

Characteristic Tensor Kernels

Zoltán Szabó – CMAP, École Polytechnique



Joint work with: Bharath K. Sriperumbudur (PSU)

CREST Statistics Seminar, ENSAE
October 9, 2017

Objects of Interest

Divergence & independence measures on kernel-endowed domains.

- Mean embedding:

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

Objects of Interest

Divergence & independence measures on kernel-endowed domains.

- Mean embedding:

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

- Maximum mean discrepancy (MMD):

$$\text{MMD}_k(\mathbb{P}, \mathbb{Q}) := \|\mu_k(\mathbb{P}) - \mu_k(\mathbb{Q})\|_{\mathcal{H}_k}.$$

Objects of Interest

Divergence & independence measures on kernel-endowed domains.

- Mean embedding:

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

- Maximum mean discrepancy (MMD):

$$\text{MMD}_k(\mathbb{P}, \mathbb{Q}) := \|\mu_k(\mathbb{P}) - \mu_k(\mathbb{Q})\|_{\mathcal{H}_k}.$$

- Hilbert-Schmidt independence criterion (HSIC), $k = \otimes_{m=1}^M k_m$:

$$\text{HSIC}_k(\mathbb{P}) := \left\| \mu_k(\mathbb{P}) - \mu_k\left(\otimes_{m=1}^M \mathbb{P}_m\right) \right\|_{\mathcal{H}_k}.$$

Objects of Interest

- Mean embedding:

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

- Maximum mean discrepancy (MMD):

$$\text{MMD}_k(\mathbb{P}, \mathbb{Q}) := \|\mu_k(\mathbb{P}) - \mu_k(\mathbb{Q})\|_{\mathcal{H}_k}.$$

- Hilbert-Schmidt independence criterion (HSIC), $k = \otimes_{m=1}^M k_m$:

$$\text{HSIC}_k(\mathbb{P}) := \left\| \mu_k(\mathbb{P}) - \mu_k\left(\otimes_{m=1}^M \mathbb{P}_m\right) \right\|_{\mathcal{H}_k}.$$

Question

Conditions on k_m -s so that MMD and HSIC are characteristic?

Intuition: Distribution Representation via Functions

- Cumulative density function:

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

Intuition: Distribution Representation via Functions

- Cumulative density function:

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

Intuition: 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:

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

Intuition: 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:

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

- Moment generating function:

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

Intuition: 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:

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

- Moment generating function:

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

Pattern

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

Ingredients

Ingredients: Domain of the Distributions (\mathcal{X})

- $\mathcal{X} = \times_{m=1}^M \mathcal{X}_m$: product space.
- \mathcal{X}_m : different modalities → images, texts, audio, ...



Ingredients: Domain of the Distributions (\mathcal{X})

- $\mathcal{X} = \times_{m=1}^M \mathcal{X}_m$: product space.
- \mathcal{X}_m : different modalities \rightarrow images, texts, audio, ...



Assumption

\mathcal{X}_m : kernel-endowed domains (\Rightarrow inner product).

Ingredients: Kernel, RKHS ($\mathcal{X} := \mathcal{X}_m$, $k := k_m$)

Given: \mathcal{X} set.

- Kernel:

$$k(a, b) = \langle \varphi(a), \varphi(b) \rangle_{\mathcal{H}}, \quad \mathcal{H} : \text{some Hilbert space.}$$

Ingredients: Kernel, RKHS ($\mathcal{X} := \mathcal{X}_m$, $k := k_m$)

Given: \mathcal{X} set.

- Kernel:

$$k(a, b) = \langle \varphi(a), \varphi(b) \rangle_{\mathcal{H}}, \quad \mathcal{H} : \text{some Hilbert space.}$$

- Reproducing kernel of a Hilbert space $\mathcal{H} \subset \mathbb{R}^{\mathcal{X}}$:

$$k(\cdot, b) \in \mathcal{H}, \quad \langle f, k(\cdot, b) \rangle_{\mathcal{H}} = f(b).$$

Ingredients: Kernel, RKHS ($\mathcal{X} := \mathcal{X}_m$, $k := k_m$)

Given: \mathcal{X} set.

- Kernel:

$$k(a, b) = \langle \varphi(a), \varphi(b) \rangle_{\mathcal{H}}, \quad \mathcal{H} : \text{some Hilbert space.}$$

- Reproducing kernel of a Hilbert space $\mathcal{H} \subset \mathbb{R}^{\mathcal{X}}$:

$$k(\cdot, b) \in \mathcal{H}, \quad \langle f, k(\cdot, b) \rangle_{\mathcal{H}} = f(b).$$

$$\xrightarrow{\text{spec.}} k(a, b) = \langle k(\cdot, a), k(\cdot, b) \rangle_{\mathcal{H}}.$$

Ingredients: Kernel, RKHS – continued

Alternatives:

- $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ symmetric is positive definite if

$$\mathbf{G} = [k(x_i, x_j)]_{i,j=1}^n \geq 0 \quad \forall n \in \mathbb{Z}^+, \forall \{x_i\}_{i=1}^n.$$

Ingredients: Kernel, RKHS – continued

Alternatives:

- $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ symmetric is positive definite if

$$\mathbf{G} = [k(x_i, x_j)]_{i,j=1}^n \geq 0 \quad \forall n \in \mathbb{Z}^+, \forall \{x_i\}_{i=1}^n.$$

- $\mathcal{H} \subset \mathbb{R}^{\mathcal{X}}$ Hilbert space,

$$\delta_x : f \in \mathcal{H} \mapsto f(x) \in \mathbb{R}$$

is continuous for all $x \in \mathcal{X}$.

Ingredients: Kernel, RKHS – continued

Alternatives:

- $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ symmetric is positive definite if

$$\mathbf{G} = [k(x_i, x_j)]_{i,j=1}^n \geq 0 \quad \forall n \in \mathbb{Z}^+, \forall \{x_i\}_{i=1}^n.$$

- $\mathcal{H} \subset \mathbb{R}^{\mathcal{X}}$ Hilbert space,

$$\delta_x : f \in \mathcal{H} \mapsto f(x) \in \mathbb{R}$$

is continuous for all $x \in \mathcal{X}$.

All these definitions are equivalent.

Kernel Examples



- $\mathcal{X} = \mathbb{R}^d, \gamma > 0$:

$$k_p(\mathbf{x}, \mathbf{y}) = (\langle \mathbf{x}, \mathbf{y} \rangle + \gamma)^p, \quad k_G(\mathbf{x}, \mathbf{y}) = e^{-\gamma \|\mathbf{x} - \mathbf{y}\|_2^2},$$

$$k_e(\mathbf{x}, \mathbf{y}) = e^{-\gamma \|\mathbf{x} - \mathbf{y}\|_2}, \quad k_C(\mathbf{x}, \mathbf{y}) = 1 + \frac{1}{\gamma \|\mathbf{x} - \mathbf{y}\|_2^2}.$$

Kernel Examples



- $\mathcal{X} = \mathbb{R}^d, \gamma > 0$:

$$k_p(\mathbf{x}, \mathbf{y}) = (\langle \mathbf{x}, \mathbf{y} \rangle + \gamma)^p, \quad k_G(\mathbf{x}, \mathbf{y}) = e^{-\gamma \|\mathbf{x}-\mathbf{y}\|_2^2},$$

$$k_e(\mathbf{x}, \mathbf{y}) = e^{-\gamma \|\mathbf{x}-\mathbf{y}\|_2}, \quad k_C(\mathbf{x}, \mathbf{y}) = 1 + \frac{1}{\gamma \|\mathbf{x}-\mathbf{y}\|_2^2}.$$

- $\mathcal{X} = \text{texts, strings}$:

- r -spectrum kernel: # of common $\leq r$ -substrings.

Kernel Examples



- $\mathcal{X} = \mathbb{R}^d, \gamma > 0$:

$$k_p(\mathbf{x}, \mathbf{y}) = (\langle \mathbf{x}, \mathbf{y} \rangle + \gamma)^p, \quad k_G(\mathbf{x}, \mathbf{y}) = e^{-\gamma \|\mathbf{x} - \mathbf{y}\|_2^2},$$

$$k_e(\mathbf{x}, \mathbf{y}) = e^{-\gamma \|\mathbf{x} - \mathbf{y}\|_2}, \quad k_C(\mathbf{x}, \mathbf{y}) = 1 + \frac{1}{\gamma \|\mathbf{x} - \mathbf{y}\|_2^2}.$$

- $\mathcal{X} = \text{texts, strings}$:

- r -spectrum kernel: # of common $\leq r$ -substrings.

- $\mathcal{X} = \text{time-series}$: dynamic time-warping.

Kernel Examples



- $\mathcal{X} = \mathbb{R}^d, \gamma > 0$:

$$k_p(\mathbf{x}, \mathbf{y}) = (\langle \mathbf{x}, \mathbf{y} \rangle + \gamma)^p, \quad k_G(\mathbf{x}, \mathbf{y}) = e^{-\gamma \|\mathbf{x}-\mathbf{y}\|_2^2},$$

$$k_e(\mathbf{x}, \mathbf{y}) = e^{-\gamma \|\mathbf{x}-\mathbf{y}\|_2}, \quad k_C(\mathbf{x}, \mathbf{y}) = 1 + \frac{1}{\gamma \|\mathbf{x}-\mathbf{y}\|_2^2}.$$

- $\mathcal{X} = \text{texts, strings:}$
 - r -spectrum kernel: # of common $\leq r$ -substrings.
- $\mathcal{X} = \text{time-series: dynamic time-warping.}$
- $\mathcal{X} = \text{trees, graphs, dynamical systems, sets, permutations, ...}$

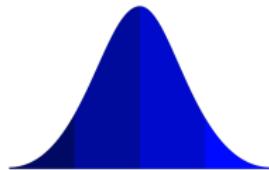
RKHS: Constructively

$k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ kernel, $\xrightarrow{1:1} \mathcal{H}_k$ RKHS.

RKHS: Constructively

$k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ kernel, $\xleftrightarrow{1:1} \mathcal{H}_k$ RKHS.

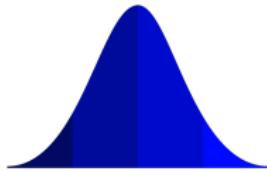
- **Elements** of \mathcal{H}_k : $\overline{\{\sum_{i=1}^n \alpha_i k(\cdot, x_i) : x_i \in \mathcal{X}, \alpha_i \in \mathbb{R}\}}$.



RKHS: Constructively

$k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ kernel, $\xleftrightarrow{1:1} \mathcal{H}_k$ RKHS.

- Elements of \mathcal{H}_k : $\overline{\{\sum_{i=1}^n \alpha_i k(\cdot, x_i) : x_i \in \mathcal{X}, \alpha_i \in \mathbb{R}\}}$.

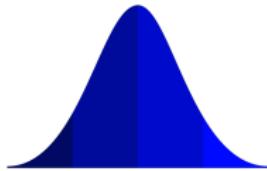


- Inner product:
 - $\langle k(\cdot, x), k(\cdot, y) \rangle_{\mathcal{H}_k} := k(x, y)$.
 - Extension: by linearity & limit.

RKHS: Constructively

$k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ kernel, $\xleftrightarrow{1:1} \mathcal{H}_k$ RKHS.

- Elements of \mathcal{H}_k : $\overline{\{\sum_{i=1}^n \alpha_i k(\cdot, x_i) : x_i \in \mathcal{X}, \alpha_i \in \mathbb{R}\}}$.



- Inner product:
 - $\langle k(\cdot, x), k(\cdot, y) \rangle_{\mathcal{H}_k} := k(x, y)$.
 - Extension: by linearity & limit.

We represent distributions in an RKHS.

Mean embedding: kernel trick \rightarrow mean trick

- Kernel: $k(x, x') = \langle k(\cdot, x), k(\cdot, x') \rangle_{\mathcal{H}_k}$, $\mu_{\mathbb{P}} := k(\cdot, x)$, $\mathbb{P} = \delta_x$.

Mean embedding: kernel trick \rightarrow mean trick

- Kernel: $k(x, x') = \langle k(\cdot, x), k(\cdot, x') \rangle_{\mathcal{H}_k}$, $\mu_{\mathbb{P}} := k(\cdot, x)$, $\mathbb{P} = \delta_x$.
- Mean embedding (feature of \mathbb{P}):

$$\mathbb{P} = \sum_{i=1}^N w_i \delta_{x_i},$$

Mean embedding: kernel trick \rightarrow mean trick

- Kernel: $k(x, x') = \langle k(\cdot, x), k(\cdot, x') \rangle_{\mathcal{H}_k}$, $\mu_{\mathbb{P}} := k(\cdot, x)$, $\mathbb{P} = \delta_x$.
- Mean embedding (feature of \mathbb{P}):

$$\mu_{\mathbb{P}} := \sum_{i=1}^N w_i k(\cdot, x_i) \in \mathcal{H}_k, \quad \mathbb{P} = \sum_{i=1}^N w_i \delta_{x_i},$$

Mean embedding: kernel trick \rightarrow mean trick

- Kernel: $k(x, x') = \langle k(\cdot, x), k(\cdot, x') \rangle_{\mathcal{H}_k}$, $\mu_{\mathbb{P}} := k(\cdot, x)$, $\mathbb{P} = \delta_x$.
- Mean embedding (feature of \mathbb{P}):

$$\mu_{\mathbb{P}} := \sum_{i=1}^N w_i k(\cdot, x_i) \in \mathcal{H}_k, \quad \mathbb{P} = \sum_{i=1}^N w_i \delta_{x_i},$$
$$\mu_{\mathbb{P}} := \underbrace{\int k(\cdot, x) d\mathbb{P}(x)}_{\text{Bochner integral}} \in \mathcal{H}_k.$$

Mean embedding: kernel trick \rightarrow mean trick

- Kernel: $k(x, x') = \langle k(\cdot, x), k(\cdot, x') \rangle_{\mathcal{H}_k}$, $\mu_{\mathbb{P}} := k(\cdot, x)$, $\mathbb{P} = \delta_x$.
- Mean embedding (feature of \mathbb{P}):

$$\mu_{\mathbb{P}} := \sum_{i=1}^N w_i k(\cdot, x_i) \in \mathcal{H}_k, \quad \mathbb{P} = \sum_{i=1}^N w_i \delta_{x_i},$$
$$\mu_{\mathbb{P}} := \underbrace{\int k(\cdot, x) d\mathbb{P}(x)}_{\text{Bochner integral}} \in \mathcal{H}_k.$$

- $\exists \mu_{\mathbb{P}} \Leftrightarrow \int \underbrace{\|k(\cdot, x)\|_{\mathcal{H}_k}}_{\sqrt{k(x,x)}} d\mathbb{P}(x) < \infty$. Assume: bounded k .

- Applications:
 - two-sample testing [Gretton et al., 2012], domain adaptation [Zhang et al., 2013],
 - kernel belief propagation [Song et al., 2011], kernel Bayes' rule [Fukumizu et al., 2013], model criticism [Lloyd et al., 2014],
 - approximate Bayesian computation [Park et al., 2016], probabilistic programming [Schölkopf et al., 2015],
 - distribution classification [Muandet et al., 2011], distribution regression [Szabó et al., 2016], topological data analysis [Kusano et al., 2016].
- Review [Muandet et al., 2017].

Mean Embedding → MMD

- k is called characteristic if

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

is injective [Fukumizu et al., 2008, Sriperumbudur et al., 2010].

Mean Embedding → MMD

- k is called **characteristic** if

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

is **injective** [Fukumizu et al., 2008, Sriperumbudur et al., 2010].

- In this case

$$\text{MMD}_k(\mathbb{P}, \mathbb{Q}) := \|\mu_k(\mathbb{P}) - \mu_k(\mathbb{Q})\|_{\mathcal{H}_k}$$

is a **metric**.

Maximum Mean Discrepancy → HSIC

- $\mathcal{X} := \times_{m=1}^M \mathcal{X}_m$. Tensor product of $(k_m)_{m=1}^M$ kernels:

$$\left(\otimes_{m=1}^M k_m \right) (x, x') = \prod_{m=1}^M k_m (x_m, x'_m), \quad x, x' \in \mathcal{X}$$

is a kernel, $\mathcal{H}_{\otimes_{m=1}^M k_m} = \otimes_{m=1}^M \mathcal{H}_{k_m}$ [Berlinet and Thomas-Agnan, 2004].

Maximum Mean Discrepancy → HSIC

- $\mathcal{X} := \times_{m=1}^M \mathcal{X}_m$. Tensor product of $(k_m)_{m=1}^M$ kernels:

$$\left(\otimes_{m=1}^M k_m \right) (x, x') = \prod_{m=1}^M k_m (x_m, x'_m), \quad x, x' \in \mathcal{X}$$

- is a kernel, $\mathcal{H}_{\otimes_{m=1}^M k_m} = \otimes_{m=1}^M \mathcal{H}_{k_m}$ [Berlinet and Thomas-Agnan, 2004].
- Choosing $k := \otimes_{m=1}^M k_m$, $\mathbb{Q} := \otimes_{m=1}^M \mathbb{P}_m$ in MMD:

$$\text{HSIC}_k (\mathbb{P}) := \left\| \mu_k (\mathbb{P}) - \mu_k \left(\otimes_{m=1}^M \mathbb{P}_m \right) \right\|_{\mathcal{H}_k}.$$

- Blind source separation [Gretton et al., 2005],
- feature selection [Song et al., 2012],
- independence testing [Gretton et al., 2008],
- post selection inference [Yamada et al., 2016],
- causal detection [Mooij et al., 2016, Pfister et al., 2017].

MMD: Easy to Estimate

Using $\{x_i\}_{i=1}^{N_x} \sim \mathbb{P}$, $\{y_j\}_{j=1}^{N_y} \sim \mathbb{Q}$,

$$\widehat{MMD}^2(\mathbb{P}, \mathbb{Q}) = MMD^2(\hat{\mathbb{P}}, \hat{\mathbb{Q}})$$

MMD: Easy to Estimate

Using $\{x_i\}_{i=1}^{N_x} \sim \mathbb{P}$, $\{y_j\}_{j=1}^{N_y} \sim \mathbb{Q}$,

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

MMD: Easy to Estimate

Using $\{x_i\}_{i=1}^{N_x} \sim \mathbb{P}$, $\{y_j\}_{j=1}^{N_y} \sim \mathbb{Q}$,

$$\begin{aligned}\widehat{\text{MMD}}^2(\mathbb{P}, \mathbb{Q}) &= \text{MMD}^2\left(\hat{\mathbb{P}}, \hat{\mathbb{Q}}\right) = \left\| \mu_{\hat{\mathbb{P}}} - \mu_{\hat{\mathbb{Q}}} \right\|_{\mathcal{H}_k}^2 \\ &= \left\| \frac{1}{N_x} \sum_{i=1}^{N_x} k(\cdot, x_i) - \frac{1}{N_y} \sum_{j=1}^{N_y} k(\cdot, y_j) \right\|_{\mathcal{H}_k}^2\end{aligned}$$

MMD: Easy to Estimate

Using $\{x_i\}_{i=1}^{N_x} \sim \mathbb{P}$, $\{y_j\}_{j=1}^{N_y} \sim \mathbb{Q}$,

$$\begin{aligned}\widehat{\text{MMD}}^2(\mathbb{P}, \mathbb{Q}) &= \text{MMD}^2\left(\hat{\mathbb{P}}, \hat{\mathbb{Q}}\right) = \left\| \mu_{\hat{\mathbb{P}}} - \mu_{\hat{\mathbb{Q}}} \right\|_{\mathcal{H}_k}^2 \\ &= \left\| \frac{1}{N_x} \sum_{i=1}^{N_x} k(\cdot, x_i) - \frac{1}{N_y} \sum_{j=1}^{N_y} k(\cdot, y_j) \right\|_{\mathcal{H}_k}^2 \\ &= \underbrace{\frac{1}{N_x^2} \sum_{i,j=1}^{N_x} k(x_i, x_j)}_{\overline{\mathbf{G}_{\mathbb{P}, \mathbb{P}}}} + \underbrace{\frac{1}{N_y^2} \sum_{i,j=1}^{N_y} k(y_i, y_j)}_{\overline{\mathbf{G}_{\mathbb{Q}, \mathbb{Q}}}} - \underbrace{2 \frac{1}{N_x N_y} \sum_{i=1}^{N_x} \sum_{j=1}^{N_y} k(x_i, y_j)}_{\overline{\mathbf{G}_{\mathbb{P}, \mathbb{Q}}}}.\end{aligned}$$

MMD: Easy to Estimate

Using $\{x_i\}_{i=1}^{N_x} \sim \mathbb{P}$, $\{y_j\}_{j=1}^{N_y} \sim \mathbb{Q}$,

$$\begin{aligned}\widehat{\text{MMD}}^2(\mathbb{P}, \mathbb{Q}) &= \text{MMD}^2\left(\hat{\mathbb{P}}, \hat{\mathbb{Q}}\right) = \left\| \mu_{\hat{\mathbb{P}}} - \mu_{\hat{\mathbb{Q}}} \right\|_{\mathcal{H}_k}^2 \\ &= \left\| \frac{1}{N_x} \sum_{i=1}^{N_x} k(\cdot, x_i) - \frac{1}{N_y} \sum_{j=1}^{N_y} k(\cdot, y_j) \right\|_{\mathcal{H}_k}^2 \\ &= \underbrace{\frac{1}{N_x^2} \sum_{i,j=1}^{N_x} k(x_i, x_j)}_{\mathbf{G}_{\mathbb{P}, \mathbb{P}}} + \underbrace{\frac{1}{N_y^2} \sum_{i,j=1}^{N_y} k(y_i, y_j)}_{\mathbf{G}_{\mathbb{Q}, \mathbb{Q}}} - \underbrace{2 \frac{1}{N_x N_y} \sum_{i=1}^{N_x} \sum_{j=1}^{N_y} k(x_i, y_j)}_{\mathbf{G}_{\mathbb{P}, \mathbb{Q}}}.\end{aligned}$$

A bit biased: $\overline{\mathbf{G}_{\mathbb{P}, \mathbb{P}}} \leftarrow \frac{1}{N_x^2} \sum_{i=1}^{N_x} \sum_{j=1, j \neq i}^{N_x} k(x_i, x_j)$, $\overline{\mathbf{G}_{\mathbb{Q}, \mathbb{Q}}} \leftarrow \dots$, if needed.

HSIC: Easy to Estimate

- Given: $\{(x_i, y_i)\}_{i=1}^N$ paired samples.
- Estimate:

$$\widehat{HSIC^2} = \frac{1}{N^2} \left\langle \tilde{\mathbf{G}}_x, \tilde{\mathbf{G}}_y \right\rangle_F$$

HSIC: Easy to Estimate

- Given: $\{(x_i, y_i)\}_{i=1}^N$ paired samples.
- Estimate:

$$\widehat{HSIC^2} = \frac{1}{N^2} \left\langle \tilde{\mathbf{G}}_x, \tilde{\mathbf{G}}_y \right\rangle_F,$$

$$\mathbf{G}_x = [k_1(x_i, x_j)]_{i,j=1}^n$$

HSIC: Easy to Estimate

- Given: $\{(x_i, y_i)\}_{i=1}^N$ paired samples.
- Estimate:

$$\widehat{HSIC^2} = \frac{1}{N^2} \left\langle \tilde{\mathbf{G}}_x, \tilde{\mathbf{G}}_y \right\rangle_F,$$

$$\mathbf{G}_x = [k_1(x_i, x_j)]_{i,j=1}^n, \quad \tilde{\mathbf{G}}_x = \mathbf{H}\mathbf{G}_x\mathbf{H}, \quad \mathbf{H} = \mathbf{I} - \frac{\mathbf{E}}{N}.$$

- Plug-in estimator; similarly easy to debias.

Central in Applications: Characteristic Property

- MMD: ' k is characteristic' means

$$MMD_k(\mathbb{P}, \mathbb{Q}) = 0 \Leftrightarrow \mathbb{P} = \mathbb{Q}.$$

Central in Applications: Characteristic Property

- MMD: ' k is characteristic' means

$$MMD_k(\mathbb{P}, \mathbb{Q}) = 0 \Leftrightarrow \mathbb{P} = \mathbb{Q}.$$

- HSIC: $k = \otimes_{m=1}^M k_M$ will be called **\mathcal{I} -characteristic** if

$$HSIC_k(\mathbb{P}) = 0 \Leftrightarrow \mathbb{P} = \otimes_{m=1}^M \mathbb{P}_m.$$

Central in Applications: Characteristic Property

- MMD: ' k is characteristic' means

$$MMD_k(\mathbb{P}, \mathbb{Q}) = 0 \Leftrightarrow \mathbb{P} = \mathbb{Q}.$$

- HSIC: $k = \otimes_{m=1}^M k_M$ will be called \mathcal{I} -characteristic if

$$HSIC_k(\mathbb{P}) = 0 \Leftrightarrow \mathbb{P} = \otimes_{m=1}^M \mathbb{P}_m.$$

- $\otimes_{m=1}^M k_M$: characteristic $\Rightarrow \mathcal{I}$ -characteristic.

Central in Applications: Characteristic Property

- MMD: ' k is characteristic' means

$$MMD_k(\mathbb{P}, \mathbb{Q}) = 0 \Leftrightarrow \mathbb{P} = \mathbb{Q}.$$

- HSIC: $k = \otimes_{m=1}^M k_M$ will be called \mathcal{I} -characteristic if

$$HSIC_k(\mathbb{P}) = 0 \Leftrightarrow \mathbb{P} = \otimes_{m=1}^M \mathbb{P}_m.$$

- $\otimes_{m=1}^M k_M$: characteristic $\Rightarrow \mathcal{I}$ -characteristic.

Wanted

- Precise relation between characteristic and \mathcal{I} -characteristic?
- Conditions in terms of k_m -s?

Characteristic property

Well-understood for

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

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

Characteristic property

Well-understood for

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

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

- In this case (Bochner's theorem):

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

$$\|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k} = \|c_{\mathbb{P}} - c_{\mathbb{Q}}\|_{L^2(\Lambda)}.$$

Characteristic property

Well-understood for

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

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

- In this case (Bochner's theorem):

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

$$\|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k} = \|c_{\mathbb{P}} - c_{\mathbb{Q}}\|_{L^2(\Lambda)}.$$

Theorem ([Sriperumbudur et al., 2010])

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

Translation-invariant kernels on \mathbb{R}

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	$\cos(\sigma x)$	$\{-\sigma, \sigma\}$

Translation-invariant kernels on \mathbb{R}

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

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 $(C(\mathcal{X}), \|\cdot\|_\infty)$.

Universal \Rightarrow Characteristic

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

Universal \Rightarrow Characteristic

- [Micchelli et al., 2006]: k is c-universal $\Leftrightarrow \mu_k$ is injective on $\mathcal{M}_b(\mathcal{X})$, the set of finite signed Borel measures on \mathcal{X} .
- If $\mathcal{X} = \text{LCP}$: c_0 -universality [Sriperumbudur et al., 2010]
 - \mathcal{H}_k dense in $C_0(\mathcal{X})$, equivalently
 - $\mu_k : \mathcal{M}_b(\mathcal{X}) \mapsto \mathcal{H}_k$ is injective.

Universal \Rightarrow Characteristic

- [Micchelli et al., 2006]: k is c-universal $\Leftrightarrow \mu_k$ is injective on $\mathcal{M}_b(\mathcal{X})$, the set of finite signed Borel measures on \mathcal{X} .
- If $\mathcal{X} = \text{LCP}$: c_0 -universality [Sriperumbudur et al., 2010]
 - \mathcal{H}_k dense in $C_0(\mathcal{X})$, equivalently
 - $\mu_k : \mathcal{M}_b(\mathcal{X}) \mapsto \mathcal{H}_k$ is injective.

LCP examples: \mathbb{R}^d , countable discrete.

- [Micchelli et al., 2006]: k is c -universal $\Leftrightarrow \mu_k$ is injective on $\mathcal{M}_b(\mathcal{X})$, the set of finite signed Borel measures on \mathcal{X} .
- If $\mathcal{X} = \text{LCP}$: c_0 -universality [Sriperumbudur et al., 2010]
 - \mathcal{H}_k dense in $C_0(\mathcal{X})$, equivalently
 - $\mu_k : \mathcal{M}_b(\mathcal{X}) \mapsto \mathcal{H}_k$ is injective.

LCP examples: \mathbb{R}^d , countable discrete.

- c_0 -universality \Rightarrow characteristic.

Local Summary

- Setup: $\mathcal{X} = \times_{m=1}^M \mathcal{X}_m$, (\mathcal{X}_m, k_m) : kernel-enriched domains.
- Mean embedding, MMD, HSIC, $\otimes_{m=1}^M k_m$.
- Characteristic, c_0 -universality.

- Setup: $\mathcal{X} = \times_{m=1}^M \mathcal{X}_m$, (\mathcal{X}_m, k_m) : kernel-enriched domains.
- Mean embedding, MMD, HSIC, $\otimes_{m=1}^M k_m$.
- Characteristic, c_0 -universality.

Little is known about the

characteristic/ \mathcal{I} -characteristic/universality of $\otimes_{m=1}^M k_m$ in terms of k_m -s.

Known ' \mathcal{I} -characteristic' Results, $M = 2$

- [Waegeman et al., 2012, Gretton, 2015]:
 $k_1 \& k_2$: universal $\Rightarrow k_1 \otimes k_2$: universal ($\Rightarrow \mathcal{I}$ -characteristic).

Known ' \mathcal{I} -characteristic' Results, $M = 2$

- [Waegeman et al., 2012, Gretton, 2015]:
 $k_1 \& k_2$: universal $\Rightarrow k_1 \otimes k_2$: universal ($\Rightarrow \mathcal{I}$ -characteristic).
- Stronger: combining [Lyons, 2013] (DCov) and
[Sejdinovic et al., 2013] (DCov \Leftrightarrow HSIC)
 $k_1 \& k_2$: characteristic $\Leftrightarrow k_1 \otimes k_2$: \mathcal{I} -characteristic.

Known ' \mathcal{I} -characteristic' Results, $M = 2$

- [Waegeman et al., 2012, Gretton, 2015]:
 $k_1 \& k_2$: universal $\Rightarrow k_1 \otimes k_2$: universal ($\Rightarrow \mathcal{I}$ -characteristic).
- Stronger: combining [Lyons, 2013] (DCov) and [Sejdinovic et al., 2013] (DCov \Leftrightarrow HSIC)
 $k_1 \& k_2$: characteristic $\Leftrightarrow k_1 \otimes k_2$: \mathcal{I} -characteristic.

Question

Extension to $M \geq 2$?

Known ' \mathcal{I} -characteristic' Results, $M = 2$

- [Waegeman et al., 2012, Gretton, 2015]:
 $k_1 \& k_2$: universal $\Rightarrow k_1 \otimes k_2$: universal ($\Rightarrow \mathcal{I}$ -characteristic).
- Stronger: combining [Lyons, 2013] (DCov) and [Sejdinovic et al., 2013] (DCov \Leftrightarrow HSIC)
 $k_1 \& k_2$: characteristic $\Leftrightarrow k_1 \otimes k_2$: \mathcal{I} -characteristic.

Question

Extension to $M \geq 2$?

Main Challenge

' $\otimes k_m$: \mathcal{I} -characteristic $\Leftrightarrow k_m$: characteristic ($\forall m$)' does NOT hold.

Idea: Characteristic Property as lspd

- Characteristic property:

$$\underbrace{\|\mu_{\mathbb{P}_1 - \mathbb{P}_2}\|_k^2}_{\int_{\mathcal{X}} \int_{\mathcal{X}} k(x, x') d\mathbb{F}(x) d\mathbb{F}(x'), \mathbb{F} := \mathbb{P}_1 - \mathbb{P}_2 \neq 0, \mathbb{F}(\mathcal{X}) = 0} > 0, \quad \forall \mathbb{P}_1, \mathbb{P}_2 \in \mathcal{M}_1^+(\mathcal{X}), \mathbb{P}_1 \neq \mathbb{P}_2.$$

Idea: Characteristic Property as lspd

- Characteristic property:

$$\underbrace{\|\mu_{\mathbb{P}_1 - \mathbb{P}_2}\|_k^2}_{\int_{\mathcal{X}} \int_{\mathcal{X}} k(x, x') d\mathbb{F}(x) d\mathbb{F}(x'), \mathbb{F} := \mathbb{P}_1 - \mathbb{P}_2 \neq 0, \mathbb{F}(\mathcal{X}) = 0} > 0, \quad \forall \mathbb{P}_1, \mathbb{P}_2 \in \mathcal{M}_1^+(\mathcal{X}), \mathbb{P}_1 \neq \mathbb{P}_2.$$

- Observation [Sriperumbudur et al., 2010]: k is characteristic iff.

$$\int_{\mathcal{X}} \int_{\mathcal{X}} k(x, x') d\mathbb{F}(x) d\mathbb{F}(x') > 0, \quad \forall \underbrace{\mathbb{F} \in \mathcal{M}_b(\mathcal{X}) \setminus \{0\} \quad \mathbb{F}(\mathcal{X}) = 0}_{\mathcal{F}_1}.$$

Idea: Characteristic Property as lspd

- Characteristic property:

$$\underbrace{\|\mu_{\mathbb{P}_1 - \mathbb{P}_2}\|_k^2}_{\int_{\mathcal{X}} \int_{\mathcal{X}} k(x, x') d\mathbb{F}(x) d\mathbb{F}(x'), \mathbb{F} := \mathbb{P}_1 - \mathbb{P}_2 \neq 0, \mathbb{F}(\mathcal{X}) = 0} > 0, \quad \forall \mathbb{P}_1, \mathbb{P}_2 \in \mathcal{M}_1^+(\mathcal{X}), \mathbb{P}_1 \neq \mathbb{P}_2.$$

- Observation [Sriperumbudur et al., 2010]: k is characteristic iff.

$$\int_{\mathcal{X}} \int_{\mathcal{X}} k(x, x') d\mathbb{F}(x) d\mathbb{F}(x') > 0, \quad \forall \underbrace{\mathbb{F} \in \mathcal{M}_b(\mathcal{X}) \setminus \{0\} \quad \mathbb{F}(\mathcal{X}) = 0}_{\mathcal{F}_1}.$$

- We have also seen: k is c_0 -universal iff.

$$\int_{\mathcal{X}} \int_{\mathcal{X}} k(x, x') d\mathbb{F}(x) d\mathbb{F}(x') > 0, \quad \forall \underbrace{\mathbb{F} \in \mathcal{M}_b(\mathcal{X}) \setminus \{0\}}_{\mathcal{F}_2}.$$

From now on: $\mathcal{X} = \bigotimes_{m=1}^M \mathcal{X}_m$.

Definition

Let $\mathcal{F} \subseteq \mathcal{M}_b(\mathcal{X})$, $0 \in \mathcal{F}$. $k = \bigotimes_{m=1}^M k_m$ is called **\mathcal{F} -ispd** if

$$\mu_k(\mathbb{F}) = 0 \Rightarrow \mathbb{F} = 0 \quad (\mathbb{F} \in \mathcal{F}), \text{ equivalently}$$

$$\|\mu_k(\mathbb{F})\|_{\mathcal{H}_k}^2 = \int_{\mathcal{X}} \int_{\mathcal{X}} k(x, x') d\mathbb{F}(x) d\mathbb{F}(x') > 0, \quad \forall \mathbb{F} \in \mathcal{F} \setminus \{0\}.$$

Examples

\mathcal{F}	\mathcal{F} -is pd k
$\mathcal{M}_b(\mathcal{X})$	c_0 -universal
$[\mathcal{M}_b(\mathcal{X})]^0$	characteristic

$$\subseteq \quad \subseteq [\mathcal{M}_b(\mathcal{X})]^0 \subseteq \mathcal{M}_b(\mathcal{X}).$$

\cup

$$\Leftarrow \quad \Leftarrow \text{ characteristic} \Leftarrow c_0\text{-universal}.$$

\Downarrow

Examples

\mathcal{F}	\mathcal{F} -is pd k
$\mathcal{M}_b(\mathcal{X})$	c_0 -universal
$[\mathcal{M}_b(\mathcal{X})]^0$	characteristic
$\mathcal{I} := \{\mathbb{P} - \otimes_{m=1}^M \mathbb{P}_m : \mathbb{P} \in \mathcal{M}_1^+(\mathcal{X})\}$	\mathcal{I} -characteristic

$$\subseteq \quad \subseteq [\mathcal{M}_b(\mathcal{X})]^0 \subseteq \mathcal{M}_b(\mathcal{X}).$$

\cup
 \mathcal{I}

$$\Leftarrow \quad \Leftarrow \text{ characteristic} \Leftarrow c_0\text{-universal}.$$

\Downarrow
 $\mathcal{I}\text{-characteristic}$

Examples

\mathcal{F}	\mathcal{F} -is pd k
$\mathcal{M}_b(\mathcal{X})$	c_0 -universal
$[\mathcal{M}_b(\mathcal{X})]^0$	characteristic
$\mathcal{I} := \{\mathbb{P} - \bigotimes_{m=1}^M \mathbb{P}_m : \mathbb{P} \in \mathcal{M}_1^+(\mathcal{X})\}$	\mathcal{I} -characteristic
$[\bigotimes_{m=1}^M \mathcal{M}_b(\mathcal{X}_m)]^0$	\otimes -characteristic

$$\subseteq [\bigotimes_{m=1}^M \mathcal{M}_b(\mathcal{X}_m)]^0 \subseteq [\mathcal{M}_b(\mathcal{X})]^0 \subseteq \mathcal{M}_b(\mathcal{X}).$$

\begin{matrix}
 \cup \\
 \mathcal{I}
 \end{matrix}

$$\begin{array}{ccccccc}
 & \Leftarrow & \otimes\text{-characteristic} & \Leftarrow & \text{characteristic} & \Leftarrow & c_0\text{-universal}. \\
 & & & & & \Downarrow & \\
 & & & & & & \mathcal{I}\text{-characteristic}
 \end{array}$$

Examples

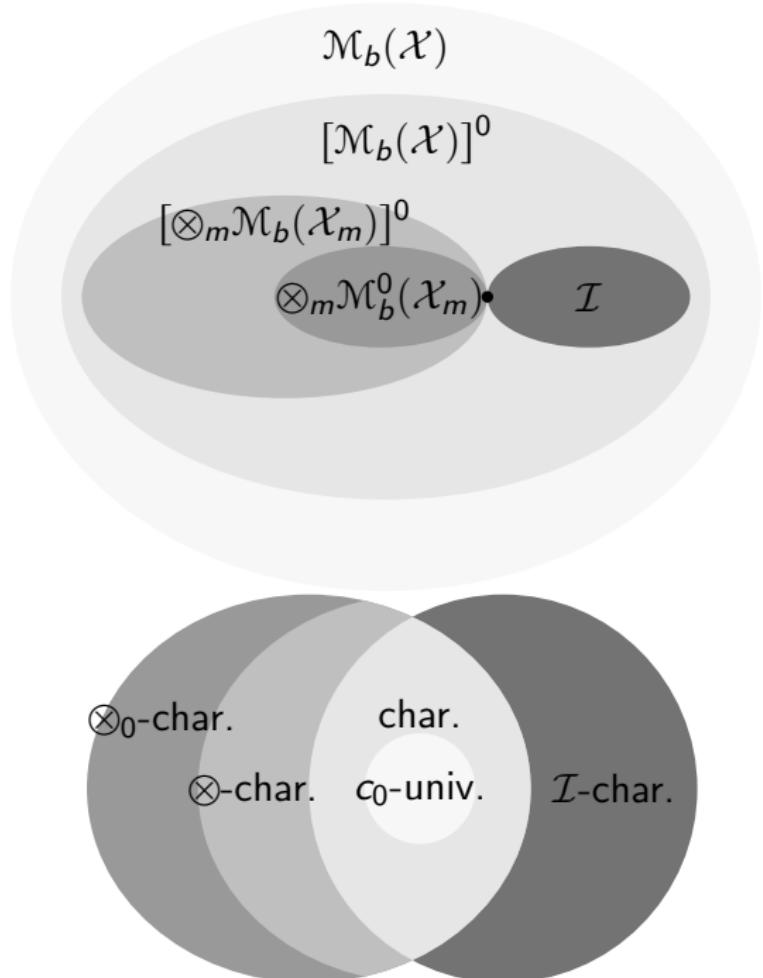
\mathcal{F}	\mathcal{F} -is pd k
$\mathcal{M}_b(\mathcal{X})$	c_0 -universal
$[\mathcal{M}_b(\mathcal{X})]^0$	characteristic
$\mathcal{I} := \{\mathbb{P} - \bigotimes_{m=1}^M \mathbb{P}_m : \mathbb{P} \in \mathcal{M}_1^+(\mathcal{X})\}$	\mathcal{I} -characteristic
$[\bigotimes_{m=1}^M \mathcal{M}_b(\mathcal{X}_m)]^0$	\otimes -characteristic
$\bigotimes_{m=1}^M \mathcal{M}_b^0(\mathcal{X}_m)$	\otimes_0 -characteristic

$$\bigotimes_{m=1}^M \mathcal{M}_b^0(\mathcal{X}_m) \subseteq [\bigotimes_{m=1}^M \mathcal{M}_b(\mathcal{X}_m)]^0 \subseteq [\mathcal{M}_b(\mathcal{X})]^0 \subseteq \mathcal{M}_b(\mathcal{X}).$$

\sqcup
 \mathcal{I}

$$\otimes_0\text{-characteristic} \Leftarrow \otimes\text{-characteristic} \Leftarrow \text{characteristic} \Leftarrow c_0\text{-universal}.$$

\Downarrow
 $\mathcal{I}\text{-characteristic}$



Results

Proposition

- (i) $\otimes_{m=1}^M k_m$: *characteristic* \Rightarrow \otimes -*characteristic*.
- (ii) $\otimes_{m=1}^M k_m$: \otimes -*characteristic* \Rightarrow \otimes_0 -*characteristic*.
- (iii) $\otimes_{m=1}^M k_m$: \otimes_0 -*characteristic* $\Leftrightarrow (k_m)_{m=1}^M$ are *characteristic*.

Proposition

- (i) $\otimes_{m=1}^M k_m$: characteristic \Rightarrow \otimes -characteristic.
- (ii) $\otimes_{m=1}^M k_m$: \otimes -characteristic \Rightarrow \otimes_0 -characteristic.
- (iii) $\otimes_{m=1}^M k_m$: \otimes_0 -characteristic $\Leftrightarrow (k_m)_{m=1}^M$ are characteristic.

(iii) remains. Proof idea: with $k = \otimes_{m=1}^M k_m$, $\mathbb{F} = \otimes_{m=1}^M \mathbb{F}_m$,

$$\|\mu_k(\mathbb{F})\|_{\mathcal{H}_k}^2 = \prod_{m=1}^M \|\mu_{k_m}(\mathbb{F}_m)\|_{\mathcal{H}_{k_m}}^2,$$

$$\mathbb{F} \in \left[\otimes_{m=1}^M \mathcal{M}_b^0(\mathcal{X}_m) \right] \setminus \{0\} \Leftrightarrow \forall m : \mathbb{F}_m \in \mathcal{M}_b^0(\mathcal{X}_m) \setminus \{0\}.$$

Reverse of (ii) does not hold.

Example

- $\mathcal{X}_m = \{1, 2\}$, $\tau_{\mathcal{X}_m} = \mathcal{P}(\{1, 2\})$, $k_m(x, x') = 2\delta_{x,x'} - 1$, $M = 2$.
- $k_1 = k_2$: characteristic, but $k_1 \otimes k_2$ is not \otimes -characteristic.
- $k_1 \otimes k_2$ is \mathcal{I} -characteristic.

Proof idea: $k_1 \otimes k_2$: not \otimes -characteristic

Goal: construct a witness $0 \neq \mathbb{F} = \mathbb{F}_1 \otimes \mathbb{F}_2 \in \otimes_{m=1}^2 \mathcal{M}_b(\mathcal{X}_m)$ s.t.

$$0 = \mathbb{F}(\mathcal{X}_1 \times \mathcal{X}_2) = \mathbb{F}_1(\mathcal{X}_1)\mathbb{F}_2(\mathcal{X}_2),$$

$$0 = \int_{\mathcal{X}_1 \times \mathcal{X}_2} \int_{\mathcal{X}_1 \times \mathcal{X}_2} k_1(x_1, x'_1) k_2(x_2, x'_2) d\mathbb{F}(x_1, x_2) d\mathbb{F}(x'_1, x'_2).$$

Proof idea: $k_1 \otimes k_2$: not \otimes -characteristic

Goal: construct a witness $0 \neq \mathbb{F} = \mathbb{F}_1 \otimes \mathbb{F}_2 \in \otimes_{m=1}^2 \mathcal{M}_b(\mathcal{X}_m)$ s.t.

$$0 = \mathbb{F}(\mathcal{X}_1 \times \mathcal{X}_2) = \mathbb{F}_1(\mathcal{X}_1)\mathbb{F}_2(\mathcal{X}_2),$$

$$0 = \int_{\mathcal{X}_1 \times \mathcal{X}_2} \int_{\mathcal{X}_1 \times \mathcal{X}_2} k_1(x_1, x'_1) k_2(x_2, x'_2) d\mathbb{F}(x_1, x_2) d\mathbb{F}(x'_1, x'_2).$$

Finite signed measures on $\{1, 2\}$:

$$\mathbb{F}_1 = \mathbb{F}_1(\mathbf{a}) = a_1 \delta_1 + a_2 \delta_2, \quad \mathbb{F}_2 = \mathbb{F}_2(\mathbf{b}) = b_1 \delta_1 + b_2 \delta_2.$$

Proof idea: $k_1 \otimes k_2$: not \otimes -characteristic

Goal: construct a witness $0 \neq \mathbb{F} = \mathbb{F}_1 \otimes \mathbb{F}_2 \in \otimes_{m=1}^2 \mathcal{M}_b(\mathcal{X}_m)$ s.t.

$$0 = \mathbb{F}(\mathcal{X}_1 \times \mathcal{X}_2) = \mathbb{F}_1(\mathcal{X}_1)\mathbb{F}_2(\mathcal{X}_2),$$

$$0 = \int_{\mathcal{X}_1 \times \mathcal{X}_2} \int_{\mathcal{X}_1 \times \mathcal{X}_2} k_1(x_1, x'_1) k_2(x_2, x'_2) d\mathbb{F}(x_1, x_2) d\mathbb{F}(x'_1, x'_2).$$

Finite signed measures on $\{1, 2\}$:

$$\mathbb{F}_1 = \mathbb{F}_1(\mathbf{a}) = a_1\delta_1 + a_2\delta_2, \quad \mathbb{F}_2 = \mathbb{F}_2(\mathbf{b}) = b_1\delta_1 + b_2\delta_2.$$

This gives

$$0 = (a_1 + a_2)(b_1 + b_2), \quad 0 = (a_1 - a_2)^2(b_1 - b_2)^2.$$

Proof idea: $k_1 \otimes k_2$: not \otimes -characteristic

Goal: construct a witness $0 \neq \mathbb{F} = \mathbb{F}_1 \otimes \mathbb{F}_2 \in \otimes_{m=1}^2 \mathcal{M}_b(\mathcal{X}_m)$ s.t.

$$0 = \mathbb{F}(\mathcal{X}_1 \times \mathcal{X}_2) = \mathbb{F}_1(\mathcal{X}_1)\mathbb{F}_2(\mathcal{X}_2),$$

$$0 = \int_{\mathcal{X}_1 \times \mathcal{X}_2} \int_{\mathcal{X}_1 \times \mathcal{X}_2} k_1(x_1, x'_1) k_2(x_2, x'_2) d\mathbb{F}(x_1, x_2) d\mathbb{F}(x'_1, x'_2).$$

Finite signed measures on $\{1, 2\}$:

$$\mathbb{F}_1 = \mathbb{F}_1(\mathbf{a}) = a_1\delta_1 + a_2\delta_2, \quad \mathbb{F}_2 = \mathbb{F}_2(\mathbf{b}) = b_1\delta_1 + b_2\delta_2.$$

This gives

$$0 = (a_1 + a_2)(b_1 + b_2), \quad 0 = (a_1 - a_2)^2(b_1 - b_2)^2.$$

\Rightarrow Two symmetric solutions ($\mathbf{a} \neq \mathbf{0}$, $\mathbf{b} \neq \mathbf{0}$):

$$a_1 + a_2 = 0,$$

$$b_1 = b_2$$

$$a_1 = a_2,$$

$$b_1 + b_2 = 0.$$

Towards \mathcal{I} -characteristicity

In the previous example:

$$k_1, k_2: \text{characteristic} \Rightarrow k_1 \otimes k_2: \mathcal{I}\text{-characteristic}.$$

In fact:

- this holds for any bounded kernel,
- +converse for any $M \geqslant 2!$ Formally, ...

\mathcal{I} -characteristic Property

Proposition

- (i) k_1, k_2 : characteristic $\Rightarrow k_1 \otimes k_2$: \mathcal{I} -characteristic.
- (ii) Suppose \mathcal{X}_m is Hausdorff ($\forall m$). Then
 $\otimes_{m=1}^M k_m$: \mathcal{I} -characteristic $\Rightarrow (k_m)_{m=1}^M$ are characteristic.

Proposition

(i) k_1, k_2 : characteristic $\Rightarrow k_1 \otimes k_2$: \mathcal{I} -characteristic.

(ii) Suppose \mathcal{X}_m is Hausdorff ($\forall m$). Then

$\otimes_{m=1}^M k_m$: \mathcal{I} -characteristic $\Rightarrow (k_m)_{m=1}^M$ are characteristic.

Proof idea:

(i) Induction: see later c_0 -universality ($\mathbb{F} = \mathbb{P} - \mathbb{P}_1 \otimes \mathbb{P}_2$).

\mathcal{I} -characteristic Property

Proposition

- (i) k_1, k_2 : characteristic $\Rightarrow k_1 \otimes k_2$: \mathcal{I} -characteristic.
- (ii) Suppose \mathcal{X}_m is Hausdorff ($\forall m$). Then
 $\otimes_{m=1}^M k_m$: \mathcal{I} -characteristic $\Rightarrow (k_m)_{m=1}^M$ are characteristic.

Proof idea:

- (i) Induction: see later c_0 -universality ($\mathbb{F} = \mathbb{P} - \mathbb{P}_1 \otimes \mathbb{P}_2$).
- (ii) If a k_m is not characteristic, then we construct a witness for the non- \mathcal{I} -characteristic property.

k_1, k_2, k_3 : characteristic $\Rightarrow \otimes_{m=1}^3 k_m$: \mathcal{I} -characteristic

Example

- $\mathcal{X}_m = \{1, 2\}$, $\tau_{\mathcal{X}_m} = \mathcal{P}(\{1, 2\})$, $k_m(x, x') = 2\delta_{x,x'} - 1$, $M = 3$.
- Then
 - $(k_m)_{m=1}^3$: characteristic.
 - $\otimes_{m=1}^3 k_m$: is **not** \mathcal{I} -characteristic. Witness:

$$p_{1,1,1} = \frac{1}{5}, \quad p_{1,1,2} = \frac{1}{10}, \quad p_{1,2,1} = \frac{1}{10}, \quad p_{1,2,2} = \frac{1}{10},$$
$$p_{2,1,1} = \frac{1}{5}, \quad p_{2,1,2} = \frac{1}{10}, \quad p_{2,2,1} = \frac{1}{10}, \quad p_{2,2,2} = \frac{1}{10},$$

Non- \mathcal{I} -characteristicity: Analytical Solution

Parameter: $\mathbf{z} = (z_0, z_1, \dots, z_5) \in [0, 1]^6$.

Non- \mathcal{I} -characteristicity: Analytical Solution

Parameter: $\mathbf{z} = (z_0, z_1, \dots, z_5) \in [0, 1]^6$. Example: $p_{1,1,1} =$

$$\begin{aligned} & z_2 + z_1 + z_4 + z_5 - 3z_2z_1 - 4z_2z_4 - 4z_1z_4 - z_2z_3 - 2z_2z_0 - 2z_1z_3 - 3z_2z_5 \\ & - 2z_4z_3 - z_1z_0 - 3z_1z_5 - 2z_4z_0 - 4z_4z_5 - z_3z_0 - z_3z_5 - z_0z_5 + 2z_2z_1^2 + 2z_2^2z_1 \\ & + 4z_2z_4^2 + 2z_2^2z_4 + 4z_1z_4^2 + 2z_1^2z_4 + 2z_2^2z_0 + 2z_1^2z_3 + 2z_2z_5^2 + 2z_2^2z_5 + 2z_4^2z_3 \\ & + 2z_1z_5^2 + 2z_1^2z_5 + 2z_4^2z_0 + 2z_4z_5^2 + 4z_4^2z_5 - z_2^2 - z_1^2 - 3z_4^2 + 2z_4^3 - z_5^2 \\ & + 6z_2z_1z_4 + 2z_2z_1z_3 + 2z_2z_4z_3 + 2z_2z_1z_0 + 4z_2z_1z_5 + 4z_2z_4z_0 + 4z_1z_4z_3 \\ & + 6z_2z_4z_5 + 2z_1z_4z_0 + 6z_1z_4z_5 + 2z_2z_3z_0 + 2z_2z_3z_5 + 2z_1z_3z_0 + 2z_2z_0z_5 \\ & + 2z_1z_3z_5 + 2z_4z_3z_0 + 2z_4z_3z_5 + 2z_1z_0z_5 + 2z_4z_0z_5 \end{aligned} - \frac{2z_2z_1 - z_1 - 2z_4 - z_3 - z_0 - 2z_5 - z_2 + 2z_2z_4 + 2z_1z_4 + 2z_2z_0 + 2z_1z_3 + 2z_2z_5}{2z_4z_3 + 2z_1z_5 + 2z_4z_0 + 4z_4z_5 + 2z_3z_0 + 2z_3z_5 + 2z_0z_5 + 2z_4^2 + 2z_5^2}.$$

Non- \mathcal{I} -characteristicity: Analytical Solution

Parameter: $\mathbf{z} = (z_0, z_1, \dots, z_5) \in [0, 1]^6$. Example: $p_{1,1,1} =$

$$\begin{aligned} & z_2 + z_1 + z_4 + z_5 - 3z_2z_1 - 4z_2z_4 - 4z_1z_4 - z_2z_3 - 2z_2z_0 - 2z_1z_3 - 3z_2z_5 \\ & - 2z_4z_3 - z_1z_0 - 3z_1z_5 - 2z_4z_0 - 4z_4z_5 - z_3z_0 - z_3z_5 - z_0z_5 + 2z_2z_1^2 + 2z_2^2z_1 \\ & + 4z_2z_4^2 + 2z_2^2z_4 + 4z_1z_4^2 + 2z_1^2z_4 + 2z_2^2z_0 + 2z_1^2z_3 + 2z_2z_5^2 + 2z_2^2z_5 + 2z_4^2z_3 \\ & + 2z_1z_5^2 + 2z_1^2z_5 + 2z_4^2z_0 + 2z_4z_5^2 + 4z_4^2z_5 - z_2^2 - z_1^2 - 3z_4^2 + 2z_4^3 - z_5^2 \\ & + 6z_2z_1z_4 + 2z_2z_1z_3 + 2z_2z_4z_3 + 2z_2z_1z_0 + 4z_2z_1z_5 + 4z_2z_4z_0 + 4z_1z_4z_3 \\ & + 6z_2z_4z_5 + 2z_1z_4z_0 + 6z_1z_4z_5 + 2z_2z_3z_0 + 2z_2z_3z_5 + 2z_1z_3z_0 + 2z_2z_0z_5 \\ & + 2z_1z_3z_5 + 2z_4z_3z_0 + 2z_4z_3z_5 + 2z_1z_0z_5 + 2z_4z_0z_5 \end{aligned} - \frac{2z_2z_1 - z_1 - 2z_4 - z_3 - z_0 - 2z_5 - z_2 + 2z_2z_4 + 2z_1z_4 + 2z_2z_0 + 2z_1z_3 + 2z_2z_5}{2z_4z_3 + 2z_1z_5 + 2z_4z_0 + 4z_4z_5 + 2z_3z_0 + 2z_3z_5 + 2z_0z_5 + 2z_4^2 + 2z_5^2}.$$

We chose: $\mathbf{z} = \left(\frac{1}{10}, \frac{1}{10}, \frac{1}{10}, \frac{1}{10}, \frac{1}{10}, \frac{1}{10}\right)$.

Proposition

Assume $k_m : \mathbb{R}^{d_m} \times \mathbb{R}^{d_m} \rightarrow \mathbb{R}$ are continuous, translation-invariant kernels. Then the followings are equivalent:

- (i) $(k_m)_{m=1}^M$ -s are characteristic.
- (ii) $\otimes_{m=1}^M k_m$: \otimes_0 -characteristic.
- (iii) $\otimes_{m=1}^M k_m$: \otimes -characteristic.
- (iv) $\otimes_{m=1}^M k_m$: \mathcal{I} -characteristic.
- (v) $\otimes_{m=1}^M k_m$: characteristic.

Proposition

Assume $k_m : \mathbb{R}^{d_m} \times \mathbb{R}^{d_m} \rightarrow \mathbb{R}$ are continuous, translation-invariant kernels. Then the followings are equivalent:

- (i) $(k_m)_{m=1}^M$ -s are characteristic.
- (ii) $\otimes_{m=1}^M k_m$: \otimes_0 -characteristic.
- (iii) $\otimes_{m=1}^M k_m$: \otimes -characteristic.
- (iv) $\otimes_{m=1}^M k_m$: \mathcal{I} -characteristic.
- (v) $\otimes_{m=1}^M k_m$: characteristic.

Proof idea: We already know

$$(v) \Rightarrow (iii) \Rightarrow (ii) \Leftrightarrow (i), \quad (v) \Rightarrow (iv) \Rightarrow (i).$$

Remains: (i) \Rightarrow (v).

$$(k_m)_{m=1}^M: \text{characteristic} \Rightarrow \otimes_{m=1}^M k_m: \text{characteristic}$$

- Bochner theorem and *supp*-characterization:

$$k_m(x_m, x'_m) = \int_{\mathbb{R}^{d_m}} e^{-i\langle \omega_m, x_m - x'_m \rangle} d\Lambda_m(\omega_m), \quad x_m, x'_m \in \mathbb{R}^{d_m},$$

where $\text{supp}(\Lambda_m) = \mathbb{R}^{d_m}$.

$$(k_m)_{m=1}^M: \text{characteristic} \Rightarrow \otimes_{m=1}^M k_m: \text{characteristic}$$

- Bochner theorem and *supp*-characterization:

$$k_m(x_m, x'_m) = \int_{\mathbb{R}^{d_m}} e^{-i\langle \omega_m, x_m - x'_m \rangle} d\Lambda_m(\omega_m), \quad x_m, x'_m \in \mathbb{R}^{d_m},$$

where $\text{supp}(\Lambda_m) = \mathbb{R}^{d_m}$.

- Tensor kernel:

$$\left(\otimes_{m=1}^M k_m \right) (x, x') = \int_{\mathbb{R}^d} e^{-i\langle \omega, x - x' \rangle} d\Lambda(\omega), \quad \Lambda := \otimes_{m=1}^M \Lambda_m.$$

$$(k_m)_{m=1}^M: \text{characteristic} \Rightarrow \otimes_{m=1}^M k_m: \text{characteristic}$$

- Bochner theorem and supp -characterization:

$$k_m(x_m, x'_m) = \int_{\mathbb{R}^{d_m}} e^{-i\langle \omega_m, x_m - x'_m \rangle} d\Lambda_m(\omega_m), \quad x_m, x'_m \in \mathbb{R}^{d_m},$$

where $\text{supp}(\Lambda_m) = \mathbb{R}^{d_m}$.

- Tensor kernel:

$$\left(\otimes_{m=1}^M k_m \right) (x, x') = \int_{\mathbb{R}^d} e^{-i\langle \omega, x - x' \rangle} d\Lambda(\omega), \quad \Lambda := \otimes_{m=1}^M \Lambda_m.$$

- $\text{supp}(\Lambda) = \times_{m=1}^M \text{supp}(\Lambda_m) = \times_{m=1}^M \mathbb{R}^{d_m} = \mathbb{R}^d$.

c_0 -universality of $\otimes_{m=1}^M k_m$

We saw: for $M \geq 3$

$$(k_m)_{m=1}^M \text{ characteristic} \Rightarrow \otimes_{m=1}^M k_m: \mathcal{I}\text{-characteristic.}$$

Proposition

Assume $(k_m)_{m=1}^M$ are c_0 -kernels on LCP spaces \mathcal{X}_m . Then

$$\otimes_{m=1}^M k_m: c_0\text{-universal} \Leftrightarrow (k_m)_{m=1}^M \text{ are } c_0\text{-universal.}$$

Checking:

- \mathcal{X}_m : LCP $\Rightarrow \times_{m=1}^M \mathcal{X}_m$: LCP.
- k_m : c_0 -kernel $\Rightarrow \otimes_{m=1}^M k_m$: c_0 -kernel.

The tricky direction: if $(k_m)_{m=1}^M$ are c_0 -universal . . .

Goal: injectivity of $\mu = \mu_{\bigotimes_{m=1}^M k_m}$ on $\mathcal{M}_b(\mathcal{X})$, i.e.

$$\underbrace{\mu(\mathbb{F})}_{\int_{\mathcal{X}} \bigotimes_{m=1}^M k_m(\cdot, x_m) d\mathbb{F}(x)} = 0 \stackrel{?}{\Rightarrow} \mathbb{F} = 0, \quad \mathbb{F} \in \mathcal{M}_b(\mathcal{X}).$$

The tricky direction: if $(k_m)_{m=1}^M$ are c_0 -universal . . .

Goal: injectivity of $\mu = \mu_{\otimes_{m=1}^M k_m}$ on $\mathcal{M}_b(\mathcal{X})$, i.e.

$$\underbrace{\mu(\mathbb{F})}_{\int_{\mathcal{X}} \otimes_{m=1}^M k_m(\cdot, x_m) d\mathbb{F}(x)} = 0 \stackrel{?}{\Rightarrow} \mathbb{F} = 0, \quad \mathbb{F} \in \mathcal{M}_b(\mathcal{X}).$$

To get $\mathbb{F} = 0$ it is enough:

$$\mathbb{F}\left(\times_{m=1}^M B_m\right) = 0, \quad \forall B_m \in \mathcal{B}(\mathcal{X}_m).$$

Proof idea

$$\otimes_{m=1}^M \mathcal{H}_{k_m} \ni 0 = \int_{\mathcal{X}} \otimes_{m=1}^M k_m(\cdot, x_m) d\mathbb{F}(x),$$

,

$$\mathbb{R} \ni 0 = \mathbb{F} \left(\times_{m=1}^M B_m \right) = \int_{\mathcal{X}} \times_{m=1}^M \chi_{B_m}(x_m) d\mathbb{F}(x), \forall B_m.$$

Proof idea

$$\otimes_{m=1}^M \mathcal{H}_{k_m} \ni 0 = \int_{\mathcal{X}} \otimes_{m=1}^M k_m(\cdot, x_m) d\mathbb{F}(x),$$

$$\otimes_{m=J+1}^M \mathcal{H}_{k_m} \ni 0 = \int_{\mathcal{X}} \prod_{m=1}^J \chi_{B_m}(x_m) \otimes_{m=J+1}^M k_m(\cdot, x_m) d\mathbb{F}(x), \forall B_m,$$

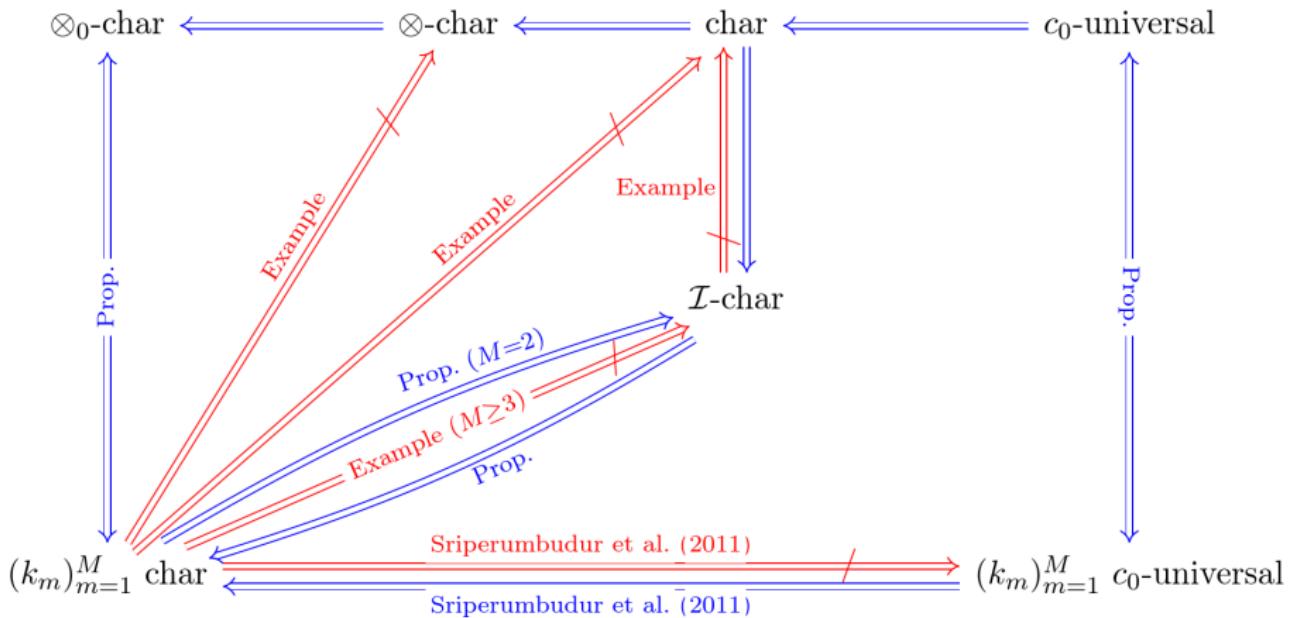
$$\mathbb{R} \ni 0 = \mathbb{F} \left(\times_{m=1}^M B_m \right) = \int_{\mathcal{X}} \times_{m=1}^M \chi_{B_m}(x_m) d\mathbb{F}(x), \forall B_m.$$

Proof idea

$$\begin{aligned}\otimes_{m=1}^M \mathcal{H}_{k_m} \ni 0 &= \int_{\mathcal{X}} \otimes_{m=1}^M k_m(\cdot, x_m) d\mathbb{F}(x), \\ \otimes_{m=J+1}^M \mathcal{H}_{k_m} \ni 0 &= \int_{\mathcal{X}} \prod_{m=1}^J \chi_{B_m}(x_m) \otimes_{m=J+1}^M k_m(\cdot, x_m) d\mathbb{F}(x), \forall B_m, \\ \mathbb{R} \ni 0 &= \mathbb{F} \left(\times_{m=1}^M B_m \right) = \int_{\mathcal{X}} \times_{m=1}^M \chi_{B_m}(x_m) d\mathbb{F}(x), \forall B_m.\end{aligned}$$

We proceed by induction ($J = 0, \dots, M$).

Visual Illustration



- Divergence & independence measures on kernel-endowed domains:
 - Maximum mean discrepancy,
 - Hilbert Schmidt independence criterion.
- Kernel: $k = \otimes_{m=1}^M k_m$.
- \mathcal{F} -ispd \Rightarrow
 - various characteristic properties.
 - relations & expressed in terms of k_m -s.

Thank you for the attention!

Acks: A part of the work was carried out while BKS was visiting ZSz at CMAP, École Polytechnique. BKS is supported by NSF-DMS-1713011.

Bochner integral: quick summary

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

Bochner integral: quick summary

- Given:
 - $(\mathcal{X}, \mathcal{A}, \mu)$: 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

- Given:
 - $(\mathcal{X}, \mathcal{A}, \mu)$: 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 $\left\| \int_{\mathcal{X}} f d\mu \right\|_B \leq \int_{\mathcal{X}} \|f\|_B d\mu$. ('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

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

In short

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

-  Berlinet, A. and Thomas-Agnan, C. (2004).
Reproducing Kernel Hilbert Spaces in Probability and Statistics.
Kluwer.
-  Fukumizu, K., Gretton, A., Sun, X., and Schölkopf, B. (2008).
Kernel measures of conditional dependence.
In *Neural Information Processing Systems (NIPS)*, pages 498–496.
-  Fukumizu, K., Song, L., and Gretton, A. (2013).
Kernel Bayes' rule: Bayesian inference with positive definite kernels.
Journal of Machine Learning Research, 14:3753–3783.
-  Gretton, A. (2015).
A simpler condition for consistency of a kernel independence test.
Technical report, University College London.
(<http://arxiv.org/abs/1501.06103>).

-  Gretton, A., Borgwardt, K. M., Rasch, M. J., Schölkopf, B., and Smola, A. (2012).
A kernel two-sample test.
Journal of Machine Learning Research, 13:723–773.
-  Gretton, A., Bousquet, O., Smola, A., and Schölkopf, B. (2005).
Measuring statistical dependence with Hilbert-Schmidt norms.
In *Algorithmic Learning Theory (ALT)*, pages 63–78.
-  Gretton, A., Fukumizu, K., Teo, C. H., Song, L., Schölkopf, B., and Smola, A. J. (2008).
A kernel statistical test of independence.
In *Neural Information Processing Systems (NIPS)*, pages 585–592.
-  Kusano, G., Fukumizu, K., and Hiraoka, Y. (2016).
Persistence weighted Gaussian kernel for topological data analysis.

In *International Conference on Machine Learning (ICML)*, pages 2004–2013.

-  Lloyd, J. R., Duvenaud, D., Grosse, R., Tenenbaum, J. B., and Ghahramani, Z. (2014).

Automatic construction and natural-language description of nonparametric regression models.

In *AAAI Conference on Artificial Intelligence*, pages 1242–1250.

-  Lyons, R. (2013).

Distance covariance in metric spaces.

The Annals of Probability, 41:3284–3305.

-  Micchelli, C. A., Xu, Y., and Zhang, H. (2006).

Universal kernels.

Journal of Machine Learning Research, 7:2651–2667.

-  Mooij, J. M., Peters, J., Janzing, D., Zscheischler, J., and Schölkopf, B. (2016).

Distinguishing cause from effect using observational data:
Methods and benchmarks.

Journal of Machine Learning Research, 17:1–102.

 Muandet, K., Fukumizu, K., Dinuzzo, F., and Schölkopf, B. (2011).

Learning from distributions via support measure machines.
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.

 Park, M., Jitkrittum, W., and Sejdinovic, D. (2016).

K2-ABC: Approximate Bayesian computation with kernel embeddings.

In *International Conference on Artificial Intelligence and Statistics (AISTATS; PMLR)*, volume 51, pages 51:398–407.

-  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).
-  Schölkopf, B., Muandet, K., Fukumizu, K., Harmeling, S., and Peters, J. (2015).
Computing functions of random variables via reproducing kernel Hilbert space representations.
Statistics and Computing, 25(4):755–766.
-  Sejdinovic, D., Sriperumbudur, B. K., Gretton, A., and Fukumizu, K. (2013).
Equivalence of distance-based and RKHS-based statistics in hypothesis testing.
Annals of Statistics, 41:2263–2291.

-  Song, L., Gretton, A., Bickson, D., Low, Y., and Guestrin, C. (2011).
Kernel belief propagation.
In International Conference on Artificial Intelligence and Statistics (AISTATS), pages 707–715.
-  Song, L., Smola, A., Gretton, A., Bedo, J., and Borgwardt, K. (2012).
Feature selection via dependence maximization.
Journal of Machine Learning Research, 13:1393–1434.
-  Sriperumbudur, B. K., Gretton, A., Fukumizu, K., Schölkopf, B., and Lanckriet, G. R. (2010).
Hilbert space embeddings and metrics on probability measures.
Journal of Machine Learning Research, 11:1517–1561.
-  Steinwart, I. (2001).
On the influence of the kernel on the consistency of support vector machines.

-  Szabó, Z., Sriperumbudur, B., Póczos, B., and Gretton, A. (2016).

Learning theory for distribution regression.

Journal of Machine Learning Research, 17(152):1–40.

-  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.

-  Yamada, M., Umezu, Y., Fukumizu, K., and Takeuchi, I. (2016).

Post selection inference with kernels.

Technical report.

(<https://arxiv.org/abs/1610.03725>).

-  Zhang, K., Schölkopf, B., Muandet, K., and Wang, Z. (2013).

Domain adaptation under target and conditional shift.

