

Approximation quality

- Hoeffding inequality + union bound
[Rahimi and Recht, 2007, Sutherland and Schneider, 2015]:

$$\|k - \hat{k}\|_{L^\infty(\mathcal{S})} = \mathcal{O}_p \left(|\mathcal{S}| \frac{\sqrt{\log(m)}}{\sqrt{m}} \right).$$

- ECFs [Csörgő and Totik, 1983]: $|\mathcal{S}_m| = e^{o(m)}$ – optimal rate, asymptotic!
- Finite-sample L^∞ -bound [Sriperumbudur and Szabó, 2015] $\xrightarrow{\text{spec.}}$

$$\|k - \hat{k}\|_{L^\infty(\mathcal{S})} = \mathcal{O}_{a.s.} \left(\frac{\sqrt{\log |\mathcal{S}|}}{\sqrt{m}} \right).$$

- RFF in ridge regression [Rudi and Rosasco, 2017], kernel PCA [Sriperumbudur and Sterge, 2018, Ullah et al., 2018], classification with 0-1 loss [Sun et al., 2018], Lipschitz losses [Li et al., 2018].

Optimal $\|k - \hat{k}\|_{L^\infty(\mathcal{S})}$: proof idea

- Empirical process form [$\mathbb{P}g := \int g d\mathbb{P}; \textcolor{brown}{g}(\omega) = \cos(\omega^T(\mathbf{x} - \mathbf{y}))$]:

$$\sup_{\mathbf{x}, \mathbf{y} \in \mathcal{S}} |k(\mathbf{x}, \mathbf{y}) - \hat{k}(\mathbf{x}, \mathbf{y})| = \sup_{\textcolor{brown}{g} \in \mathcal{G}} |\Lambda \textcolor{brown}{g} - \Lambda_m \textcolor{brown}{g}| .$$

Optimal $\|k - \hat{k}\|_{L^\infty(\mathcal{S})}$: proof idea

- Empirical process form [$\mathbb{P}g := \int g d\mathbb{P}; \textcolor{brown}{g}(\omega) = \cos(\omega^T(\mathbf{x} - \mathbf{y}))$]:

$$\sup_{\mathbf{x}, \mathbf{y} \in \mathcal{S}} |k(\mathbf{x}, \mathbf{y}) - \hat{k}(\mathbf{x}, \mathbf{y})| = \sup_{\textcolor{brown}{g} \in \mathcal{G}} |\Lambda \textcolor{brown}{g} - \Lambda_m \textcolor{brown}{g}| =: \|\Lambda - \Lambda_m\|_{\mathcal{G}}.$$

Optimal $\|k - \hat{k}\|_{L^\infty(\mathcal{S})}$: proof idea

- Empirical process form [$\mathbb{P}g := \int g d\mathbb{P}; \textcolor{brown}{g}(\omega) = \cos(\omega^T(\mathbf{x} - \mathbf{y}))$]:

$$\sup_{\mathbf{x}, \mathbf{y} \in \mathcal{S}} |k(\mathbf{x}, \mathbf{y}) - \hat{k}(\mathbf{x}, \mathbf{y})| = \sup_{\textcolor{brown}{g} \in \mathcal{G}} |\Lambda \textcolor{brown}{g} - \Lambda_m \textcolor{brown}{g}| =: \|\Lambda - \Lambda_m\|_{\mathcal{G}}.$$

- $f(\omega_{1:m}) = \|\Lambda - \Lambda_m\|_{\mathcal{G}}$ concentrates (bounded difference):

$$\|\Lambda - \Lambda_m\|_{\mathcal{G}} \lesssim \mathbb{E}_{\omega_{1:m}} \|\Lambda - \Lambda_m\|_{\mathcal{G}} + \frac{1}{\sqrt{m}}.$$

Optimal $\|k - \hat{k}\|_{L^\infty(\mathcal{S})}$: proof idea

- Empirical process form $[\mathbb{P}g := \int g d\mathbb{P}; g(\omega) = \cos(\omega^T(\mathbf{x} - \mathbf{y}))]$:

$$\sup_{\mathbf{x}, \mathbf{y} \in \mathcal{S}} |k(\mathbf{x}, \mathbf{y}) - \hat{k}(\mathbf{x}, \mathbf{y})| = \sup_{g \in \mathcal{G}} |\Lambda g - \Lambda_m g| =: \|\Lambda - \Lambda_m\|_{\mathcal{G}}.$$

- $f(\omega_{1:m}) = \|\Lambda - \Lambda_m\|_{\mathcal{G}}$ concentrates (bounded difference):

$$\|\Lambda - \Lambda_m\|_{\mathcal{G}} \lesssim \mathbb{E}_{\omega_{1:m}} \|\Lambda - \Lambda_m\|_{\mathcal{G}} + \frac{1}{\sqrt{m}}.$$

- \mathcal{G} is 'nice' (uniformly bounded, separable Carathéodory) \Rightarrow

$$\mathbb{E}_{\omega_{1:m}} \|\Lambda - \Lambda_m\|_{\mathcal{G}} \lesssim \mathbb{E}_{\omega_{1:m}} \underbrace{\mathcal{R}(\mathcal{G}, \omega_{1:m})}_{\mathbb{E}_\epsilon \sup_{g \in \mathcal{G}} \left| \frac{1}{m} \sum_{j=1}^m \epsilon_j g(\omega_j) \right|} .$$

Proof idea – continued

- Using Dudley's entropy bound:

$$\mathcal{R}(\mathcal{G}, \omega_{1:m}) \lesssim \frac{1}{\sqrt{m}} \int_0^{|\mathcal{G}|_{L^2(\Lambda_m)}} \sqrt{\log \mathcal{N}(\mathcal{G}, L^2(\Lambda_m), r)} dr.$$

Proof idea – continued

- Using Dudley's entropy bound:

$$\mathcal{R}(\mathcal{G}, \omega_{1:m}) \lesssim \frac{1}{\sqrt{m}} \int_0^{|\mathcal{G}|_{L^2(\Lambda_m)}} \sqrt{\log \mathcal{N}(\mathcal{G}, L^2(\Lambda_m), r)} dr.$$

- \mathcal{G} is smoothly parameterized by a compact set \Rightarrow

$$\mathcal{N}(\mathcal{G}, L^2(\Lambda_m), r) \leq \left(\frac{4|S|A}{r} + 1 \right)^d, \quad A(\omega_{1:m}) = \sqrt{\frac{1}{m} \sum_{j=1}^m \|\omega_j\|_2^2}.$$

Proof idea – continued

- Using Dudley's entropy bound:

$$\mathcal{R}(\mathcal{G}, \omega_{1:m}) \lesssim \frac{1}{\sqrt{m}} \int_0^{|\mathcal{G}|_{L^2(\Lambda_m)}} \sqrt{\log \mathcal{N}(\mathcal{G}, L^2(\Lambda_m), r)} dr.$$

- \mathcal{G} is smoothly parameterized by a compact set \Rightarrow

$$\mathcal{N}(\mathcal{G}, L^2(\Lambda_m), r) \leq \left(\frac{4|S|A}{r} + 1 \right)^d, \quad A(\omega_{1:m}) = \sqrt{\frac{1}{m} \sum_{j=1}^m \|\omega_j\|_2^2}.$$

- Putting together $[|\mathcal{G}|_{L^2(\Lambda_m)} \leq 2, \text{ Jensen inequality}]$ we get ...

Theorem (Finite-sample, asymptotically optimal uniform bound for RFF)

Let k be continuous, bounded, shift-invariant, and
 $\sigma^2 := \int \|\omega\|^2 d\Lambda(\omega) < \infty$. Then for $\forall \tau > 0$ and compact set
 $\mathcal{S} \subset \mathbb{R}^d$

$$\Lambda^m \left(\|\hat{k} - k\|_{L^\infty(\mathcal{S})} \geq \frac{h(d, |\mathcal{S}|, \sigma) + \sqrt{2\tau}}{\sqrt{m}} \right) \leq e^{-\tau},$$

$$h(d, |\mathcal{S}|, \sigma) := 32\sqrt{2d \log(2|\mathcal{S}| + 1)} + 16\sqrt{\frac{2d}{\log(2|\mathcal{S}| + 1)}} + \\ 32\sqrt{2d \log(\sigma + 1)}.$$

Empirical process theory: motivation

The object of interest:

$$\sup_{f \in \mathcal{F}} |\mathbb{P}f - \mathbb{P}_n f|.$$

Empirical process theory: motivation

The object of interest:

$$\sup_{f \in \mathcal{F}} |\mathbb{P}f - \mathbb{P}_n f|.$$

Original motivation:

- F : cdf, F_n : empirical cdf.

Empirical process theory: motivation

The object of interest:

$$\sup_{f \in \mathcal{F}} |\mathbb{P}f - \mathbb{P}_n f|.$$

Original motivation:

- F : cdf, F_n : empirical cdf.
- Glivenko-Cantelli theorem:

$$0 \xleftarrow{\infty \leftarrow n} \|F - F_n\|_\infty = \sup_x |F(x) - F_n(x)|$$

Empirical process theory: motivation

The object of interest:

$$\sup_{f \in \mathcal{F}} |\mathbb{P}f - \mathbb{P}_n f|.$$

Original motivation:

- F : cdf, F_n : empirical cdf.
- Glivenko-Cantelli theorem:

$$\begin{aligned} 0 &\xleftarrow{\text{---}} \|F - F_n\|_\infty = \sup_x |F(x) - F_n(x)| \\ &= \sup_{f \in \mathcal{F}} |\mathbb{P}f - \mathbb{P}_n f|, \quad \mathcal{F} = \{\chi_{(\infty, x)} : x \in \mathbb{R}^d\}. \end{aligned}$$

Ref: [van der Vaart and Wellner, 1996, van der Vaart, 1998, van de Geer, 2009].

- One can also get:
 - $L^p(\mathcal{S})$ results (\Leftarrow uniform bound, type of L^p).
 - bounds for $\partial k^{\mathbf{p}, \mathbf{q}}$ [Szabó and Sriperumbudur, 2019].

Notes on RFF: L^p bounds, kernel derivatives

- One can also get:
 - $L^p(\mathcal{S})$ results (\Leftarrow uniform bound, type of L^p).
 - bounds for $\partial k^{p,q}$ [Szabó and Sriperumbudur, 2019].
- Kernel derivatives: $\frac{\partial^{p,q} f(x,y)}{\partial^p x \partial^q y}$,
 - nonlinear variable selection [Rosasco et al., 2010, Rosasco et al., 2013],

Notes on RFF: L^p bounds, kernel derivatives

- One can also get:
 - $L^p(\mathcal{S})$ results (\Leftarrow uniform bound, type of L^p).
 - bounds for $\partial k^{p,q}$ [Szabó and Sriperumbudur, 2019].
- Kernel derivatives: $\frac{\partial^{p,q} f(x,y)}{\partial^p x \partial^q y}$,
 - nonlinear variable selection [Rosasco et al., 2010, Rosasco et al., 2013],
 - infinite-dimensional exponential family fitting [Sriperumbudur et al., 2017].

Let us look at the examples!

Nonlinear variable selection

- Objective function, $\lambda > 0$:

$$J(f) = \frac{1}{n} \sum_{i=1}^n [f(x_i) - y_i]^2 + \lambda \sum_{j=1}^d \|\partial_j f\| \rightarrow \min_{f \in \mathcal{H}_k},$$

$$\|g\| := \sqrt{\frac{1}{n} \sum_{i=1}^n |g(x_i)|^2}.$$

- Objective function, $\lambda > 0$:

$$J(f) = \frac{1}{n} \sum_{i=1}^n [f(x_i) - y_i]^2 + \lambda \sum_{j=1}^d \|\partial_j f\| \rightarrow \min_{f \in \mathcal{H}_k},$$

$$\|g\| := \sqrt{\frac{1}{n} \sum_{i=1}^n |g(x_i)|^2}.$$

- Intuition:

- if f does not depend on variable j , then $\partial_j f = 0$.

Infinite-dimensional exponential family (\mathbb{R}^d)

- Exponential family:

$$p_{\theta}(\mathbf{x}) \propto e^{\langle \theta, T(\mathbf{x}) \rangle},$$

where θ : natural parameter, $T(\mathbf{x})$: sufficient statistics.

Infinite-dimensional exponential family (\mathbb{R}^d)

- Exponential family:

$$p_{\theta}(\mathbf{x}) \propto e^{\langle \theta, T(\mathbf{x}) \rangle},$$

where θ : natural parameter, $T(\mathbf{x})$: sufficient statistics.

- Examples: normal, exponential, gamma, χ^2 , beta, ...

Infinite-dimensional exponential family (\mathbb{R}^d)

- Exponential family:

$$p_{\theta}(\mathbf{x}) \propto e^{\langle \theta, T(\mathbf{x}) \rangle},$$

where θ : natural parameter, $T(\mathbf{x})$: sufficient statistics.

- Examples: normal, exponential, gamma, χ^2 , beta, ...
- InfiniteD generalization:

$$p_f(\mathbf{x}) \propto e^{f(\mathbf{x})} = e^{\langle f, k(\cdot, \mathbf{x}) \rangle_{\mathcal{H}_k}}.$$

Infinite-dimensional exponential family (\mathbb{R}^d)

- Exponential family:

$$p_{\theta}(\mathbf{x}) \propto e^{\langle \theta, T(\mathbf{x}) \rangle},$$

where θ : natural parameter, $T(\mathbf{x})$: sufficient statistics.

- Examples: normal, exponential, gamma, χ^2 , beta, ...
- InfiniteD generalization:

$$p_f(\mathbf{x}) \propto e^{f(\mathbf{x})} = e^{\langle f, k(\cdot, \mathbf{x}) \rangle_{\mathcal{H}_k}}.$$

Fitting idea (score matching, Fischer divergence):

$$J(p_*, p_f) := \int p_*(\mathbf{x}) \left\| \frac{\partial \log p_*(\mathbf{x})}{\partial \mathbf{x}} - \frac{\partial \log p_f(\mathbf{x})}{\partial \mathbf{x}} \right\|_2^2 d\mathbf{x} \rightarrow \min_{f \in \mathcal{H}_k} .$$

Notes on RFF: operator-valued extension

- Standard setup: $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$

$$\mathcal{H}_k = \{f : \mathcal{X} \rightarrow \mathbb{R} \mid \dots\}.$$

Notes on RFF: operator-valued extension

- Standard setup: $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$

$$\mathcal{H}_k = \{f : \mathcal{X} \rightarrow \mathbb{R} \mid \dots\}.$$

- Operator-valued case:

$$\mathcal{H}_k = \{f : \mathcal{X} \rightarrow \mathbb{Y} \mid \dots\}$$

Notes on RFF: operator-valued extension

- Standard setup: $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$

$$\mathcal{H}_k = \{f : \mathcal{X} \rightarrow \mathbb{R} \mid \dots\}.$$

- Operator-valued case:

$$\mathcal{H}_k = \{f : \mathcal{X} \rightarrow \mathcal{Y} \mid \dots\}, \quad k : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y}).$$

\mathcal{Y} : (separable) Hilbert. Example: $\mathcal{Y} = \mathbb{R}^d$, $\mathcal{L}(\mathcal{Y}) = \mathbb{R}^{d \times d}$.

Notes on RFF: operator-valued extension

- Standard setup: $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$

$$\mathcal{H}_k = \{f : \mathcal{X} \rightarrow \mathbb{R} \mid \dots\}.$$

- Operator-valued case:

$$\mathcal{H}_k = \{f : \mathcal{X} \rightarrow \mathcal{Y} \mid \dots\}, \quad k : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{L}(\mathcal{Y}).$$

\mathcal{Y} : (separable) Hilbert. Example: $\mathcal{Y} = \mathbb{R}^d$, $\mathcal{L}(\mathcal{Y}) = \mathbb{R}^{d \times d}$.

- RFF idea

- works [Brault et al., 2016]; $(\mathbb{R}^d, +) \rightarrow \text{LCA}$: ✓
- open question: 'optimal' rates.

Nyström method, RFF: the end.

Linear-time two-sample testing: analytic representations.

- Recall:

$$\text{MMD}(\mathbb{P}, \mathbb{Q}) = \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k}.$$

Linear-time 2-sample test [Chwialkowski et al., 2015]

- Recall:

$$\text{MMD}(\mathbb{P}, \mathbb{Q}) = \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k}.$$

- Idea: change the norm

$$\rho(\mathbb{P}, \mathbb{Q}) := \rho \left(\mathbb{P}, \mathbb{Q}; \{\mathbf{v}_j\}_{j=1}^J \right) := \sqrt{\frac{1}{J} \sum_{j=1}^J [\mu_{\mathbb{P}}(\mathbf{v}_j) - \mu_{\mathbb{Q}}(\mathbf{v}_j)]^2}$$

with random $\{\mathbf{v}_j\}_{j=1}^J$ test locations.

Linear-time 2-sample test [Chwialkowski et al., 2015]

- Recall:

$$\text{MMD}(\mathbb{P}, \mathbb{Q}) = \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k}.$$

- Idea: change the norm

$$\rho(\mathbb{P}, \mathbb{Q}) := \rho \left(\mathbb{P}, \mathbb{Q}; \{\mathbf{v}_j\}_{j=1}^J \right) := \sqrt{\frac{1}{J} \sum_{j=1}^J [\mu_{\mathbb{P}}(\mathbf{v}_j) - \mu_{\mathbb{Q}}(\mathbf{v}_j)]^2}$$

with random $\{\mathbf{v}_j\}_{j=1}^J$ test locations.

Is ρ a random metric? How do we estimate it? Distribution under H_0 ?

What is a random metric?

In short

It is a metric almost surely (assumptions: next slide).

What is a random metric?

In short

It is a **metric almost surely** (assumptions: next slide).

In other words,

- $\rho(\mathbb{P}, \mathbb{Q}) \geq 0$, $\rho(\mathbb{P}, \mathbb{Q}) = 0 \Leftrightarrow \mathbb{P} = \mathbb{Q}$ **almost surely**.

What is a random metric?

In short

It is a **metric almost surely** (assumptions: next slide).

In other words,

- $\rho(\mathbb{P}, \mathbb{Q}) \geq 0$, $\rho(\mathbb{P}, \mathbb{Q}) = 0 \Leftrightarrow \mathbb{P} = \mathbb{Q}$ **almost surely**.
- $\rho(\mathbb{P}, \mathbb{Q}) = \rho(\mathbb{Q}, \mathbb{P})$ **almost surely**.

What is a random metric?

In short

It is a metric almost surely (assumptions: next slide).

In other words,

- $\rho(\mathbb{P}, \mathbb{Q}) \geq 0$, $\rho(\mathbb{P}, \mathbb{Q}) = 0 \Leftrightarrow \mathbb{P} = \mathbb{Q}$ almost surely.
- $\rho(\mathbb{P}, \mathbb{Q}) = \rho(\mathbb{Q}, \mathbb{P})$ almost surely.
- $\rho(\mathbb{P}, \mathbb{Q}) \leq \rho(\mathbb{P}, \mathbb{D}) + \rho(\mathbb{D}, \mathbb{Q})$ almost surely.

What is a random metric?

In short

It is a metric almost surely (assumptions: next slide).

In other words,

- $\rho(\mathbb{P}, \mathbb{Q}) \geq 0$, $\rho(\mathbb{P}, \mathbb{Q}) = 0 \Leftrightarrow \mathbb{P} = \mathbb{Q}$ almost surely.
- $\rho(\mathbb{P}, \mathbb{Q}) = \rho(\mathbb{Q}, \mathbb{P})$ almost surely.
- $\rho(\mathbb{P}, \mathbb{Q}) \leq \rho(\mathbb{P}, \mathbb{D}) + \rho(\mathbb{D}, \mathbb{Q})$ almost surely.

$\mathcal{V} = \{\mathbf{v}_j\}_{j=1}^J \subset \mathbb{R}^d$: reason of randomness.

Theorem

If $\mathcal{X} \subset \mathbb{R}^d$ is connected open, and k is

- bounded: $\sup_{\mathbf{x}, \mathbf{x}'} k(\mathbf{x}, \mathbf{x}') \leq B_k < \infty$,

Theorem

If $\mathcal{X} \subset \mathbb{R}^d$ is connected open, and k is

- bounded: $\sup_{\mathbf{x}, \mathbf{x}'} k(\mathbf{x}, \mathbf{x}') \leq B_k < \infty$,
- analytic: $\mathbf{x} \mapsto k(\mathbf{x}, \mathbf{y})$ is analytic for any $\mathbf{y} \in \mathbb{R}^d$.

Theorem

If $\mathcal{X} \subset \mathbb{R}^d$ is connected open, and k is

- bounded: $\sup_{\mathbf{x}, \mathbf{x}'} k(\mathbf{x}, \mathbf{x}') \leq B_k < \infty$,
- analytic: $\mathbf{x} \mapsto k(\mathbf{x}, \mathbf{y})$ is analytic for any $\mathbf{y} \in \mathbb{R}^d$.
- characteristic: μ_k is injective,

Theorem

If $\mathcal{X} \subset \mathbb{R}^d$ is connected open, and k is

- bounded: $\sup_{\mathbf{x}, \mathbf{x}'} k(\mathbf{x}, \mathbf{x}') \leq B_k < \infty$,
- analytic: $\mathbf{x} \mapsto k(\mathbf{x}, \mathbf{y})$ is analytic for any $\mathbf{y} \in \mathbb{R}^d$.
- characteristic: μ_k is injective,

then

$$\rho(\mathbb{P}, \mathbb{Q}) := \sqrt{\frac{1}{J} \sum_{j=1}^J [\mu_{\mathbb{P}}(\mathbf{v}_j) - \mu_{\mathbb{Q}}(\mathbf{v}_j)]^2}$$

is a metric a.s. w.r.t. $\{\mathbf{v}_j\}_{j=1}^J$.

Why do analytic features work? – proof idea

- μ is injective and maps to analytic functions:
 - k : bounded, analytic \Rightarrow elements of \mathcal{H}_k : analytic.
 - k : characteristic, bounded $\Rightarrow \mu = \mu_k$: well-defined, injective.

Why do analytic features work? – proof idea

- μ is injective and maps to analytic functions:
 - k : bounded, analytic \Rightarrow elements of \mathcal{H}_k : analytic.
 - k : characteristic, bounded $\Rightarrow \mu = \mu_k$: well-defined, injective.
- μ : characteristic \Rightarrow for $\mathbb{P} \neq \mathbb{Q}$, $f := \mu_{\mathbb{P}} - \mu_{\mathbb{Q}} \neq 0$.

Why do analytic features work? – proof idea

- μ is injective and maps to analytic functions:
 - k : bounded, analytic \Rightarrow elements of \mathcal{H}_k : analytic.
 - k : characteristic, bounded $\Rightarrow \mu = \mu_k$: well-defined, injective.
- μ : characteristic \Rightarrow for $\mathbb{P} \neq \mathbb{Q}$, $f := \mu_{\mathbb{P}} - \mu_{\mathbb{Q}} \neq 0$.
- f : analytic, thus

$$\rho(\mathbb{P}, \mathbb{Q}) = \sqrt{\sum_{j=1}^J [\mu_{\mathbb{P}}(\mathbf{v}_j) - \mu_{\mathbb{Q}}(\mathbf{v}_j)]^2}$$

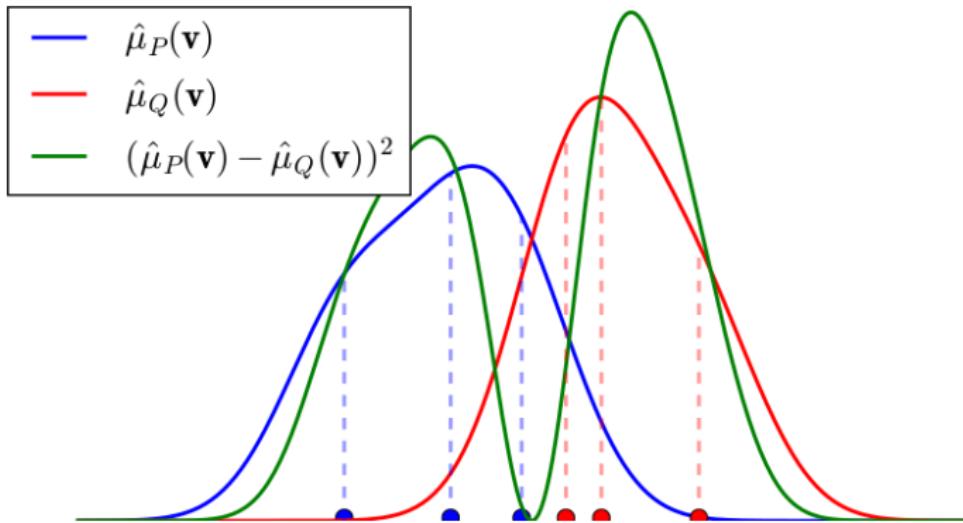
is a metric, a.s. w.r.t. $(\mathbf{v}_j \stackrel{i.i.d.}{\sim} \cdot)$ $m \ll \lambda$. Reason: for an analytic $f \neq 0$, $m\{\mathbf{v} : f(\mathbf{v}) = 0\} = 0$.

Estimation

Compute

$$\hat{\rho}^2(\mathbb{P}, \mathbb{Q}) = \frac{1}{J} \sum_{j=1}^J [\hat{\mu}_{\mathbb{P}}(\mathbf{v}_j) - \hat{\mu}_{\mathbb{Q}}(\mathbf{v}_j)]^2,$$

where $\hat{\mu}_{\mathbb{P}}(\mathbf{v}) = \frac{1}{n} \sum_{i=1}^n k(\mathbf{x}_i, \mathbf{v})$. Example using $k(\mathbf{x}, \mathbf{v}) = e^{-\frac{\|\mathbf{x}-\mathbf{v}\|^2}{2\sigma^2}}$:



Estimation – continued

$$\hat{\rho}^2(\mathbb{P}, \mathbb{Q}) = \frac{1}{J} \sum_{j=1}^J [\hat{\mu}_{\mathbb{P}}(\mathbf{v}_j) - \hat{\mu}_{\mathbb{Q}}(\mathbf{v}_j)]^2$$

Estimation – continued

$$\begin{aligned}\hat{\rho}^2(\mathbb{P}, \mathbb{Q}) &= \frac{1}{J} \sum_{j=1}^J [\hat{\mu}_{\mathbb{P}}(\mathbf{v}_j) - \hat{\mu}_{\mathbb{Q}}(\mathbf{v}_j)]^2 \\ &= \frac{1}{J} \sum_{j=1}^J \left[\frac{1}{n} \sum_{i=1}^n k(\mathbf{x}_i, \mathbf{v}_j) - \frac{1}{n} \sum_{i=1}^n k(\mathbf{y}_i, \mathbf{v}_j) \right]^2\end{aligned}$$

Estimation – continued

$$\begin{aligned}\hat{\rho}^2(\mathbb{P}, \mathbb{Q}) &= \frac{1}{J} \sum_{j=1}^J [\hat{\mu}_{\mathbb{P}}(\mathbf{v}_j) - \hat{\mu}_{\mathbb{Q}}(\mathbf{v}_j)]^2 \\ &= \frac{1}{J} \sum_{j=1}^J \left[\frac{1}{n} \sum_{i=1}^n k(\mathbf{x}_i, \mathbf{v}_j) - \frac{1}{n} \sum_{i=1}^n k(\mathbf{y}_i, \mathbf{v}_j) \right]^2 = \frac{1}{J} \sum_{j=1}^J (\bar{\mathbf{z}}_n)_j^2 = \frac{1}{J} \bar{\mathbf{z}}_n^T \bar{\mathbf{z}}_n,\end{aligned}$$

where $\bar{\mathbf{z}}_n = \frac{1}{n} \sum_{i=1}^n \underbrace{[k(\mathbf{x}_i, \mathbf{v}_j) - k(\mathbf{y}_i, \mathbf{v}_j)]_{j=1}^J}_{=: \mathbf{z}_i} \in \mathbb{R}^J$.

Estimation – continued

$$\begin{aligned}\hat{\rho}^2(\mathbb{P}, \mathbb{Q}) &= \frac{1}{J} \sum_{j=1}^J [\hat{\mu}_{\mathbb{P}}(\mathbf{v}_j) - \hat{\mu}_{\mathbb{Q}}(\mathbf{v}_j)]^2 \\ &= \frac{1}{J} \sum_{j=1}^J \left[\frac{1}{n} \sum_{i=1}^n k(\mathbf{x}_i, \mathbf{v}_j) - \frac{1}{n} \sum_{i=1}^n k(\mathbf{y}_i, \mathbf{v}_j) \right]^2 = \frac{1}{J} \sum_{j=1}^J (\bar{\mathbf{z}}_n)_j^2 = \frac{1}{J} \bar{\mathbf{z}}_n^T \bar{\mathbf{z}}_n,\end{aligned}$$

where $\bar{\mathbf{z}}_n = \frac{1}{n} \sum_{i=1}^n \underbrace{[k(\mathbf{x}_i, \mathbf{v}_j) - k(\mathbf{y}_i, \mathbf{v}_j)]_{j=1}^J}_{=: \mathbf{z}_i} \in \mathbb{R}^J$.

- Good news: estimation is linear in n !
- Bad news: intractable null distr. $= \sqrt{n} \hat{\rho}^2(\mathbb{P}, \mathbb{P}) \xrightarrow{d}$ sum of J correlated χ^2 .

- Modified test statistic:

$$\hat{\lambda}_n = n\bar{\mathbf{z}}_n^T \boldsymbol{\Sigma}_n^{-1} \bar{\mathbf{z}}_n,$$

where $\boldsymbol{\Sigma}_n = \text{cov}(\{\mathbf{z}_i\}_{i=1}^n)$.

- Under H_0 :
 - $\hat{\lambda}_n \xrightarrow{d} \chi^2(J)$. \Rightarrow Easy to get the $(1 - \alpha)$ -quantile!

- Characteristic functions – 'poor' choice:

$$\rho_2(\mathbb{P}, \mathbb{Q}) := \sqrt{\frac{1}{J} \sum_{j=1}^J [\phi_{\mathbb{P}}(\mathbf{v}_j) - \phi_{\mathbb{Q}}(\mathbf{v}_j)]^2}.$$

[Chwialkowski et al., 2015, Prop. 1]: It fails to distinguish a large class of distributions.

- Characteristic functions – 'poor' choice:

$$\rho_2(\mathbb{P}, \mathbb{Q}) := \sqrt{\frac{1}{J} \sum_{j=1}^J [\phi_{\mathbb{P}}(\mathbf{v}_j) - \phi_{\mathbb{Q}}(\mathbf{v}_j)]^2}.$$

[Chwialkowski et al., 2015, Prop. 1]: It fails to distinguish a large class of distributions.

- [Moulines et al., 2007]:

$$\rho_3(\mathbb{P}, \mathbb{Q}) := \frac{n_x n_y}{n} \left\| C^{-\frac{1}{2}} (\mu_{\mathbb{Q}} - \mu_{\mathbb{P}}) \right\|_{\mathcal{H}_k},$$

$$C = \frac{n_x}{n_x + n_y} C_{xx} + \frac{n_y}{n_x + n_y} C_{yy} : \text{pooled covariance operator.}$$

- Characteristic functions – 'poor' choice:

$$\rho_2(\mathbb{P}, \mathbb{Q}) := \sqrt{\frac{1}{J} \sum_{j=1}^J [\phi_{\mathbb{P}}(\mathbf{v}_j) - \phi_{\mathbb{Q}}(\mathbf{v}_j)]^2}.$$

[Chwialkowski et al., 2015, Prop. 1]: It fails to distinguish a large class of distributions.

- [Moulines et al., 2007]:

$$\rho_3(\mathbb{P}, \mathbb{Q}) := \frac{n_x n_y}{n} \left\| C^{-\frac{1}{2}} (\mu_{\mathbb{Q}} - \mu_{\mathbb{P}}) \right\|_{\mathcal{H}_k},$$

$$C = \frac{n_x}{n_x + n_y} C_{xx} + \frac{n_y}{n_x + n_y} C_{yy} : \text{pooled covariance operator.}$$

Computational cost: **high** (cubic).

- Until now: spatial domain.
- Smoothed characteristic functions:

$$\psi_{\mathbb{P}}(\mathbf{t}) = \int_{\mathbb{R}^d} \phi_{\mathbb{P}}(\boldsymbol{\omega}) \ell(\mathbf{t} - \boldsymbol{\omega}) d\boldsymbol{\omega}, \quad \mathbf{t} \in \mathbb{R}^d,$$

$$\rho_4(\mathbb{P}, \mathbb{Q}) := \sqrt{\frac{1}{J} \sum_{j=1}^J [\psi_{\mathbb{P}}(\mathbf{v}_j) - \psi_{\mathbb{Q}}(\mathbf{v}_j)]^2}.$$

- Until now: spatial domain.
- Smoothed characteristic functions:

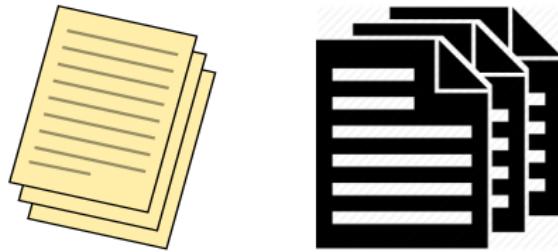
$$\psi_{\mathbb{P}}(\mathbf{t}) = \int_{\mathbb{R}^d} \phi_{\mathbb{P}}(\boldsymbol{\omega}) \ell(\mathbf{t} - \boldsymbol{\omega}) d\boldsymbol{\omega}, \quad \mathbf{t} \in \mathbb{R}^d,$$
$$\rho_4(\mathbb{P}, \mathbb{Q}) := \sqrt{\frac{1}{J} \sum_{j=1}^J [\psi_{\mathbb{P}}(\mathbf{v}_j) - \psi_{\mathbb{Q}}(\mathbf{v}_j)]^2}.$$

- Notes:
 - For analytic smoothing kernels (ℓ), it works.
 - It is more sensitive to differences in the frequency domain.

Linear-time **high-power** two-sample testing

Example-1: NLP

- Given: two categories of documents (Bayesian inference, neuroscience).
- Task:
 - test their distinguishability,
 - most discriminative words → interpretability.



Example-2: computer vision



- Given: two sets of faces (happy, angry).
- Task:
 - check if they are different,
 - determine the most discriminative features/regions.

- We get a nonparametric t-test.
- It gives a reason why H_0 is rejected.
- It is
 - adaptive → high test power.
 - fast (linear time).

- We get a nonparametric t-test.
- It gives a reason why H_0 is rejected.
- It is
 - adaptive → high test power.
 - fast (linear time).

Code:

- <https://github.com/wittawatj/interpretable-test>

- Until this point: test locations (\mathcal{V}) are **fixed**.
- Instead: choose $\theta = \{\mathcal{V}, \sigma\}$ to
maximize lower bound on the test power.

- Until this point: test locations (\mathcal{V}) are **fixed**.
- Instead: choose $\theta = \{\mathcal{V}, \sigma\}$ to
maximize lower bound on the test power.

Theorem (Lower bound on power, for large n)

Test power $\geq L(\lambda_n)$; L : explicit function, monotonically increasing.

- Here,
 - $\lambda_n = n\mu^T \Sigma^{-1} \mu$: population version of $\hat{\lambda}_n = n\bar{z}_n^T \Sigma_n^{-1} \bar{z}_n$.
 - $\mu = \mathbb{E}_{xy}[z_1]$, $\Sigma = \mathbb{E}_{xy}[(z_1 - \mu)(z_1 - \mu)^T]$.

Convergence of the λ_n estimator

But λ_n is **unknown**. \Rightarrow Split (X, Y) into (X_{tr}, Y_{tr}) and (X_{te}, Y_{te}) .

- Locations, kernel parameter: $\hat{\theta} = \arg \max_{\theta} \hat{\lambda}_{\frac{n}{2}}^{tr}(\theta)$.

Convergence of the λ_n estimator

But λ_n is **unknown**. \Rightarrow Split (X, Y) into (X_{tr}, Y_{tr}) and (X_{te}, Y_{te}) .

- Locations, kernel parameter: $\hat{\theta} = \arg \max_{\theta} \hat{\lambda}_{\frac{n}{2}}^{tr}(\theta)$.
- Test statistic: $\hat{\lambda}_{\frac{n}{2}}^{te}(\hat{\theta})$.

Convergence of the λ_n estimator

Theorem (Guarantee on objective approximation, $\gamma_n \rightarrow 0$)

$$\sup_{\mathcal{V}, \mathcal{K}} |\bar{\mathbf{z}}_n^T (\boldsymbol{\Sigma}_n + \gamma_n)^{-1} \bar{\mathbf{z}}_n - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}| = \mathcal{O}(n^{-\frac{1}{4}}).$$

Convergence of the λ_n estimator

Theorem (Guarantee on objective approximation, $\gamma_n \rightarrow 0$)

$$\sup_{\mathcal{V}, \mathcal{K}} |\bar{\mathbf{z}}_n^T (\boldsymbol{\Sigma}_n + \gamma_n)^{-1} \bar{\mathbf{z}}_n - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}| = \mathcal{O}(n^{-\frac{1}{4}}).$$

Examples:

$$\mathcal{K} = \left\{ k_\sigma(\mathbf{x}, \mathbf{y}) = e^{-\frac{\|\mathbf{x}-\mathbf{y}\|^2}{2\sigma^2}} : \sigma > 0 \right\},$$

$$\mathcal{K} = \left\{ k_{\mathbf{A}}(\mathbf{x}, \mathbf{y}) = e^{-(\mathbf{x}-\mathbf{y})^T \mathbf{A} (\mathbf{x}-\mathbf{y})} : \mathbf{A} > 0 \right\}.$$

- Lower bound on the test power:
 - $|\hat{\lambda}_n - \lambda_n| \lesssim \|\bar{\mathbf{z}}_n - \boldsymbol{\mu}\|_2 + \|\boldsymbol{\Sigma}_n - \boldsymbol{\Sigma}\|_F$.
 - Bound the r.h.s. by Hoeffding inequality $\Rightarrow P(|\hat{\lambda}_n - \lambda_n| \geq t)$.
 - By reparameterization: $P(\hat{\lambda}_n \geq T_\alpha)$ bound.

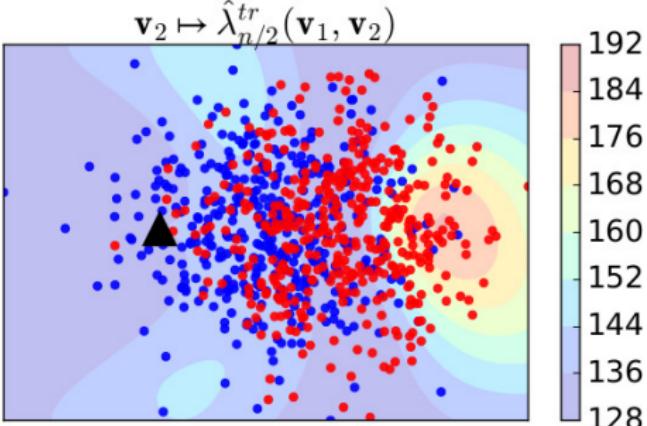
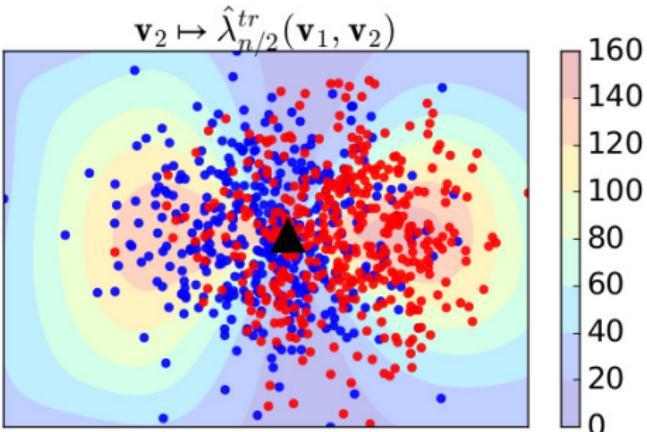
- Lower bound on the test power:
 - $|\hat{\lambda}_n - \lambda_n| \lesssim \|\bar{\mathbf{z}}_n - \boldsymbol{\mu}\|_2 + \|\boldsymbol{\Sigma}_n - \boldsymbol{\Sigma}\|_F$.
 - Bound the r.h.s. by Hoeffding inequality $\Rightarrow P(|\hat{\lambda}_n - \lambda_n| \geq t)$.
 - By reparameterization: $P(\hat{\lambda}_n \geq T_\alpha)$ bound.
- Uniformly $\hat{\lambda}_n \approx \lambda_n$:
 - Reduction to bounding $\sup_{\mathcal{V}, \mathcal{S}} \|\bar{\mathbf{z}}_n - \boldsymbol{\mu}\|_2, \sup_{\mathcal{V}, \mathcal{S}} \|\boldsymbol{\Sigma}_n - \boldsymbol{\Sigma}\|_F$.
 - Empirical processes, Dudley entropy bound.

Non-convexity, informative features

- 2D problem:

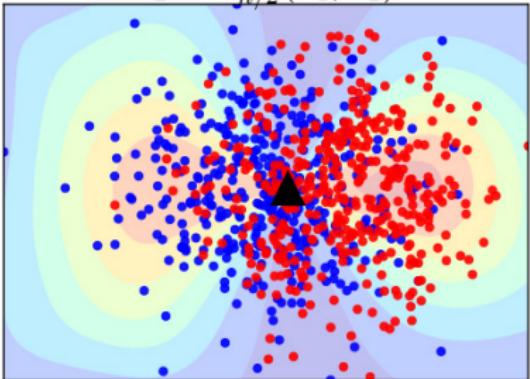
$$\mathbb{P} := \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad \mathbb{Q} := \mathcal{N}(\mathbf{e}_1, \mathbf{I}).$$

- $\mathcal{V} = \{\mathbf{v}_1, \mathbf{v}_2\}$. Fix \mathbf{v}_1 to the triangle.
- $\mathbf{v}_2 \mapsto \hat{\lambda}_n(\{\mathbf{v}_1, \mathbf{v}_2\})$: contour plot.

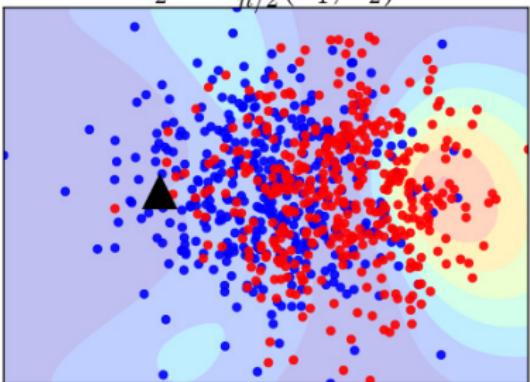


Non-convexity, informative features

$$\mathbf{v}_2 \mapsto \hat{\lambda}_{n/2}^{tr}(\mathbf{v}_1, \mathbf{v}_2)$$



$$\mathbf{v}_2 \mapsto \hat{\lambda}_{n/2}^{tr}(\mathbf{v}_1, \mathbf{v}_2)$$



- **Nearby locations:** do not increase discriminability.
- **Non-convexity:** reveals multiple ways to capture the difference.

Computational complexity

- Optimization & testing: linear in n .
- Testing: $\mathcal{O}(ndJ + nJ^2 + J^3)$.
- Optimization: $\mathcal{O}(ndJ^2 + J^3)$ per gradient ascent.

- Small J :

- often enough to detect the difference of \mathbb{P} & \mathbb{Q} .
- few distinguishing regions to reject H_0 .
- faster test.

Number of locations (J)

- Very large J :
 - test power need not increase monotonically in J (more locations \Rightarrow statistic can gain in variance).
 - defeats the purpose of a linear-time test.

Numerical demos

Parameter settings

- Gaussian kernel (σ). $\alpha = 0.01$. $J = 1$. Repeat 500 trials.
- Report

$$P(\text{reject } H_0) \approx \frac{\#\text{times } \hat{\lambda}_n > T_\alpha \text{ holds}}{\#\text{trials}}.$$

- Compare 4 methods
 - **ME-full**: Optimize \mathcal{V} and Gaussian bandwidth σ .
 - **ME-grid**: Optimize σ . Random \mathcal{V} [Chwialkowski et al., 2015].
 - **MMD-quad**: Test with quadratic-time MMD [Gretton et al., 2012].
 - **MMD-lin**: Test with linear-time MMD [Gretton et al., 2012].
- Optimize kernels to power in MMD-lin, MMD-quad.

NLP: discrimination of document categories

- 5903 NIPS papers (1988-2015).
- Keyword-based category assignment into 4 groups:
 - Bayesian inference, Deep learning, Learning theory, Neuroscience
- $d = 2000$ nouns. TF-IDF representation.

Problem	n^{te}	ME-full	ME-grid	MMD-quad	MMD-lin
1. Bayes-Bayes	215	.012	.018	.022	.008
2. Bayes-Deep	216	.954	.034	.906	.262
3. Bayes-Learn	138	.990	.774	1.00	.238
4. Bayes-Neuro	394	1.00	.300	.952	.972
5. Learn-Deep	149	.956	.052	.876	.500
6. Learn-Neuro	146	.960	.572	1.00	.538

- Performance of ME-full [$\mathcal{O}(n)$] is comparable to MMD-quad [$\mathcal{O}(n^2)$].

- Aggregating over trials; example: 'Bayes-Neuro'.
- Most discriminative words:
spike, markov, cortex, dropout, recur, iii, gibb.
 - learned test locations: highly interpretable,
 - '**markov**', '**gibb**' (\Leftarrow Gibbs): **Bayes**ian inference,
 - '**spike**', '**cortexneuroscience**.

- Aggregating over trials; example: 'Bayes-Neuro'.
- Least discriminatory ones:
circumfer, bra, dominiqu, rhino, mitra, kid, impostor.

Distinguish positive/negative emotions

- Karolinska Directed Emotional Faces (KDEF) [Lundqvist et al., 1998].
- 70 actors = 35 females and 35 males.
- $d = 48 \times 34 = 1632$. Grayscale. Pixel features.



Problem	n^{te}	ME-full	ME-grid	MMD-quad	MMD-lin
\pm vs. \pm	201	.010	.012	.018	.008
$+$ vs. $-$	201	.998	.656	1.00	.578

- Learned test location (averaged) = 

Linear-time high-power two-sample testing:
finished

Linear-time **high-power** independence testing

Example: dependency testing of media annotations

- We are given **paired samples**. Task: test **independence**.
- Examples:
 - (song, year of release) pairs

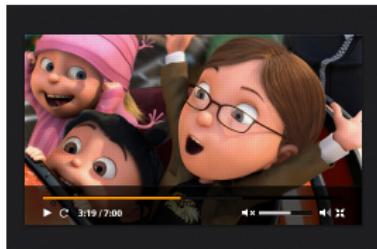


Example: dependency testing of media annotations

- We are given **paired samples**. Task: test **independence**.
- Examples:
 - (song, year of release) pairs



- (video, caption) pairs

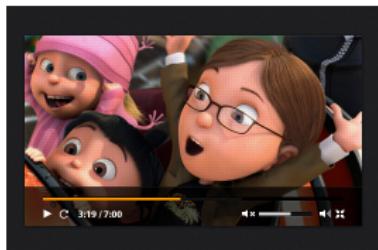


Example: dependency testing of media annotations

- We are given **paired samples**. Task: test **independence**.
- Examples:
 - (song, year of release) pairs



- (video, caption) pairs



- $\{(x_i, y_i)\}_{i=1}^n \xrightarrow{?} H_0 : \mathbb{P}_{xy} = \mathbb{P}_x \mathbb{P}_y, H_1 : \mathbb{P}_{xy} \neq \mathbb{P}_x \mathbb{P}_y.$

2-sample test → independence test

Until now:

- adaptive linear-time 2-sample test (automatic parameter tuning).

2-sample test → independence test

2-sample test:

$$\text{MMD}(\mathbb{P}, \mathbb{Q}) = \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k}, \quad \rho(\mathbb{P}, \mathbb{Q}) = \sqrt{\frac{1}{J} \sum_{j=1}^J [\mu_{\mathbb{P}}(\mathbf{v}_j) - \mu_{\mathbb{Q}}(\mathbf{v}_j)]^2},$$

2-sample test → independence test

2-sample test:

$$\text{MMD}(\mathbb{P}, \mathbb{Q}) = \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k}, \quad \rho(\mathbb{P}, \mathbb{Q}) = \sqrt{\frac{1}{J} \sum_{j=1}^J [\mu_{\mathbb{P}}(\mathbf{v}_j) - \mu_{\mathbb{Q}}(\mathbf{v}_j)]^2},$$

Independence test [Jitkrittum et al., 2016b]:

$$\text{HSIC}(x, y) = \|\mu_{xy} - \mu_x \otimes \mu_y\|_{\mathcal{H}_k \otimes \mathcal{H}_\ell}, \quad \text{FSIC}(\mathbf{x}, \mathbf{y}) = \sqrt{\frac{1}{J} \sum_{j=1}^J u^2(\mathbf{v}_j, \mathbf{w}_j)}$$

2-sample test → independence test

2-sample test:

$$\text{MMD}(\mathbb{P}, \mathbb{Q}) = \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k}, \quad \rho(\mathbb{P}, \mathbb{Q}) = \sqrt{\frac{1}{J} \sum_{j=1}^J [\mu_{\mathbb{P}}(\mathbf{v}_j) - \mu_{\mathbb{Q}}(\mathbf{v}_j)]^2},$$

Independence test [Jitkrittum et al., 2016b]:

$$\text{HSIC}(x, y) = \|\mu_{xy} - \mu_x \otimes \mu_y\|_{\mathcal{H}_k \otimes \mathcal{H}_\ell}, \quad \text{FSIC}(\mathbf{x}, \mathbf{y}) = \sqrt{\frac{1}{J} \sum_{j=1}^J u^2(\mathbf{v}_j, \mathbf{w}_j)},$$

with $u(\mathbf{v}, \mathbf{w}) = \mu_{xy}(\mathbf{v}, \mathbf{w}) - \mu_x(\mathbf{v})\mu_y(\mathbf{w})$ witness function.

FSIC: covariance view

(\mathbf{v}, \mathbf{w}) : fixed. By rewriting

$$\begin{aligned} u(\mathbf{v}, \mathbf{w}) &= \mu_{\mathbf{x}\mathbf{y}}(\mathbf{v}, \mathbf{w}) - \mu_{\mathbf{x}}(\mathbf{v})\mu_{\mathbf{y}}(\mathbf{w}) \\ &= \mathbb{E}_{\mathbf{x}\mathbf{y}}[k(\mathbf{x}, \mathbf{v})\ell(\mathbf{y}, \mathbf{w})] - \mathbb{E}_{\mathbf{x}}[k(\mathbf{x}, \mathbf{v})]\mathbb{E}_{\mathbf{y}}[\ell(\mathbf{y}, \mathbf{w})] \\ &= \text{cov}_{\mathbf{x}\mathbf{y}}(k(\mathbf{x}, \mathbf{v}), \ell(\mathbf{y}, \mathbf{w})). \end{aligned}$$

\Rightarrow We picked the $(\mathbf{v}, \mathbf{w})^{th}$ entry of

$$\begin{aligned} C_{\mathbf{x}\mathbf{y}}^c &= \mathbb{E}_{\mathbf{x}\mathbf{y}} [\varphi(\mathbf{x}) \otimes \psi(\mathbf{y})] - \mu_{\mathbf{x}} \otimes \mu_{\mathbf{y}}, \\ \text{HSIC} &= \|C_{\mathbf{x}\mathbf{y}}^c\|_{HS}. \end{aligned}$$

FSIC is an independence measure

Theorem

If $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ are bounded, characteristic, analytic kernels [$\mathcal{X} \subseteq \mathbb{R}^{d_x}$, $\mathcal{Y} \subseteq \mathbb{R}^{d_y}$: connected open], then almost surely

$$\text{FSIC}(\mathbf{x}, \mathbf{y}) = 0 \Leftrightarrow \mathbf{x} \perp \mathbf{y}.$$

FSIC is an independence measure

Theorem

If $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ are bounded, characteristic, analytic kernels [$\mathcal{X} \subseteq \mathbb{R}^{d_x}$, $\mathcal{Y} \subseteq \mathbb{R}^{d_y}$: connected open], then almost surely

$$\text{FSIC}(\mathbf{x}, \mathbf{y}) = 0 \Leftrightarrow \mathbf{x} \perp \mathbf{y}.$$

Consequence

FSIC can be applied in ISA, feature selection, outlier-robust image registration, ...

Empirical estimator for FSIC

$$\text{FSIC}^2(\mathbf{x}, \mathbf{y}) = \frac{1}{J} \sum_{j=1}^J u^2(\mathbf{v}_j, \mathbf{w}_j), \quad u(\mathbf{v}, \mathbf{w}) = \mu_{xy}(\mathbf{v}, \mathbf{w}) - \mu_x(\mathbf{v})\mu_y(\mathbf{w}),$$

$$\begin{aligned}\widehat{\text{FSIC}}^2(\mathbf{x}, \mathbf{y}) &= \frac{1}{J} \sum_{j=1}^J \hat{u}^2(\mathbf{v}_j, \mathbf{w}_j), \quad \hat{u}(\mathbf{v}, \mathbf{w}) = \widehat{\mu_{xy}}(\mathbf{v}, \mathbf{w}) - (\widehat{\mu_x \mu_y})(\mathbf{v}, \mathbf{w}), \\ &= \frac{1}{J} \|\mathbf{u}\|_2^2\end{aligned}$$

Empirical estimator for FSIC

$$\text{FSIC}^2(\mathbf{x}, \mathbf{y}) = \frac{1}{J} \sum_{j=1}^J u^2(\mathbf{v}_j, \mathbf{w}_j), \quad u(\mathbf{v}, \mathbf{w}) = \mu_{xy}(\mathbf{v}, \mathbf{w}) - \mu_x(\mathbf{v})\mu_y(\mathbf{w}),$$

$$\begin{aligned}\widehat{\text{FSIC}}^2(\mathbf{x}, \mathbf{y}) &= \frac{1}{J} \sum_{j=1}^J \hat{u}^2(\mathbf{v}_j, \mathbf{w}_j), \quad \hat{u}(\mathbf{v}, \mathbf{w}) = \widehat{\mu_{xy}}(\mathbf{v}, \mathbf{w}) - (\widehat{\mu_x \mu_y})(\mathbf{v}, \mathbf{w}), \\ &= \frac{1}{J} \|\mathbf{u}\|_2^2,\end{aligned}$$

where we use the unbiased estimators [2nd = ' $\mu_x(\mathbf{v})\mu_y(\mathbf{w})$ - diag']:

$$\widehat{\mu_{xy}}(\mathbf{v}, \mathbf{w}) = \frac{1}{n} \sum_{i=1}^n k(\mathbf{x}_i, \mathbf{v})\ell(\mathbf{y}_i, \mathbf{w}),$$

$$\widehat{\mu_x \mu_y}(\mathbf{v}, \mathbf{w}) = \frac{1}{n(n-1)} \sum_{i \neq j} k(\mathbf{x}_i, \mathbf{v})\ell(\mathbf{y}_j, \mathbf{w}).$$

Asymptotic distribution of $\hat{\mathbf{u}}$

For fixed (\mathbf{v}, \mathbf{w}) :

$$\hat{u}(\mathbf{v}, \mathbf{w}) = \frac{2}{n(n-1)} \sum_{i < j} h_{\mathbf{v}, \mathbf{w}} ((\mathbf{x}_i, \mathbf{y}_i), (\mathbf{x}_j, \mathbf{y}_j)),$$

$$h_{\mathbf{v}, \mathbf{w}} ((\mathbf{x}, \mathbf{y}), (\mathbf{x}', \mathbf{y}')) = \frac{1}{2} [k(\mathbf{x}, \mathbf{v}) - k(\mathbf{x}', \mathbf{v})] [\ell(\mathbf{y}, \mathbf{w}) - \ell(\mathbf{y}', \mathbf{w})]$$

Asymptotic distribution of $\hat{\mathbf{u}}$

For fixed (\mathbf{v}, \mathbf{w}) :

$$\hat{u}(\mathbf{v}, \mathbf{w}) = \frac{2}{n(n-1)} \sum_{i < j} h_{\mathbf{v}, \mathbf{w}} ((\mathbf{x}_i, \mathbf{y}_i), (\mathbf{x}_j, \mathbf{y}_j)),$$

$$h_{\mathbf{v}, \mathbf{w}} ((\mathbf{x}, \mathbf{y}), (\mathbf{x}', \mathbf{y}')) = \frac{1}{2} [k(\mathbf{x}, \mathbf{v}) - k(\mathbf{x}', \mathbf{v})] [\ell(\mathbf{y}, \mathbf{w}) - \ell(\mathbf{y}', \mathbf{w})],$$

thus $\xrightarrow{\text{theory of U-statistics}}$

Theorem (Asymptotic normality)

For any fixed locations $\mathcal{V} = \{(\mathbf{v}_j, \mathbf{w}_j)\}_{j=1}^J$, $\hat{\mathbf{u}} := [\hat{u}(\mathbf{v}_j, \mathbf{w}_j)]_{j=1}^J$

$$\sqrt{n} (\hat{\mathbf{u}} - \mathbf{u}) \xrightarrow{d} N(\mathbf{0}, \boldsymbol{\Sigma}),$$

$$\Sigma_{ij} = cov_{\mathbf{xy}} (\hat{u}(\mathbf{v}_i, \mathbf{w}_i), \hat{u}(\mathbf{v}_j, \mathbf{w}_j)).$$

$$\text{NFSIC} = \text{FSIC} + \text{whitening}$$

- $\widehat{n\text{FSIC}}^2(x, y) = n \frac{\|\mathbf{u}\|_2^2}{J}$: asymptotically **sum of correlated χ^2 -s.**

$$\text{NFSIC} = \text{FSIC} + \text{whitening}$$

- $n\widehat{\text{FSIC}}^2(x, y) = n\frac{\|\mathbf{u}\|_2^2}{J}$: asymptotically **sum of correlated χ^2 -s**.
- Quantile: **hard**. \Rightarrow With the **whitening** trick:

Theorem

- Under H_0 : with $\gamma_n \rightarrow 0$

$$\hat{\lambda}_n = n\hat{\mathbf{u}}^T \left(\hat{\Sigma}_n + \gamma_n \mathbf{I}_J \right)^{-1} \hat{\mathbf{u}} \xrightarrow{d} \chi^2(J).$$

$$\text{NFSIC} = \text{FSIC} + \text{whitening}$$

- $n\widehat{\text{FSIC}}^2(x, y) = n\frac{\|\mathbf{u}\|_2^2}{J}$: asymptotically sum of correlated χ^2 -s.
- Quantile: hard. \Rightarrow With the whitening trick:

Theorem

- Under H_0 : with $\gamma_n \rightarrow 0$

$$\hat{\lambda}_n = n\hat{\mathbf{u}}^T \left(\hat{\Sigma}_n + \gamma_n \mathbf{I}_J \right)^{-1} \hat{\mathbf{u}} \xrightarrow{d} \chi^2(J).$$

- Under H_1 : we get a consistent test (i.e., power $\rightarrow 1$).

NFSIC can be estimated **easily**

Test statistic:

$$\hat{\lambda}_n = n\hat{\mathbf{u}}^T \left(\hat{\Sigma}_n + \gamma_n \mathbf{I}_J \right)^{-1} \hat{\mathbf{u}}.$$

Estimator: **no $n \times n$ Gram matrix**

- $\mathbf{K} := [k(\mathbf{v}_i, \mathbf{x}_j)] \in \mathbb{R}^{J \times n}$, $\mathbf{L} := [\ell(\mathbf{w}_i, \mathbf{y}_j)] \in \mathbb{R}^{J \times n}$,
- $\hat{\Sigma}_n = \frac{\Gamma \Gamma^T}{n}$, $\Gamma = (\mathbf{K} \mathbf{H}_n) \circ (\mathbf{L} \mathbf{H}_n) - \hat{\mathbf{u}} \mathbf{1}_n^T$, $\hat{\mathbf{u}} := \frac{(\mathbf{K} \mathbf{1}_n) \circ (\mathbf{L} \mathbf{1}_n)}{n-1} - \frac{(\mathbf{K} \mathbf{1}_n) \circ (\mathbf{L} \mathbf{1}_n)}{n(n-1)}$.

Computational time:

$$\mathcal{O}(J^3 + J^2 \textcolor{blue}{n} + (d_x + d_y) J \textcolor{blue}{n}) .$$

NFSIC can be estimated **easily**

Test statistic:

$$\hat{\lambda}_n = n\hat{\mathbf{u}}^T \left(\hat{\Sigma}_n + \gamma_n \mathbf{I}_J \right)^{-1} \hat{\mathbf{u}}.$$

Estimator: **no $n \times n$ Gram matrix**

- $\mathbf{K} := [k(\mathbf{v}_i, \mathbf{x}_j)] \in \mathbb{R}^{J \times n}$, $\mathbf{L} := [\ell(\mathbf{w}_i, \mathbf{y}_j)] \in \mathbb{R}^{J \times n}$,
- $\hat{\Sigma}_n = \frac{\Gamma \Gamma^T}{n}$, $\Gamma = (\mathbf{K} \mathbf{H}_n) \circ (\mathbf{L} \mathbf{H}_n) - \hat{\mathbf{u}} \mathbf{1}_n^T$, $\hat{\mathbf{u}} := \frac{(\mathbf{K} \mathbf{1}_n) \mathbf{1}_n^T}{n-1} - \frac{(\mathbf{K} \mathbf{1}_n) \circ (\mathbf{L} \mathbf{1}_n)}{n(n-1)}$.

Computational time:

$$\mathcal{O}(J^3 + J^2 \textcolor{blue}{n} + (d_x + d_y) J \textcolor{blue}{n}).$$

Code with demos:

<https://github.com/wittawatj/fsic-test>

Choosing the locations & kernel parameters

- Consistent test: for $\forall \mathcal{V} = \{(\mathbf{v}_j, \mathbf{w}_j)\}_{j=1}^J$ and kernel parameters.

Choosing the locations & kernel parameters

- Consistent test: for $\forall \mathcal{V} = \{(\mathbf{v}_j, \mathbf{w}_j)\}_{j=1}^J$ and kernel parameters.
- Choose the **power proxy maximizer**.

Theorem

Let $\text{NFSIC}^2(x, y) = \lambda_n = n\mathbf{u}^T \boldsymbol{\Sigma}^{-1} \mathbf{u}$. For large n ,
test power $\geq L(\lambda_n)$,

L : monotonically increasing.

Choosing the locations & kernel parameters

- Consistent test: for $\forall \mathcal{V} = \{(\mathbf{v}_j, \mathbf{w}_j)\}_{j=1}^J$ and kernel parameters.
- Choose the **power proxy maximizer**.

Theorem

Let $\text{NFSIC}^2(x, y) = \lambda_n = n\mathbf{u}^T \boldsymbol{\Sigma}^{-1} \mathbf{u}$. For large n ,
test power $\geq L(\lambda_n)$,

L : monotonically increasing.

- In practice: data-splitting (a la 2-sample testing).

Question

Which one to choose?

- HSIC = $\|u\|_{\mathcal{H}_k \otimes \mathcal{H}_\ell}$.
- FSIC = $\|u\|_{L^2(\mathcal{V})}$, $\mathcal{V} = \{(\mathbf{v}_j, \mathbf{w}_j)\}_{j=1}^J$.

Question

Which one to choose?

- HSIC = $\|u\|_{\mathcal{H}_k \otimes \mathcal{H}_\ell}$.
 - When $p_{xy} - p_x p_y$ is diffuse, close to flat.
- FSIC = $\|u\|_{L^2(\mathcal{V})}$, $\mathcal{V} = \{(\mathbf{v}_j, \mathbf{w}_j)\}_{j=1}^J$.

Question

Which one to choose?

- HSIC = $\|u\|_{\mathcal{H}_k \otimes \mathcal{H}_\ell}$.
 - When $p_{xy} - p_x p_y$ is **diffuse**, close to flat.
- FSIC = $\|u\|_{L^2(\mathcal{V})}$, $\mathcal{V} = \{(\mathbf{v}_j, \mathbf{w}_j)\}_{j=1}^J$.
 - When $p_{xy} - p_x p_y$ is local, with **many peaks**.

Demo settings

- k, ℓ : Gaussian. $J = 10$.
- Report: rejection rate of H_0 .
- Compare 6 methods:

Method	Description	Tuning	Test size	Complexity
NFSIC-opt	Studied	Gradient descent	$n/2$	$\mathcal{O}(n)$
NFSIC-med	No tuning	Random locations	n	$\mathcal{O}(n)$
QHSIC	Full HSIC	Median heuristic	n	$\mathcal{O}(n^2)$
NyHSIC	Nyström + HSIC	Median heuristic	n	$\mathcal{O}(n)$
FHSIC	RFF + HSIC	Median heuristic	n	$\mathcal{O}(n)$
RDC	RFF + CCA	Median heuristic	n	$\mathcal{O}(n \log n)$

Demo-1: million song data

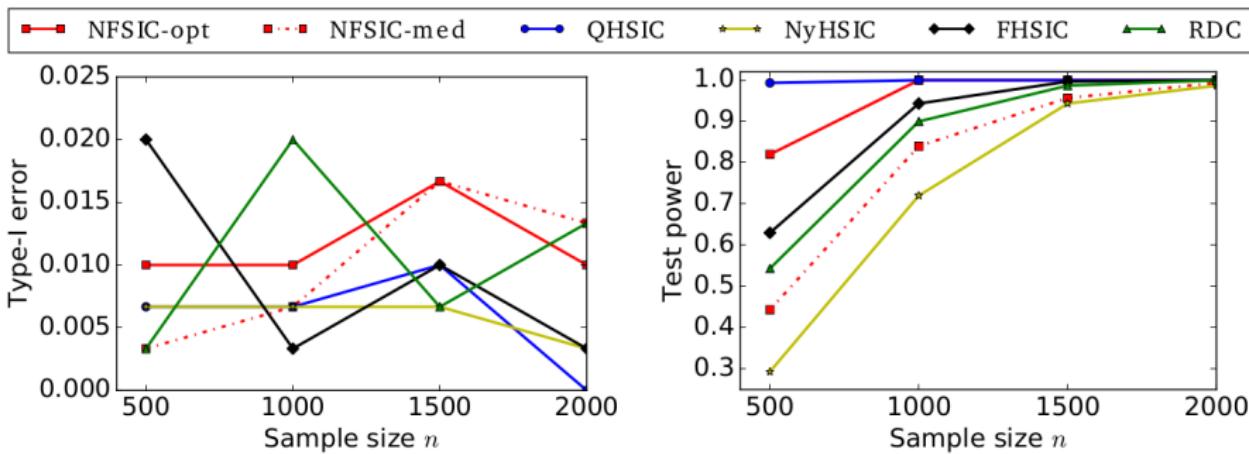
(Song, year of release) =: (\mathbf{x}, y).

- Western commercial tracks from 1922 to 2011 [Bertin-Mahieux et al., 2011].
- $\mathbf{x} \in \mathbb{R}^{90=d_x}$: audio features.
- **Left**: break (\mathbf{x}, y) pairs, i.e. H_0 holds; **right**: H_1 is true.

Demo-1: million song data

(Song, year of release) =: (\mathbf{x}, y) .

- Western commercial tracks from 1922 to 2011 [Bertin-Mahieux et al., 2011].
- $\mathbf{x} \in \mathbb{R}^{90=d_x}$: audio features.
- **Left:** break (\mathbf{x}, y) pairs, i.e. H_0 holds; **right:** H_1 is true.



Demo-2: videos and captions

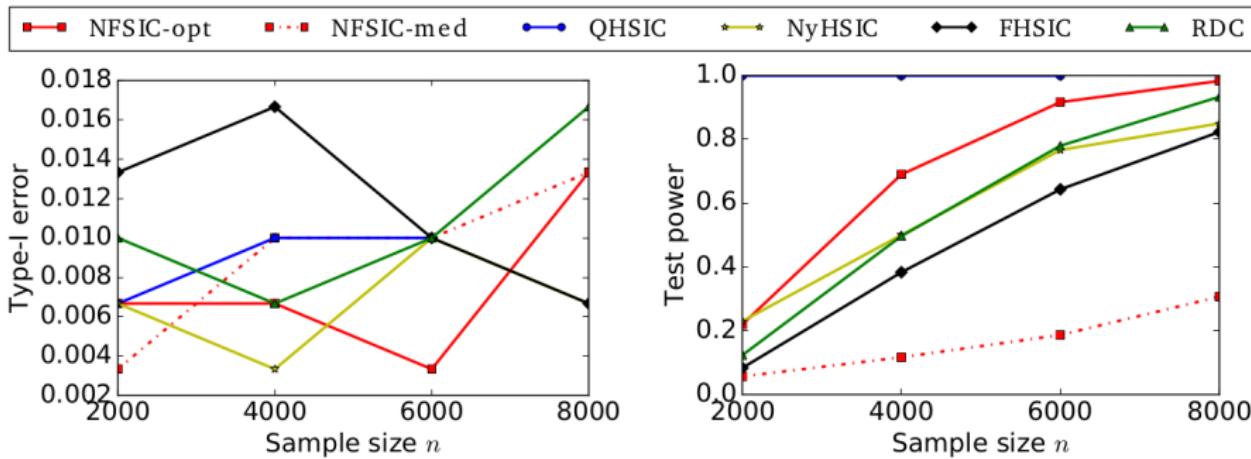
(Youtube video, caption) =: (\mathbf{x}, \mathbf{y}) .

- VideoStory46K [Habibian et al., 2014]
- $\mathbf{x} \in \mathbb{R}^{2000=d_x}$: Fisher vector encoding of motion boundary histograms [Wang and Schmid, 2013].
- $\mathbf{y} \in \mathbb{R}^{1878=d_y}$: bag of words. TF.
- **Left**: break (\mathbf{x}, \mathbf{y}) pairs, i.e. H_0 holds; **right**: H_1 is true.

Demo-2: videos and captions

(Youtube video, caption) =: (\mathbf{x}, \mathbf{y}) .

- VideoStory46K [Habibian et al., 2014]
- $\mathbf{x} \in \mathbb{R}^{2000=d_x}$: Fisher vector encoding of motion boundary histograms [Wang and Schmid, 2013].
- $\mathbf{y} \in \mathbb{R}^{1878=d_y}$: bag of words. TF.
- **Left**: break (\mathbf{x}, \mathbf{y}) pairs, i.e. H_0 holds; **right**: H_1 is true.



Linear-time goodness-of-fit testing: Stein operator & analytical kernels [Jitkrittum et al., 2017]

Given:

- Density/model: p .

Linear-time goodness-of-fit testing: Stein operator & analytical kernels [Jitkrittum et al., 2017]

Given:

- Density/model: p .
- Samples: $\mathbf{X} = \{x_i\}_{i=1}^n \sim q$ (unknown).



Linear-time goodness-of-fit testing: Stein operator & analytical kernels [Jitkrittum et al., 2017]

Given:

- Density/model: p .
- Samples: $X = \{x_i\}_{i=1}^n \sim q$ (unknown).

Problem: using p, X test

$$H_0 : p = q, \text{ vs}$$

$$H_1 : p \neq q.$$

Quick summary:

- Best paper award (NIPS-2017, 3/3240).
- Demo: criminal data analysis.
- Code: <https://github.com/wittawatj/kernel-gof>



Summary

- Dependency measures, distances: KCCA, HSIC, MMD.
- Mean embedding, cross-covariance operator.

- Dependency measures, distances: KCCA, HSIC, MMD.
- Mean embedding, cross-covariance operator.
- Applications:
 - ISA, distribution regression, image registration, feature selection,
 - hypothesis testing.

- Dependency measures, distances: KCCA, HSIC, MMD.
- Mean embedding, cross-covariance operator.
- Applications:
 - ISA, distribution regression, image registration, feature selection,
 - hypothesis testing.
- Hypothesis testing:
 - quadratic methods,
 - scaling: block-variants, Nyström, RFF,
 - linear-time adaptive nonparametric tests.

Thank you for the attention!



-  Altun, Y. and Smola, A. (2006).
Unifying divergence minimization and statistical inference via convex duality.
In *Conference on Learning Theory (COLT)*, pages 139–153.
-  Bach, F. R. and Jordan, M. I. (2002).
Kernel independent component analysis.
Journal of Machine Learning Research, 3:1–48.
-  Baker, C. R. (1973).
Joint measures and cross-covariance operators.
Transactions of the American Mathematical Society, 186:273–289.
-  Baringhaus, L. and Franz, C. (2004).
On a new multivariate two-sample test.
Journal of Multivariate Analysis, 88:190–206.
-  Berg, C., Christensen, J. P. R., and Ressel, P. (1984).
Harmonic Analysis on Semigroups.
Springer-Verlag.

-  Berlinet, A. and Thomas-Agnan, C. (2004).
Reproducing Kernel Hilbert Spaces in Probability and Statistics.
Kluwer.
-  Bertin-Mahieux, T., Ellis, D. P., Whitman, B., and Lamere, P. (2011).
The million song dataset.
In *International Conference on Music Information Retrieval (ISMIR)*.
-  Blanchard, G., Lee, G., and Scott, C. (2011).
Generalizing from several related classification tasks to a new unlabeled sample.
In *Advances in Neural Information Processing Systems (NIPS)*, pages 2178–2186.
-  Brault, R., Heinonen, M., and d'Alché-Buc, F. (2016).
Random Fourier features for operator-valued kernels.

In *Asian Conference in Machine Learning (ACML; JMLR W&CP)*, volume 63, pages 110–125.

-  Caponnetto, A. and De Vito, E. (2007).
Optimal rates for regularized least-squares algorithm.
Foundations of Computational Mathematics, 7:331–368.
-  Cardoso, J.-F. (1998).
Multidimensional independent component analysis.
In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1941–1944.
-  Carmeli, C., Vito, E. D., Toigo, A., and Umanitá, V. (2010).
Vector valued reproducing kernel Hilbert spaces and universality.
Analysis and Applications, 8:19–61.
-  Christmann, A. and Steinwart, I. (2010).
Universal kernels on non-standard input spaces.
In *Advances in Neural Information Processing Systems (NIPS)*, pages 406–414.

-  Chwialkowski, K. and Gretton, A. (2014).
A kernel independence test for random processes.
In *International Conference on Machine Learning (ICML; JMLR W&CP)*, volume 32, page 14221430.
-  Chwialkowski, K., Ramdas, A., Sejdinovic, D., and Gretton, A. (2015).
Fast two-sample testing with analytic representations of probability measures.
In *Advances in Neural Information Processing Systems (NIPS)*, pages 1972–1980.
-  Chwialkowski, K., Sejdinovic, D., and Gretton, A. (2014).
A wild bootstrap for degenerate kernel tests.
In *Advances in Neural Information Processing Systems (NIPS)*, pages 3608–3616.
-  Chwialkowski, K., Strathmann, H., and Gretton, A. (2016).
A kernel test of goodness of fit.

In *International Conference on Machine Learning (ICML)*,
pages 2606–2615.

-  Collins, M. and Duffy, N. (2001).
Convolution kernels for natural language.
In *Advances in Neural Information Processing Systems (NIPS)*,
pages 625–632.
-  Csörgo, S. and Totik, V. (1983).
On how long interval is the empirical characteristic function
uniformly consistent?
Acta Scientiarum Mathematicarum, 45:141–149.
-  Cuturi, M. (2011).
Fast global alignment kernels.
In *International Conference on Machine Learning (ICML)*,
pages 929–936.
-  Cuturi, M., Fukumizu, K., and Vert, J.-P. (2005).
Semigroup kernels on measures.
Journal of Machine Learning Research, 6:1169–1198.

-  Diestel, J. and Uhl, J. J. (1977).
Vector Measures.
American Mathematical Society. Providence.
-  Dinculeanu, N. (2000).
Vector Integration and Stochastic Integration in Banach Spaces.
Wiley.
-  Drineas, P. and Mahoney, M. W. (2005).
On the Nyström method for approximating a Gram matrix for improved kernel-based learning.
Journal of Machine Learning Research, 6:2153–2175.
-  Dudley, R. M. (2004).
Real Analysis and Probability.
Cambridge University Press.
-  Fang, K.-T., Kotz, S., and Ng, K. W. (1990).
Symmetric multivariate and related distributions.

-  Fukumizu, K., Bach, F., and Jordan, M. (2004). Dimensionality reduction for supervised learning with reproducing kernel Hilbert spaces. *Journal of Machine Learning Research*, 5:73–99.
-  Fukumizu, K., Bach, F., and Jordan, M. (2009a). Kernel dimension reduction in regression. *The Annals of Statistics*, 37(4):1871–1905.
-  Fukumizu, K., Bach, F. R., and Gretton, A. (2007). Statistical consistency of kernel canonical correlation analysis. *Journal of Machine Learning Research*, 8:361–383.
-  Fukumizu, K., Gretton, A., Schölkopf, B., and Sriperumbudur, B. K. (2009b). Characteristic kernels on groups and semigroups. In *Advances in Neural Information Processing Systems (NIPS)*, pages 473–480.

-  Fukumizu, K., Gretton, A., Sun, X., and Schölkopf, B. (2008). Kernel measures of conditional dependence. In *Advances in Neural Information Processing Systems (NIPS)*, pages 498–496.
-  Gärtner, T., Flach, P. A., Kowalczyk, A., and Smola, A. (2002). Multi-instance kernels. In *International Conference on Machine Learning (ICML)*, pages 179–186.
-  Gretton, A. (2015). A simpler condition for consistency of a kernel independence test. Technical report, University College London.
[\(https://arxiv.org/abs/1501.06103\).](https://arxiv.org/abs/1501.06103)
-  Gretton, A., Borgwardt, K. M., Rasch, M. J., Schölkopf, B., and Smola, A. (2012). A kernel two-sample test.

-  Gretton, A., Bousquet, O., Smola, A., and Schölkopf, B. (2005a). Measuring statistical dependence with Hilbert-Schmidt norms. In *Algorithmic Learning Theory (ALT)*, pages 63–78.
-  Gretton, A., Fukumizu, K., Harchaoui, Z., and Sriperumbudur, B. K. (2009). A fast, consistent kernel two-sample test. In *Advances in Neural Information Processing Systems (NIPS)*, pages 673–681.
-  Gretton, A., Fukumizu, K., Teo, C. H., Song, L., Schölkopf, B., and Smola, A. J. (2007). A kernel statistical test of independence. In *Advances in Neural Information Processing Systems (NIPS)*, pages 585–592.
-  Gretton, A., Fukumizu, K., Teo, C. H., Song, L., Schölkopf, B., and Smola, A. J. (2008).

A kernel statistical test of independence.

In *Advances in Neural Information Processing Systems (NIPS)*, pages 585–592.

 Gretton, A., Herbrich, R., Smola, A., Bousquet, O., and Schölkopf, B. (2005b).

Kernel methods for measuring independence.

Journal of Machine Learning Research, 6:2075–2129.

 Guevara, J., Hirata, R., and Canu, S. (2017).

Cross product kernels for fuzzy set similarity.

In *IEEE International Conference on Fuzzy Systems (FUZZ-IEEE)*, pages 1–6.

 Habibian, A., Mensink, T., and Snoek, C. G. (2014).

Videostory: A new multimedia embedding for few-example recognition and translation of events.

In *ACM International Conference on Multimedia*, pages 17–26.

 Haussler, D. (1999).

Convolution kernels on discrete structures.

Technical report, Department of Computer Science, University of California at Santa Cruz.
[\(http://cbse.soe.ucsc.edu/sites/default/files/convolutions.pdf\).](http://cbse.soe.ucsc.edu/sites/default/files/convolutions.pdf)

-  Hein, M. and Bousquet, O. (2005).
Hilbertian metrics and positive definite kernels on probability measures.
In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 136–143.
-  Jebara, T., Kondor, R., and Howard, A. (2004).
Probability product kernels.
Journal of Machine Learning Research, 5:819–844.
-  Jiao, Y. and Vert, J.-P. (2016).
The Kendall and Mallows kernels for permutations.
In *International Conference on Machine Learning (ICML; PMLR)*, volume 37, pages 2982–2990.

-  Jitkrittum, W., Szabó, Z., Chwialkowski, K., and Gretton, A. (2016a).
Interpretable distribution features with maximum testing power.
In *Advances in Neural Information Processing Systems (NIPS)*, pages 181–189.
-  Jitkrittum, W., Szabó, Z., and Gretton, A. (2016b).
An adaptive test of independence with analytic kernel embeddings.
Technical report.
(<https://arxiv.org/abs/1610.04782>).
-  Jitkrittum, W., Xu, W., Szabó, Z., Fukumizu, K., and Gretton, A. (2017).
A linear-time kernel goodness-of-fit test.
In *Advances in Neural Information Processing Systems (NIPS)*, pages 261–270.
-  Johnson, N. L., Kotz, S., and Balakrishnan, N. (1994).

Continuous univariate distributions (volume 1).

John Wiley & Sons.

-  Kashima, H. and Koyanagi, T. (2002).
Kernels for semi-structured data.
In *International Conference on Machine Learning (ICML)*,
pages 291–298.
-  Klebanov, L. (2005).
N-Distances and Their Applications.
Charles University, Prague.
-  Kondor, R. and Pan, H. (2016).
The multiscale Laplacian graph kernel.
In *Advances in Neural Information Processing Systems (NIPS)*,
pages 2982–2990.
-  Kybic, J. (2004).
High-dimensional mutual information estimation for image
registration.

In *IEEE International Conference on Image Processing (ICIP)*,
pages 1779–1782.

-  Lancaster, H. O. (1969).
The Chi-squared Distribution.
John Wiley and Sons Inc.
-  Leucht, A. and H. Neumann, M. (2013).
Dependent wild bootstrap for degenerate U- and V-statistics.
Journal of Multivariate Analysis, 117:257–280.
-  Leurgans, S. E., Moyeed, R. A., and Silverman, B. W. (1993).
Canonical correlation analysis when the data are curves.
Journal of the Royal Statistical Society, Series B (Methodological), 55(3):725–740.
-  Li, Z., Ton, J.-F., Ogle, D., and Sejdinovic, D. (2018).
A unified analysis of random Fourier features.
Technical report.
(<https://arxiv.org/abs/1806.09178>).

-  Liu, Q., Lee, J., and Jordan, M. (2016).
A Kernelized Stein Discrepancy for Goodness-of-fit Tests.
In *International Conference on Machine Learning (ICML)*,
pages 276–284.
-  Lodhi, H., Saunders, C., Shawe-Taylor, J., Cristianini, N., and Watkins, C. (2002).
Text classification using string kernels.
Journal of Machine Learning Research, 2:419–444.
-  Lopez-Paz, D. (2016).
From Dependence to Causation.
PhD thesis, University of Cambridge.
-  Lopez-Paz, D., Sra, S., Smola, A., Ghahramani, Z., and Schölkopf, B. (2014).
Randomized nonlinear component analysis.
In *International Conference on Machine Learning (ICML)*,
pages 1359–1367.

-  Lundqvist, D., Flykt, A., and Öhman, A. (1998).
The Karolinska directed emotional faces-KDEF.
Technical report, ISBN 91-630-7164-9.
-  Lyons, R. (2013).
Distance covariance in metric spaces.
Annals of Probability, 41:3284–3305.
-  Martins, A. F. T., Smith, N. A., Xing, E. P., Aguiar, P. M. Q., and Figueiredo, M. A. T. (2009).
Nonextensive information theoretic kernels on measures.
The Journal of Machine Learning Research, 10:935–975.
-  Micchelli, C. A., Xu, Y., and Zhang, H. (2006).
Universal kernels.
Journal of Machine Learning Research, 7:2651–2667.
-  Moulines, É., Bach, F. R., and Harchaoui, Z. (2007).
Testing for homogeneity with kernel Fisher discriminant analysis.

In *Advances in Neural Information Processing Systems (NIPS)*, pages 609–616.

-  Muandet, K., Fukumizu, K., Dinuzzo, F., and Schölkopf, B. (2011).

Learning from distributions via support measure machines.

In *Advances in Neural Information Processing Systems (NIPS)*, pages 10–18.

-  Muandet, K., Fukumizu, K., Sriperumbudur, B., and Schölkopf, B. (2017).

Kernel mean embedding of distributions: A review and beyond.

Foundations and Trends in Machine Learning, 10(1-2):1–141.

-  Müller, A. (1997).

Integral probability metrics and their generating classes of functions.

Advances in Applied Probability, 29:429–443.

 Neemuchwala, H., Hero, A., Zabuawala, S., and Carson, P. (2007).

Image registration methods in high dimensional space.

International Journal of Imaging Systems and Technology, 16:130–145.

 Nishiyama, Y. and Fukumizu, K. (2016).

Characteristic kernels and infinity divisibility.

Journal of Machine Learning Research, 17:1–28.

 Peng, H., Long, F., and Ding, C. (2005).

Feature selection based on mutual information: criteria of max-dependency, max-relevance, and min-redundancy.

IEEE Transactions on Pattern Analysis and Machine Intelligence, 27(8):1226–1238.

 Pfister, N., Bühlmann, P., Schölkopf, B., and Peters, J. (2017).

Kernel-based tests for joint independence.

Journal of the Royal Statistical Society: Series B (Statistical Methodology).

-  Póczos, B., Singh, A., Rinaldo, A., and Wasserman, L. (2013). Distribution-free distribution regression.
In *International Conference on AI and Statistics (AISTATS; JMLR W&CP)*, volume 31, pages 507–515.
-  Póczos, B., Xiong, L., Sutherland, D., and Schneider, J. (2012). Support distribution machines.
Technical report, Carnegie Mellon University.
(<http://arxiv.org/abs/1202.0302>).
-  Rahimi, A. and Recht, B. (2007). Random features for large-scale kernel machines.
In *Advances in Neural Information Processing Systems (NIPS)*, pages 1177–1184.
-  Reed, M. and Simon, B. (1980).

Methods of Modern Mathematical Physics, I: Functional Analysis.

Academic Press.

 Rényi, A. (1959).

On measures of dependence.

Acta Mathematica Academiae Scientiarum Hungaricae,
10:441–451.

 Rosasco, L., Santoro, M., Mosci, S., Verri, A., and Villa, S. (2010).

A regularization approach to nonlinear variable selection.

JMLR W&CP – International Conference on Artificial Intelligence and Statistics (AISTATS), 9:653–660.

 Rosasco, L., Villa, S., Mosci, S., Santoro, M., and Verri, A. (2013).

Nonparametric sparsity and regularization.

Journal of Machine Learning Research, 14:1665–1714.

-  Rubenstein, P. K., Chwialkowski, K. P., and Gretton, A. (2016).
A kernel test for three-variable interactions with random processes.
In *Conference on Uncertainty in Artificial Intelligence (UAI)*, pages 637–646.
-  Rudi, A. and Rosasco, L. (2017).
Generalization properties of learning with random features.
In *Advances in Neural Information Processing Systems (NIPS)*, pages 3215–3225.
-  Rudin, W. (1991).
Functional Analysis.
McGraw-Hill, USA.
-  Sasvári, Z. (2013).
Multivariate Characteristic and Correlation Functions.
Walter de Gruyter.
-  Sejdinovic, D., Gretton, A., and Bergsma, W. (2013a).

A kernel test for three-variable interactions.

In *Advances in Neural Information Processing Systems (NIPS)*,
pages 1124–1132.

 Sejdinovic, D., Sriperumbudur, B., Gretton, A., and Fukumizu, K. (2013b).

Equivalence of distance-based and RKHS-based statistics in hypothesis testing.

Annals of Statistics, 41:2263–2291.

 Serfling, R. J. (1980).

Approximation Theorems of Mathematical Statistics.

John Wiley & Sons.

 Simon-Gabriel, C.-J. and Schölkopf, B. (2018).

Kernel distribution embeddings: Universal kernels, characteristic kernels and kernel metrics on distributions.

Journal of Machine Learning Research, 44:1–29.

 Snelson, E. and Ghahramani, Z. (2006).

Sparse Gaussian processes using pseudo-inputs.

In *Advances in Neural Information Processing Systems (NIPS)*, pages 1257–1264.

-  Sriperumbudur, B. and Sterge, N. (2018). Approximate kernel PCA using random features: Computational vs. statistical trade-off. Technical report, Pennsylvania State University. (<https://arxiv.org/abs/1706.06296>).
-  Sriperumbudur, B. K., Fukumizu, K., Gretton, A., Hyvärinen, A., and Kumar, R. (2017). Density estimation in infinite dimensional exponential families. *Journal of Machine Learning Research*, 18(57):1–59.
-  Sriperumbudur, B. K., Fukumizu, K., Gretton, A., Schölkopf, B., and Lanckriet, G. R. G. (2012). On the empirical estimation of integral probability metrics. *Electronic Journal of Statistics*, 6:1550–1599.
-  Sriperumbudur, B. K., Fukumizu, K., and Lanckriet, G. R. G. (2010a).

On the relation between universality, characteristic kernels and rkhs embedding of measures.

In *International Conference on AI and Statistics (AISTATS; JMLR W&CP)*, volume 9, pages 781–788.

 Sriperumbudur, B. K., Fukumizu, K., and Lanckriet, G. R. G. (2011).

Learning in Hilbert vs. Banach spaces: A measure embedding viewpoint.

In *Advances in Neural Information Processing Systems (NIPS)*, pages 1773–1781.

 Sriperumbudur, B. K., Gretton, A., Fukumizu, K., and Lanckriet, G. R. G. (2010b).

Hilbert space embeddings and metrics on probability measures.

Journal of Machine Learning Research, 11:1517–1561.

 Sriperumbudur, B. K. and Szabó, Z. (2015).
Optimal rates for random Fourier features.

In *Advances in Neural Information Processing Systems (NIPS)*, pages 1144–1152.

 Steinwart, I. (2001).

On the influence of the kernel on the consistency of support vector machines.

Journal of Machine Learning Research, 2:67–93.

 Steinwart, I. and Christmann, A. (2008).

Support Vector Machines.

Springer.

 Sun, Y., Gilbert, A., and Tewari, A. (2018).

But how does it work in theory? Linear SVM with random features.

In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 3383–3392.

 Sutherland, D. J. and Schneider, J. (2015).

On the error of random Fourier features.

In *Conference on Uncertainty in Artificial Intelligence (UAI)*, pages 862–871.

-  Szabó, Z., Gretton, A., Póczos, B., and Sriperumbudur, B. (2015).

Two-stage sampled learning theory on distributions.

In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 948–957.

-  Szabó, Z., Póczos, B., and Lörincz, A. (2012).
Separation theorem for independent subspace analysis and its consequences.

Pattern Recognition, 45(4):1782–1791.

-  Szabó, Z. and Sriperumbudur, B. (2017).
Characteristic and universal tensor product kernels.
Technical report.
(<http://arxiv.org/abs/1708.08157>).

-  Szabó, Z. and Sriperumbudur, B. K. (2019).

On kernel derivative approximation with random Fourier features.

In *International Conference on Artificial Intelligence and Statistics (AISTATS)*.

 Szabó, Z., Sriperumbudur, B. K., Póczos, B., and Gretton, A. (2016).

Learning theory for distribution regression.

Journal of Machine Learning Research, 17(152):1–40.

 Székely, G. J. and Rizzo, M. L. (2004).

Testing for equal distributions in high dimension.

InterStat, 5.

 Székely, G. J. and Rizzo, M. L. (2005).

A new test for multivariate normality.

Journal of Multivariate Analysis, 93:58–80.

 Székely, G. J. and Rizzo, M. L. (2009).

Brownian distance covariance.

The Annals of Applied Statistics, 3:1236–1265.

-  Székely, G. J., Rizzo, M. L., and Bakirov, N. K. (2007). Measuring and testing dependence by correlation of distances. *The Annals of Statistics*, 35:2769–2794.
-  Titsias, M. K. (2009). Variational learning of inducing variables in sparse Gaussian processes. *Journal of Machine Learning Research*, 5:567–574.
-  Ullah, E., Mianjy, P., Marinov, T. V., and Arora, R. (2018). Streaming kernel PCA with $\tilde{O}(\sqrt{n})$ random features. Technical report. (<https://arxiv.org/abs/1808.00934>).
-  van de Geer, S. A. (2009). *Empirical Processes in M-Estimation*. Cambridge University Press.
-  van der Vaart, A. W. (1998). *Asymptotic Statistics*.

-  van der Vaart, A. W. and Wellner, J. A. (1996).
Weak Convergence and Empirical Processes.
Springer-Verlag.
-  Vishwanathan, S. N., Schraudolph, N. N., Kondor, R., and Borgwardt, K. M. (2010).
Graph kernels.
Journal of Machine Learning Research, 11:1201–1242.
-  Waegeman, W., Pahikkala, T., Airola, A., Salakoski, T., Stock, M., and Baets, B. D. (2012).
A kernel-based framework for learning graded relations from data.
IEEE Transactions on Fuzzy Systems, 20:1090–1101.
-  Wang, H. and Schmid, C. (2013).
Action recognition with improved trajectories.
In *IEEE International Conference on Computer Vision (ICCV)*, pages 3551–3558.

-  Wendland, H. (2005).
Scattered Data Approximation.
Cambridge University Press.
-  Williams, C. K. I. and Seeger, M. (2001).
Using the Nyström method to speed up kernel machines.
In *Advances in Neural Information Processing Systems (NIPS)*,
pages 682–688.
-  Zaremba, W., Gretton, A., and Blaschko, M. (2013).
B-tests: Low variance kernel two-sample tests.
In *Advances in Neural Information Processing Systems (NIPS)*,
pages 755–763.
-  Zhang, H., Xu, Y., and Zhang, J. (2009).
Reproducing kernel Banach spaces for machine learning.
Journal of Machine Learning Research, 10:2741–2775.
-  Zhang, Q., Filippi, S., Gretton, A., and Sejdinovic, D. (2017).
Large-scale kernel methods for independence testing.

-  Zhao, J. and Meng, D. (2015).
FastMMD: Ensemble of circular discrepancy for efficient two-sample test.
Neural Computation, 27:1345–1372.
-  Zinger, A. A., Kakosyan, A. V., and Klebanov, L. B. (1992).
A characterization of distributions by mean values of statistics and certain probabilistic metrics.
Journal of Soviet Mathematics.
-  Zolotarev, V. M. (1983).
Probability metrics.
Theory of Probability and its Applications, 28:278–302.