

# Kernel Methods, Divergence and Independence Measures, Hypothesis Testing

Zoltán Szabó – CMAP, École Polytechnique

Summer School on Data Science for Document Analysis and Understanding  
INRIA, Palaiseau  
July 18 & 20, 2018

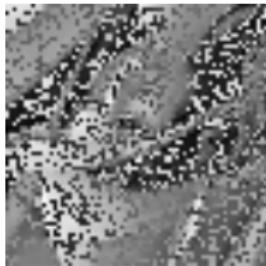
# Outline

- Applications:
  - Information theoretical objectives.
  - Testing.
- Classical information theory:  $\mathbb{R}^d \xrightarrow{\text{diverse set of domains}}$
- Kernels, RKHS:
  - Linear → non-linear techniques.
  - Classification, regression, dimensionality reduction.
  - KCCA, MMD, HSIC.
- Hypothesis testing.

# Information Theoretical Objectives

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

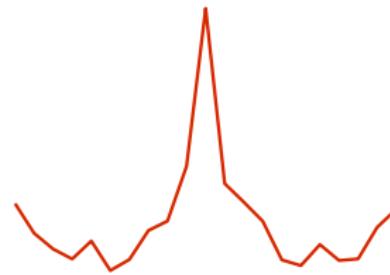
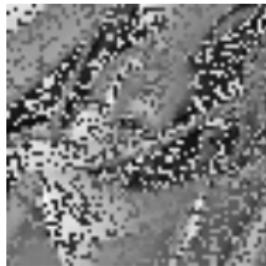
Given two images:



**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:  $\mathbf{y}_{\text{ref}}$ ,
- test image:  $\mathbf{y}_{\text{test}}$ ,
- possible transformations:  $\Theta$ .

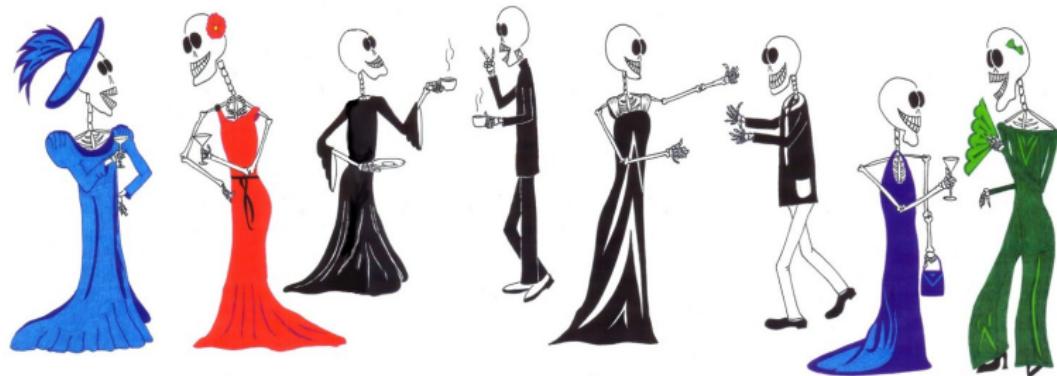
Objective:

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

In the example:  $I=KCCA$ .

Cocktail party problem:

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



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

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

$$\mathbf{y} = \mathbf{Wx} = \left[ \mathbf{y}^1; \dots; \mathbf{y}^M \right],$$

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

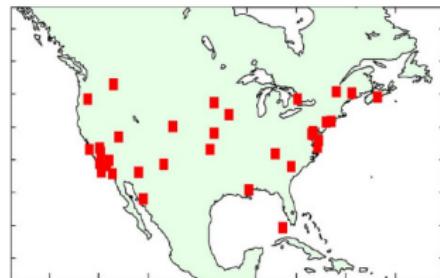
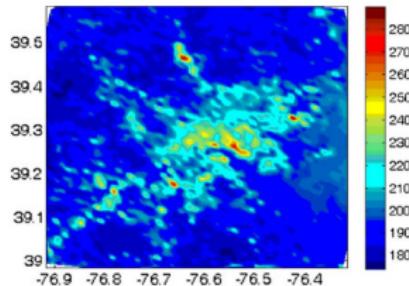
# Distribution regression

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

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



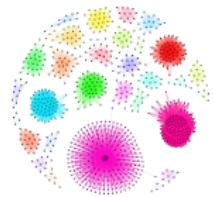
- Prediction using labelled bags:
  - bag := multi-spectral satellite measurements over an area,
  - label := local aerosol value.



# Objects in the bags



time series

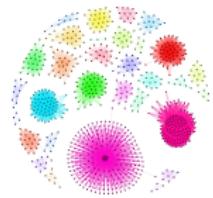


- Examples:
  - time-series modelling: user = set of **time-series**,
  - computer vision: image = collection of patch **vectors**,
  - NLP: corpus = bag of **documents**,
  - network analysis: group of people = bag of friendship **graphs**, ...

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

# Regression on labelled bags

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

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

# Regression on labelled bags

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

- Estimator:

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

- Prediction:

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

# Regression on labelled bags

- Given:

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

- Estimator:

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

- Prediction:

$$\hat{y}(\hat{\mathbb{P}}) = \mathbf{g}^T (\mathbf{G} + \ell \lambda \mathbf{I})^{-1} \mathbf{y},$$

$$\mathbf{g} = [K(\mu_{\hat{\mathbb{P}}}, \mu_{\hat{\mathbb{P}}_i})], \mathbf{G} = [K(\mu_{\hat{\mathbb{P}}_i}, \mu_{\hat{\mathbb{P}}_j})], \mathbf{y} = [y_i].$$

Inner product of distributions

$$K(\mu_{\hat{\mathbb{P}}_i}, \mu_{\hat{\mathbb{P}}_j}) = ?$$

# 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]:

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

# Hypothesis Testing

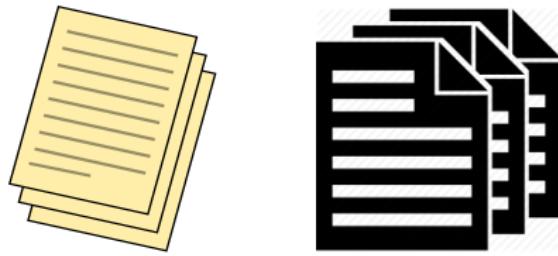
# Example-1 (2-sample testing): NLP

- Given: 2 categories of documents. Examples:
  - 1 Bayesian inference, neuroscience.
  - 2 adult attachment classes.
- Task:
  - test their distinguishability,
  - most discriminative words → interpretability.



# Example-1 (2-sample testing): NLP

- Given: 2 categories of documents. Examples:
  - 1 Bayesian inference, neuroscience.
  - 2 adult attachment classes.
- Task:
  - test their distinguishability,
  - most discriminative words → interpretability.



Do  $\{x_i\}$  and  $\{y_j\}$  come from the same distribution, i.e.  $P_x = P_y$ ?

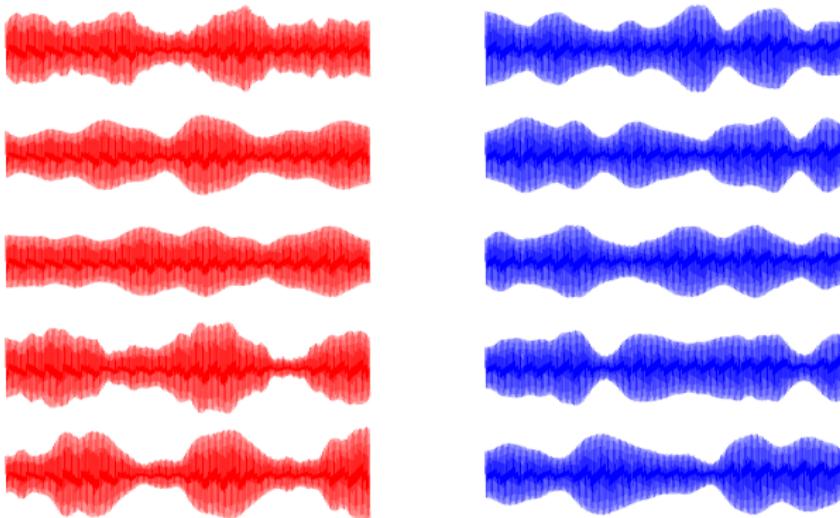
## Example-2 (2-sample testing): computer vision



- Given: two sets of faces (happy, angry).
- Task:
  - check if they are different,
  - determine the most discriminative features/regions.

## Example-3 (2-sample testing): audio

- Amplitude modulation:
  - simple technique to transmit voice over radio.
  - in the example: 2 songs.
- Fragments from song<sub>1</sub> ~  $\mathbb{P}_x$ , song<sub>2</sub> ~  $\mathbb{P}_y$ .



## Example: independence testing-1

- We are given **paired samples**. Task: test **independence**.
- Examples:
  - (song, year of release) pairs

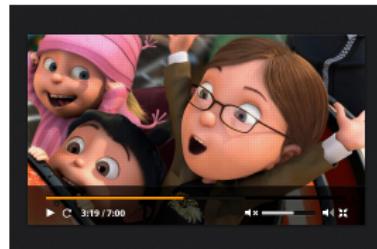


# Example: independence testing-1

- We are given **paired samples**. Task: test **independence**.
- Examples:
  - (song, year of release) pairs



- (video, caption) pairs

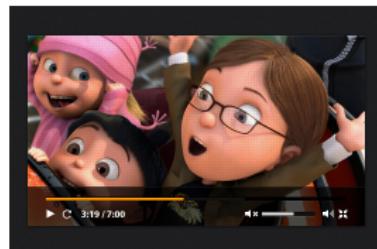


# Example: independence testing-1

- We are given **paired samples**. Task: test **independence**.
- Examples:
  - (song, year of release) pairs



- (video, caption) pairs



- $\{(x_i, y_i)\}_{i=1}^n \xrightarrow{?} \mathbb{P}_{xy} = \mathbb{P}_x \mathbb{P}_y$ .

## Example: independence testing-2

- How do we detect dependency? (**paired** samples)

*x<sub>1</sub>*: 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<sub>2</sub>*: 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<sub>1</sub>*: 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<sub>2</sub>*: 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.

...

## Example: independence testing-2

- 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.  $\mathbb{P}_{xy} = \mathbb{P}_x \mathbb{P}_y$ ?

## Example: goodness-of-fit testing

- Demo: criminal data analysis.
  - Given:
    - Density/model:  $p$ .



# Example: goodness-of-fit testing

- Demo: criminal data analysis.
- Given:
  - Density/model:  $p$ .
  - Samples:  $X = \{x_i\}_{i=1}^n \sim q$  (unknown).



## Example: goodness-of-fit testing

- Demo: criminal data analysis.
  - Given:
    - Density/model:  $p$ .
    - Samples:  $X = \{x_i\}_{i=1}^n \sim q$  (unknown).
  - Task: using  $p$ ,  $X$  test

$$H_0 : p = q, \text{ vs}$$

$$H_1 : p \neq q.$$



# 'Classical' information theory

- Kullback-Leibler divergence:

$$\text{KL}(\mathbb{P}, \mathbb{Q}) = \int_{\mathbb{R}^d} p(x) \log \left[ \frac{p(x)}{q(x)} \right] dx.$$

# 'Classical' information theory

- Kullback-Leibler divergence:

$$\text{KL}(\mathbb{P}, \mathbb{Q}) = \int_{\mathbb{R}^d} p(x) \log \left[ \frac{p(x)}{q(x)} \right] dx.$$

- Mutual information:

$$I(\mathbb{P}) = \text{KL}\left(\mathbb{P}, \otimes_{m=1}^M \mathbb{P}_m\right).$$

# 'Classical' information theory

- Kullback-Leibler divergence:

$$\text{KL}(\mathbb{P}, \mathbb{Q}) = \int_{\mathbb{R}^d} p(x) \log \left[ \frac{p(x)}{q(x)} \right] dx.$$

- Mutual information:

$$I(\mathbb{P}) = \text{KL}\left(\mathbb{P}, \otimes_{m=1}^M \mathbb{P}_m\right).$$

Properties:

①  $I(\mathbb{P}) \geq 0$ .

# 'Classical' information theory

- Kullback-Leibler divergence:

$$\text{KL}(\mathbb{P}, \mathbb{Q}) = \int_{\mathbb{R}^d} p(x) \log \left[ \frac{p(x)}{q(x)} \right] dx.$$

- Mutual information:

$$I(\mathbb{P}) = \text{KL}\left(\mathbb{P}, \otimes_{m=1}^M \mathbb{P}_m\right).$$

Properties:

- ①  $I(\mathbb{P}) \geq 0$ .
- ②  $I(\mathbb{P}) = 0 \Leftrightarrow \mathbb{P} = \otimes_{m=1}^M \mathbb{P}_m$ .

# 'Classical' information theory

- Kullback-Leibler divergence:

$$\text{KL}(\mathbb{P}, \mathbb{Q}) = \int_{\mathbb{R}^d} p(x) \log \left[ \frac{p(x)}{q(x)} \right] dx.$$

- Mutual information:

$$I(\mathbb{P}) = \text{KL}\left(\mathbb{P}, \otimes_{m=1}^M \mathbb{P}_m\right).$$

Properties:

- ①  $I(\mathbb{P}) \geq 0$ .
- ②  $I(\mathbb{P}) = 0 \Leftrightarrow \mathbb{P} = \otimes_{m=1}^M \mathbb{P}_m$ .

Alternatives: Rényi, Tsallis,  $L^2$  divergence...  $\mathcal{X} = \mathbb{R}^d$ .

Euclidean space → inner product → kernel

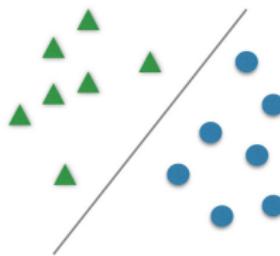
Extension of  $k(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y}$  leads to kernels.

Euclidean space → inner product → kernel

Extension of  $k(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y}$  leads to kernels. Why?

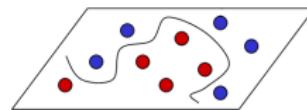
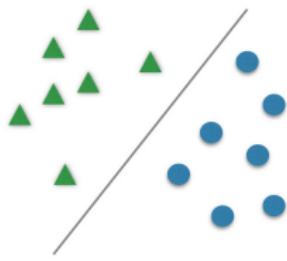
Extension of  $k(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y}$  leads to kernels. Why?

- Classification (SVM):



Extension of  $k(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y}$  leads to kernels. Why?

- Classification (SVM):

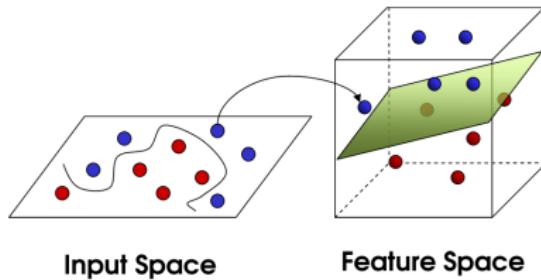
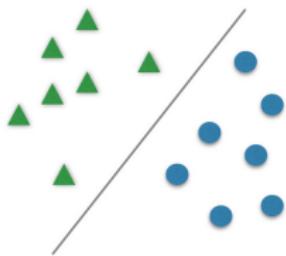


Input Space

# Euclidean space → inner product → kernel

Extension of  $k(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y}$  leads to kernels. Why?

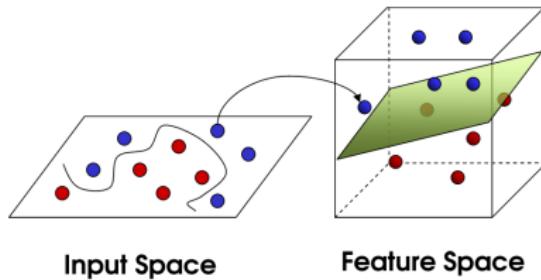
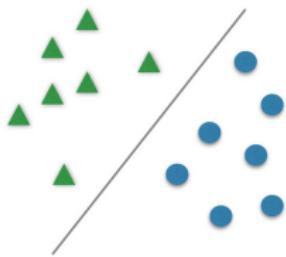
- Classification (SVM):



# Euclidean space $\rightarrow$ inner product $\rightarrow$ kernel

Extension of  $k(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y}$  leads to kernels. Why?

- Classification (SVM):



- Representation of distributions:

$$\mathbb{P} \mapsto \mathbb{E}_{\mathbf{x} \sim \mathbb{P}} \varphi(\mathbf{x}).$$

Example:  $\varphi(\mathbf{x}) = \mathbf{x}$ : mean.

# Distribution representation via functions

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

# Distribution representation via functions

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

- Cdf:

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

# Distribution representation via functions

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

- Cdf:

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

- Characteristic function:

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

# Distribution representation via functions

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

- Cdf:

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

- Characteristic function:

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

- Moment generating function:

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

# Distribution representation via functions

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

- Cdf:

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

- Characteristic function:

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

- Moment generating function:

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

Trick

$\varphi$ : on any kernel-endowed domain!

- **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]  $\xrightarrow{\text{spec.}}$  **permutations** [Jiao and Vert, 2016],
- **graphs** [Vishwanathan et al., 2010, Kondor and Pan, 2016].

# Objects of Interest

'KL divergence & mutual information' on kernel-endowed domains.

- Mean embedding:

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

# Objects of Interest

'KL divergence & mutual information' on kernel-endowed domains.

- Mean embedding:

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

- Maximum mean discrepancy:

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

# Objects of Interest

'KL divergence & mutual information' on kernel-endowed domains.

- Mean embedding:

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

- Maximum mean discrepancy:

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

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

$$\text{HSIC}_k(\mathbb{P}) := \text{MMD}_k \left( \mathbb{P}, \otimes_{m=1}^M \mathbb{P}_m \right).$$

- Applications:
  - two-sample testing [Borgwardt et al., 2006, Gretton et al., 2012],
  - domain adaptation [Zhang et al., 2013], -generalization [Blanchard et al., 2017],
  - kernel Bayesian inference [Song et al., 2011, Fukumizu et al., 2013]
  - approximate Bayesian computation [Park et al., 2016], probabilistic programming [Schölkopf et al., 2015],
  - model criticism [Lloyd et al., 2014, Kim et al., 2016], goodness-of-fit [Balasubramanian et al., 2017],
  - distribution classification [Muandet et al., 2011, Lopez-Paz et al., 2015], [Zaheer et al., 2017], distribution regression [Szabó et al., 2016], [Law et al., 2018],
  - topological data analysis [Kusano et al., 2016].
- Review [Muandet et al., 2017].

MMD with  $k = \otimes_{m=1}^M k_m$ :

$$\text{HSIC}_k(\mathbb{P}) = \text{MMD}_k\left(\mathbb{P}, \otimes_{m=1}^M \mathbb{P}_m\right).$$

MMD with  $k = \otimes_{m=1}^M k_m$ :

$$\text{HSIC}_k(\mathbb{P}) = \text{MMD}_k\left(\mathbb{P}, \otimes_{m=1}^M \mathbb{P}_m\right).$$

### Applications :

- blind source separation [Gretton et al., 2005a],
- feature selection [Song et al., 2012], post selection inference [Yamada et al., 2018],
- independence testing [Gretton et al., 2008], causal inference [Mooij et al., 2016, Pfister et al., 2017, Strobl et al., 2017].

MMD with  $k = \otimes_{m=1}^M k_m$ :

$$\text{HSIC}_k(\mathbb{P}) = \text{MMD}_k\left(\mathbb{P}, \otimes_{m=1}^M \mathbb{P}_m\right).$$

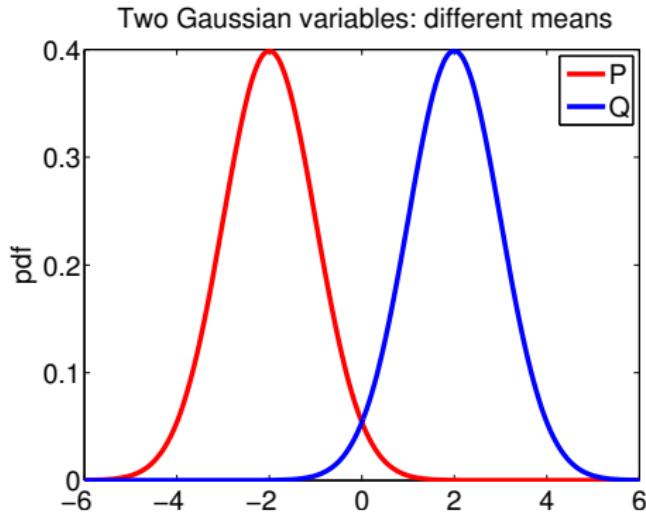
### Applications :

- blind source separation [Gretton et al., 2005a],
- feature selection [Song et al., 2012], post selection inference [Yamada et al., 2018],
- independence testing [Gretton et al., 2008], causal inference [Mooij et al., 2016, Pfister et al., 2017, Strobl et al., 2017].

MMD, HSIC: Easy to Estimate!

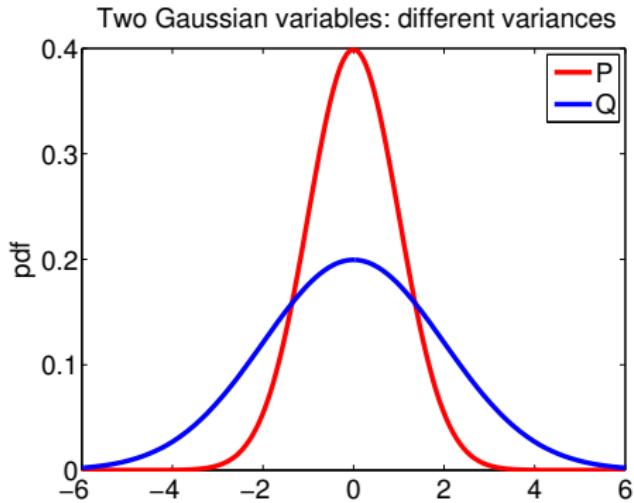
# Representations of distributions: $\mathbb{E}X$

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



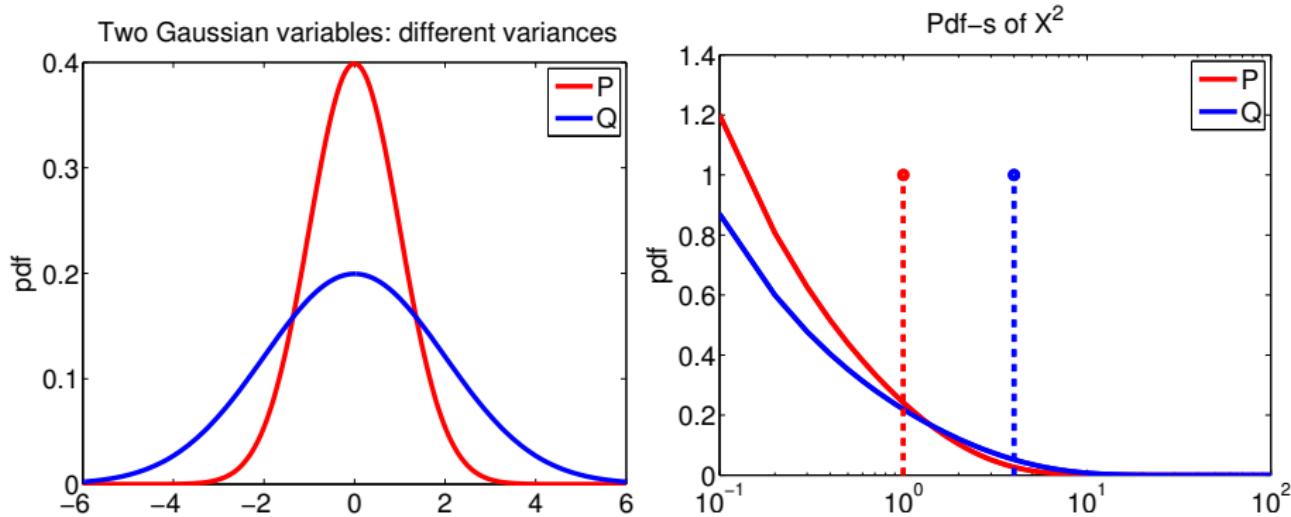
# Representations of distributions: $\mathbb{E}X^2$

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



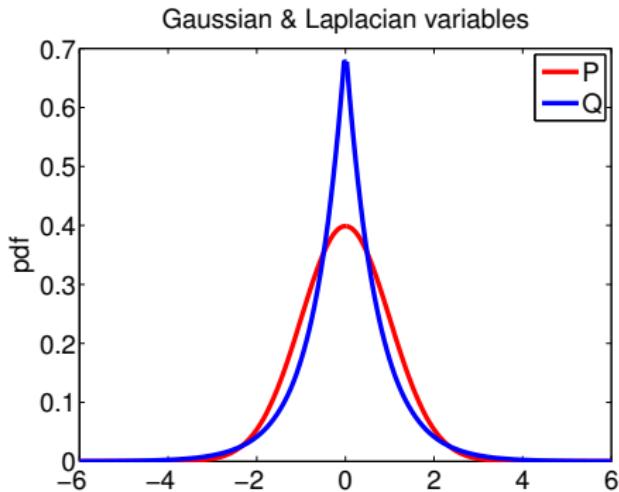
# Representations of distributions: $\mathbb{E}X^2$

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



## 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.



$\varphi(\mathbf{x}) = e^{i\langle \cdot, \mathbf{x} \rangle}$ : characteristic function,  $\mathcal{X} = \mathbb{R}^d$ .

# Kernels: why? – continued

Idea:  $\mathbb{P}_{xy} \mapsto C_{xy}$ .

- Covariance matrix

$$C_{xy} = \mathbb{E}_{xy} \left[ (x - \mathbb{E}x) (y - \mathbb{E}y)^T \right]$$

# Kernels: why? – continued

Idea:  $\mathbb{P}_{xy} \mapsto C_{xy}$ .

- Covariance matrix

$$C_{xy} = \mathbb{E}_{xy} \left[ (x - \mathbb{E}x) (y - \mathbb{E}y)^T \right],$$

$$S = \|C_{xy}\|_F$$

# Kernels: why? – continued

Idea:  $\mathbb{P}_{xy} \mapsto C_{xy}$ .

- Covariance matrix

$$C_{xy} = \mathbb{E}_{xy} \left[ (x - \mathbb{E}x) (y - \mathbb{E}y)^T \right],$$

$$S = \|C_{xy}\|_F \stackrel{?}{=} 0$$

# Kernels: why? – continued

Idea:  $\mathbb{P}_{xy} \mapsto C_{xy}$ .

- Covariance matrix

$$C_{xy} = \mathbb{E}_{xy} \left[ (x - \mathbb{E}x) (y - \mathbb{E}y)^T \right],$$

$$S = \|C_{xy}\|_F \stackrel{?}{=} 0 \leftrightarrow \text{linear dependence}.$$

# Kernels: why? – continued

Idea:  $\mathbb{P}_{xy} \mapsto C_{xy}$ .

- Covariance matrix

$$C_{xy} = \mathbb{E}_{xy} \left[ (x - \mathbb{E}x) (y - \mathbb{E}y)^T \right],$$
$$S = \|C_{xy}\|_F \stackrel{?}{=} 0 \leftrightarrow \text{linear dependence.}$$

- Covariance operator: take features of  $x$  and  $y$

$$C_{xy} = \mathbb{E}_{xy} \left[ \underbrace{(\varphi(x) - \mathbb{E}_x \varphi(x))}_{\text{centering in feature space}} \otimes (\psi(y) - \mathbb{E}_y \psi(y)) \right]$$

# Kernels: why? – continued

Idea:  $\mathbb{P}_{xy} \mapsto C_{xy}$ .

- Covariance matrix

$$C_{xy} = \mathbb{E}_{xy} \left[ (x - \mathbb{E}x) (y - \mathbb{E}y)^T \right],$$
$$S = \|C_{xy}\|_F \stackrel{?}{=} 0 \leftrightarrow \text{linear dependence.}$$

- Covariance operator: take features of  $x$  and  $y$

$$C_{xy} = \mathbb{E}_{xy} \left[ \underbrace{(\varphi(x) - \mathbb{E}_x \varphi(x))}_{\text{centering in feature space}} \otimes (\psi(y) - \mathbb{E}_y \psi(y)) \right],$$
$$S = \|C_{xy}\|_{HS} =: \text{HSIC}(\mathbb{P}_{xy}).$$

We capture non-linear dependencies via  $\varphi, \psi!$

- Kernel ( $k$ ), RKHS ( $\mathcal{H}_k$ ) → classification, regression (ridge), PCA.
- Mean embedding ( $\mu_{\mathbb{P}}$ ): characteristic property, universality,
- $\otimes_m k_m$ ,  $\otimes_m \mathcal{H}_{k_m}$ , covariance operator,
- MMD, HSIC, KCCA,
- with applications.

# Kernels & Friends

# Kernel: similarity between features

- Given:  $x$  and  $x'$  objects (images or texts).

# Kernel: similarity between features

- Given:  $x$  and  $x'$  objects (images or texts).
- Question: how similar they are?

# Kernel: similarity between features

- Given:  $x$  and  $x'$  objects (images or texts).
- Question: how similar they are?
- Define **features** of the objects:

$$\begin{aligned}\varphi(x) &: \text{features of } x, \\ \varphi(x') &: \text{features of } x'.\end{aligned}$$

- Kernel:** inner product of these features

$$k(x, x') := \langle \varphi(x), \varphi(x') \rangle.$$

# Kernel examples on $\mathbb{R}^d$ ( $\gamma > 0, p \in \mathbb{Z}^+$ )

- Polynomial kernel:

$$k(\mathbf{x}, \mathbf{y}) = (\langle \mathbf{x}, \mathbf{y} \rangle + \gamma)^p.$$

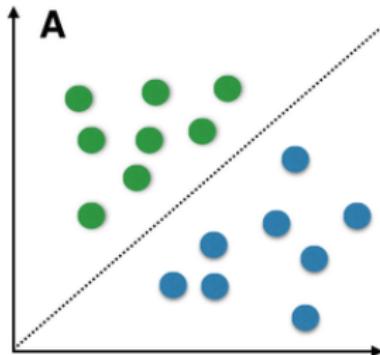
- Gaussian kernel:

$$k(\mathbf{x}, \mathbf{y}) = e^{-\gamma \|\mathbf{x} - \mathbf{y}\|_2^2}.$$

Non-linear features: why?

# Classification motivation: linear separability

Idealized situation

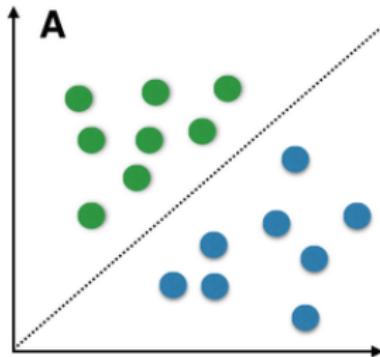


Decision surface:

$$\{\mathbf{x} : \langle \mathbf{w}, \mathbf{x} \rangle = 0\}$$

# Classification motivation: linear separability

Idealized situation



Decision surface:

$$\{\mathbf{x} : \langle \mathbf{w}, \mathbf{x} \rangle = 0\} \Rightarrow$$

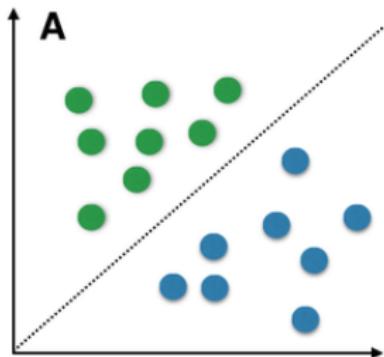
classes:

$$\{\mathbf{x} : \langle \mathbf{w}, \mathbf{x} \rangle \geq 0\}$$

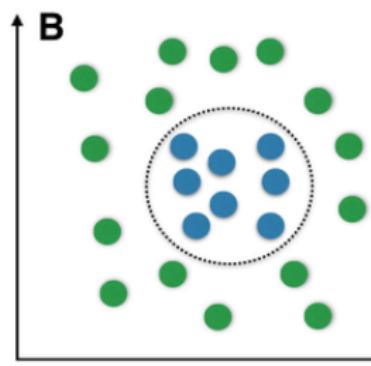
$$\{\mathbf{x} : \langle \mathbf{w}, \mathbf{x} \rangle < 0\}$$

# Classification motivation: non-linear separability

Idealized situation



Real world



Decision surface (left):

$$\{\mathbf{x} : \langle \mathbf{w}, \mathbf{x} \rangle = 0\} \Rightarrow$$

classes:

$$\{\mathbf{x} : \langle \mathbf{w}, \mathbf{x} \rangle \geq 0\}$$

$$\{\mathbf{x} : \langle \mathbf{w}, \mathbf{x} \rangle < 0\}.$$

# Non-linear separability – continued

On the ellipse

$$\left\{ \mathbf{x} : \frac{(x_1 - c_1)^2}{a^2} + \frac{(x_2 - c_2)^2}{b^2} = 1 \right\}$$

# Non-linear separability – continued

On the **ellipse**, outside

$$\left\{ \mathbf{x} : \frac{(x_1 - c_1)^2}{a^2} + \frac{(x_2 - c_2)^2}{b^2} = 1 \right\},$$
$$\left\{ \mathbf{x} : \frac{(x_1 - c_1)^2}{a^2} + \frac{(x_2 - c_2)^2}{b^2} > 1 \right\}$$

# Non-linear separability – continued

On the **ellipse**, **outside**, **inside**:

$$\left\{ \mathbf{x} : \frac{(x_1 - c_1)^2}{a^2} + \frac{(x_2 - c_2)^2}{b^2} = 1 \right\},$$

$$\left\{ \mathbf{x} : \frac{(x_1 - c_1)^2}{a^2} + \frac{(x_2 - c_2)^2}{b^2} > 1 \right\},$$

$$\left\{ \mathbf{x} : \frac{(x_1 - c_1)^2}{a^2} + \frac{(x_2 - c_2)^2}{b^2} < 1 \right\}.$$

# Non-linear separability – continued

On the **ellipse**, **outside**, **inside**:

$$\left\{ \mathbf{x} : \frac{(x_1 - c_1)^2}{a^2} + \frac{(x_2 - c_2)^2}{b^2} = 1 \right\},$$

$$\left\{ \mathbf{x} : \frac{(x_1 - c_1)^2}{a^2} + \frac{(x_2 - c_2)^2}{b^2} > 1 \right\},$$

$$\left\{ \mathbf{x} : \frac{(x_1 - c_1)^2}{a^2} + \frac{(x_2 - c_2)^2}{b^2} < 1 \right\}.$$

With polynomial feature:  $\varphi(\mathbf{x}) = (x_1^2, x_1, 1, x_2^2, x_2)$ :

- Decision surface:  $\{\mathbf{x} : \langle \mathbf{w}, \varphi(\mathbf{x}) \rangle = 0\}$ .

# Non-linear separability – continued

On the **ellipse**, **outside**, **inside**:

$$\left\{ \mathbf{x} : \frac{(x_1 - c_1)^2}{a^2} + \frac{(x_2 - c_2)^2}{b^2} = 1 \right\},$$

$$\left\{ \mathbf{x} : \frac{(x_1 - c_1)^2}{a^2} + \frac{(x_2 - c_2)^2}{b^2} > 1 \right\},$$

$$\left\{ \mathbf{x} : \frac{(x_1 - c_1)^2}{a^2} + \frac{(x_2 - c_2)^2}{b^2} < 1 \right\}.$$

With polynomial feature:  $\varphi(\mathbf{x}) = (x_1^2, x_1, 1, x_2^2, x_2)$ :

- Decision surface:  $\{\mathbf{x} : \langle \mathbf{w}, \varphi(\mathbf{x}) \rangle = 0\}$ .
- Classes:  $\{\mathbf{x} : \langle \mathbf{w}, \varphi(\mathbf{x}) \rangle > 0\}$ ,  $\{\mathbf{x} : \langle \mathbf{w}, \varphi(\mathbf{x}) \rangle < 0\}$ .

# Quadratic & polynomial features

Still in  $\mathbb{R}^2$ :

$$\varphi(\mathbf{x}) = \left( x_1^2, \sqrt{2}x_1x_2, x_2^2 \right),$$

$$\langle \varphi(\mathbf{x}), \varphi(\mathbf{x}') \rangle = ?$$

# Quadratic & polynomial features

Still in  $\mathbb{R}^2$ :

$$\varphi(\mathbf{x}) = \left( x_1^2, \sqrt{2}x_1x_2, x_2^2 \right),$$
$$\langle \varphi(\mathbf{x}), \varphi(\mathbf{x}') \rangle = \left\langle \begin{bmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{bmatrix}, \begin{bmatrix} (x'_1)^2 \\ \sqrt{2}(x'_1)(x'_2) \\ (x'_2)^2 \end{bmatrix} \right\rangle$$

# Quadratic & polynomial features

Still in  $\mathbb{R}^2$ :

$$\begin{aligned}\varphi(\mathbf{x}) &= \left( x_1^2, \sqrt{2}x_1x_2, x_2^2 \right), \\ \langle \varphi(\mathbf{x}), \varphi(\mathbf{x}') \rangle &= \left\langle \begin{bmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{bmatrix}, \begin{bmatrix} (x'_1)^2 \\ \sqrt{2}(x'_1)(x'_2) \\ (x'_2)^2 \end{bmatrix} \right\rangle \\ &= x_1^2(x'_1)^2 + \underbrace{\sqrt{2}\sqrt{2}}_2 x_1x_2(x'_1)(x'_2) + x_2^2(x'_2)^2\end{aligned}$$

# Quadratic & polynomial features

Still in  $\mathbb{R}^2$ :

$$\begin{aligned}\varphi(\mathbf{x}) &= \left( x_1^2, \sqrt{2}x_1x_2, x_2^2 \right), \\ \langle \varphi(\mathbf{x}), \varphi(\mathbf{x}') \rangle &= \left\langle \begin{bmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{bmatrix}, \begin{bmatrix} (x'_1)^2 \\ \sqrt{2}(x'_1)(x'_2) \\ (x'_2)^2 \end{bmatrix} \right\rangle \\ &= x_1^2(x'_1)^2 + \underbrace{\sqrt{2}\sqrt{2}}_2 x_1x_2(x'_1)(x'_2) + x_2^2(x'_2)^2 \\ &= (x_1x'_1 + x_2x'_2)^2\end{aligned}$$

# Quadratic & polynomial features

Still in  $\mathbb{R}^2$ :

$$\begin{aligned}\varphi(\mathbf{x}) &= \left( x_1^2, \sqrt{2}x_1x_2, x_2^2 \right), \\ \langle \varphi(\mathbf{x}), \varphi(\mathbf{x}') \rangle &= \left\langle \begin{bmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{bmatrix}, \begin{bmatrix} (x'_1)^2 \\ \sqrt{2}(x'_1)(x'_2) \\ (x'_2)^2 \end{bmatrix} \right\rangle \\ &= x_1^2(x'_1)^2 + \underbrace{\sqrt{2}\sqrt{2}}_2 x_1x_2(x'_1)(x'_2) + x_2^2(x'_2)^2 \\ &= (x_1x'_1 + x_2x'_2)^2 \\ &= \left\langle \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \begin{bmatrix} x'_1 \\ x'_2 \end{bmatrix} \right\rangle^2 = \langle \mathbf{x}, \mathbf{x}' \rangle^2 =: k(\mathbf{x}, \mathbf{x}').\end{aligned}$$

# Quadratic & polynomial features

Still in  $\mathbb{R}^2$ :

$$\begin{aligned}\varphi(\mathbf{x}) &= \left( x_1^2, \sqrt{2}x_1x_2, x_2^2 \right), \\ \langle \varphi(\mathbf{x}), \varphi(\mathbf{x}') \rangle &= \left\langle \begin{bmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{bmatrix}, \begin{bmatrix} (x'_1)^2 \\ \sqrt{2}(x'_1)(x'_2) \\ (x'_2)^2 \end{bmatrix} \right\rangle \\ &= x_1^2(x'_1)^2 + \underbrace{\sqrt{2}\sqrt{2}}_2 x_1x_2(x'_1)(x'_2) + x_2^2(x'_2)^2 \\ &= (x_1x'_1 + x_2x'_2)^2 \\ &= \left\langle \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \begin{bmatrix} x'_1 \\ x'_2 \end{bmatrix} \right\rangle^2 = \langle \mathbf{x}, \mathbf{x}' \rangle^2 =: k(\mathbf{x}, \mathbf{x}').\end{aligned}$$

$\langle \mathbf{x}, \mathbf{x}' \rangle^d = \langle \varphi(\mathbf{x}), \varphi(\mathbf{x}') \rangle$ :  $\varphi(\mathbf{x})$  =  $d$ -order polynomial.  $\Rightarrow$

# Quadratic & polynomial features

Still in  $\mathbb{R}^2$ :

$$\begin{aligned}\varphi(\mathbf{x}) &= \left( x_1^2, \sqrt{2}x_1x_2, x_2^2 \right), \\ \langle \varphi(\mathbf{x}), \varphi(\mathbf{x}') \rangle &= \left\langle \begin{bmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{bmatrix}, \begin{bmatrix} (x'_1)^2 \\ \sqrt{2}(x'_1)(x'_2) \\ (x'_2)^2 \end{bmatrix} \right\rangle \\ &= x_1^2(x'_1)^2 + \underbrace{\sqrt{2}\sqrt{2}}_2 x_1x_2(x'_1)(x'_2) + x_2^2(x'_2)^2 \\ &= (x_1x'_1 + x_2x'_2)^2 \\ &= \left\langle \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \begin{bmatrix} x'_1 \\ x'_2 \end{bmatrix} \right\rangle^2 = \langle \mathbf{x}, \mathbf{x}' \rangle^2 =: k(\mathbf{x}, \mathbf{x}').\end{aligned}$$

$\langle \mathbf{x}, \mathbf{x}' \rangle^d = \langle \varphi(\mathbf{x}), \varphi(\mathbf{x}') \rangle$ :  $\varphi(\mathbf{x})$  =  $d$ -order polynomial.  $\Rightarrow$  Explicit computation would be heavy!

## Maximum correlation: KCCA

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

$$Q(\mathbb{P}_{xy}) = \sup_{f \in \mathcal{F}, g \in \mathcal{G}} \text{corr}(f(x), g(y))$$

## Maximum correlation: KCCA

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

$$Q(\mathbb{P}_{xy}) = \sup_{f \in \mathcal{F}, g \in \mathcal{G}} \text{corr}(f(x), g(y)),$$

- $\mathcal{F} = C_b(\mathcal{X})$ ,  $\mathcal{G} = C_b(\mathcal{Y})$ : enough, but computationally...

## Maximum correlation: KCCA

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

$$Q(\mathbb{P}_{xy}) = \sup_{f \in \mathcal{F}, g \in \mathcal{G}} \text{corr}(f(x), g(y)),$$

- $\mathcal{F} = C_b(\mathcal{X})$ ,  $\mathcal{G} = C_b(\mathcal{Y})$ : enough, but computationally...
- Trick:  $\mathcal{H}_k$  dense in  $C_b(\mathcal{X})$ , similarly  $\mathcal{H}_\ell$  dense in  $C_b(\mathcal{Y})$ .
  - This universality: captures independence.
  - Computationally tractable.

# Kernels: why?

Linear → non-linear transition:

- Inner product of non-linear features:  $k(x, x')$ , implicit (highD: ✓).

Linear → non-linear transition:

- Inner product of non-linear features:  $k(x, x')$ , implicit (highD: ✓).
- Classification: non-linear separability.

# Kernels: why?

Linear → non-linear transition:

- Inner product of non-linear features:  $k(x, x')$ , implicit (highD: ✓).
- Classification: non-linear separability.
- Information theory:
  - $\sup_{f,g} \text{corr}(f(x), g(y))$ : can capture independence (KCCA),
  - $\mathbb{E}[\text{'rich' features}] \leftarrow \text{encodes distributions}$  (MMD).

# Kernels: why?

Linear → non-linear transition:

- Inner product of non-linear features:  $k(x, x')$ , implicit (highD: ✓).
- Classification: non-linear separability.
- Information theory:
  - $\sup_{f,g} \text{corr}(f(x), g(y))$ : can capture independence (KCCA),
  - $\mathbb{E}[\text{'rich' features}] \leftarrow \text{encodes distributions}$  (MMD).
- Computational tractability & accelerations.

Linear → non-linear transition:

- Inner product of non-linear features:  $k(x, x')$ , implicit (highD: ✓).
- Classification: non-linear separability.
- Information theory:
  - $\sup_{f,g} \text{corr}(f(x), g(y))$ : can capture independence (KCCA),
  - $\mathbb{E}[\text{'rich' features}] \leftarrow \text{encodes distributions}$  (MMD).
- Computational tractability & accelerations.
- Hilbert space: enables analysis.

- **ITE** toolbox (KL, MI, MMD, HSIC, KCCA & many more):  
<https://bitbucket.org/szzoli/ite-in-python/>

- **ITE** toolbox (KL, MI, MMD, HSIC, KCCA & many more):  
<https://bitbucket.org/szzoli/ite-in-python/>
- Linear-time hypothesis testing
  - **two-sample** (NIPS-2016, oral):  
<https://github.com/wittawatj/interpretable-test>

- **ITE** toolbox (KL, MI, MMD, HSIC, KCCA & many more):  
<https://bitbucket.org/szzoli/ite-in-python/>
- Linear-time hypothesis testing
  - **two-sample** (NIPS-2016, oral):  
<https://github.com/wittawatj/interpretable-test>
  - **independence** (ICML-2017):  
.../fsic-test

- **ITE** toolbox (KL, MI, MMD, HSIC, KCCA & many more):  
<https://bitbucket.org/szzoli/ite-in-python/>
- Linear-time hypothesis testing
  - **two-sample** (NIPS-2016, oral):  
<https://github.com/wittawatj/interpretable-test>
  - **independence** (ICML-2017):  
.../fsic-test
  - **goodness-of-fit** (NIPS-2017, best paper award):  
.../kernel-gof

# Kernels

- Given:  $\mathcal{X}$  set.
- $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is called **kernel** if  $\exists \varphi : \mathcal{X} \rightarrow \mathcal{H}$ ilbert space

$$k(a, b) = \langle \varphi(a), \varphi(b) \rangle_{\mathcal{H}}.$$

- Given:  $\mathcal{X}$  set.
- $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is called **kernel** if  $\exists \varphi : \mathcal{X} \rightarrow \mathcal{H}$ ilbert space

$$k(a, b) = \langle \varphi(a), \varphi(b) \rangle_{\mathcal{H}}.$$

- $\varphi/\mathcal{H}$  might not be unique:  $k(x, y) = xy$

$$\varphi_1(x) = x \in \mathbb{R}, \quad \varphi_2(x) = \frac{1}{\sqrt{2}}(x, x) \in \mathbb{R}^2.$$

- Given:  $\mathcal{X}$  set.
- $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is called **kernel** if  $\exists \varphi : \mathcal{X} \rightarrow \mathcal{H}$  Hilbert space

$$k(a, b) = \langle \varphi(a), \varphi(b) \rangle_{\mathcal{H}}.$$

- $\varphi/\mathcal{H}$  might not be unique:  $k(x, y) = xy$

$$\varphi_1(x) = x \in \mathbb{R}, \quad \varphi_2(x) = \frac{1}{\sqrt{2}}(x, x) \in \mathbb{R}^2.$$

- $\mathcal{H}$  intuition: vectors, inner product, complete ('no holes').

A bit of functional analysis follows  $\approx$  linalg, geometry!

# Vector space: $(V, +, \lambda \cdot)$

- Points = vectors.
- 2 operations:  $\mathbf{v}_1 + \mathbf{v}_2$ ,  $\lambda\mathbf{v}$  with the 'natural' properties

## Vector space: $(V, +, \lambda \cdot)$

- Points = vectors.
- 2 operations:  $\mathbf{v}_1 + \mathbf{v}_2$ ,  $\lambda \mathbf{v}$  with the 'natural' properties

$$(\mathbf{v}_1 + \mathbf{v}_2) + \mathbf{v}_3 = \mathbf{v}_1 + (\mathbf{v}_2 + \mathbf{v}_3), \text{ (associativity)}$$

# Vector space: $(V, +, \lambda \cdot)$

- Points = vectors.
- 2 operations:  $\mathbf{v}_1 + \mathbf{v}_2$ ,  $\lambda \mathbf{v}$  with the 'natural' properties

$$(\mathbf{v}_1 + \mathbf{v}_2) + \mathbf{v}_3 = \mathbf{v}_1 + (\mathbf{v}_2 + \mathbf{v}_3), \text{ (associativity)}$$

$$\mathbf{v}_1 + \mathbf{v}_2 = \mathbf{v}_2 + \mathbf{v}_1, \quad \text{ (commutativity)}$$

# Vector space: $(V, +, \lambda \cdot)$

- Points = vectors.
- 2 operations:  $\mathbf{v}_1 + \mathbf{v}_2$ ,  $\lambda \mathbf{v}$  with the 'natural' properties

$$(\mathbf{v}_1 + \mathbf{v}_2) + \mathbf{v}_3 = \mathbf{v}_1 + (\mathbf{v}_2 + \mathbf{v}_3), \text{ (associativity)}$$

$$\mathbf{v}_1 + \mathbf{v}_2 = \mathbf{v}_2 + \mathbf{v}_1, \quad \text{ (commutativity)}$$

$$\exists \mathbf{0} : \mathbf{v} + \mathbf{0} = \mathbf{v},$$

# Vector space: $(V, +, \lambda \cdot)$

- Points = vectors.
- 2 operations:  $\mathbf{v}_1 + \mathbf{v}_2$ ,  $\lambda\mathbf{v}$  with the 'natural' properties

$$(\mathbf{v}_1 + \mathbf{v}_2) + \mathbf{v}_3 = \mathbf{v}_1 + (\mathbf{v}_2 + \mathbf{v}_3), \text{ (associativity)}$$

$$\mathbf{v}_1 + \mathbf{v}_2 = \mathbf{v}_2 + \mathbf{v}_1, \quad \text{ (commutativity)}$$

$$\exists \mathbf{0} : \mathbf{v} + \mathbf{0} = \mathbf{v},$$

$$\exists -\mathbf{v} : \mathbf{v} + (-\mathbf{v}) = \mathbf{0},$$

# Vector space: $(V, +, \lambda \cdot)$

- Points = vectors.
- 2 operations:  $\mathbf{v}_1 + \mathbf{v}_2$ ,  $\lambda\mathbf{v}$  with the 'natural' properties

$$(\mathbf{v}_1 + \mathbf{v}_2) + \mathbf{v}_3 = \mathbf{v}_1 + (\mathbf{v}_2 + \mathbf{v}_3), \text{ (associativity)}$$

$$\mathbf{v}_1 + \mathbf{v}_2 = \mathbf{v}_2 + \mathbf{v}_1, \quad \text{ (commutativity)}$$

$$\exists \mathbf{0} : \mathbf{v} + \mathbf{0} = \mathbf{v},$$

$$\exists -\mathbf{v} : \mathbf{v} + (-\mathbf{v}) = \mathbf{0},$$

$$a(b\mathbf{v}) = (ab)\mathbf{v},$$

# Vector space: $(V, +, \lambda \cdot)$

- Points = vectors.
- 2 operations:  $\mathbf{v}_1 + \mathbf{v}_2$ ,  $\lambda\mathbf{v}$  with the 'natural' properties

$$(\mathbf{v}_1 + \mathbf{v}_2) + \mathbf{v}_3 = \mathbf{v}_1 + (\mathbf{v}_2 + \mathbf{v}_3), \text{ (associativity)}$$

$$\mathbf{v}_1 + \mathbf{v}_2 = \mathbf{v}_2 + \mathbf{v}_1, \quad \text{ (commutativity)}$$

$$\exists \mathbf{0} : \mathbf{v} + \mathbf{0} = \mathbf{v},$$

$$\exists -\mathbf{v} : \mathbf{v} + (-\mathbf{v}) = \mathbf{0},$$

$$a(b\mathbf{v}) = (ab)\mathbf{v},$$

$$1\mathbf{v} = \mathbf{v},$$

# Vector space: $(V, +, \lambda \cdot)$

- Points = vectors.
- 2 operations:  $\mathbf{v}_1 + \mathbf{v}_2$ ,  $\lambda\mathbf{v}$  with the 'natural' properties

$$(\mathbf{v}_1 + \mathbf{v}_2) + \mathbf{v}_3 = \mathbf{v}_1 + (\mathbf{v}_2 + \mathbf{v}_3), \text{ (associativity)}$$

$$\mathbf{v}_1 + \mathbf{v}_2 = \mathbf{v}_2 + \mathbf{v}_1, \quad \text{ (commutativity)}$$

$$\exists \mathbf{0} : \mathbf{v} + \mathbf{0} = \mathbf{v},$$

$$\exists -\mathbf{v} : \mathbf{v} + (-\mathbf{v}) = \mathbf{0},$$

$$a(b\mathbf{v}) = (ab)\mathbf{v},$$

$$1\mathbf{v} = \mathbf{v},$$

$$a(\mathbf{v}_1 + \mathbf{v}_2) = a\mathbf{v}_1 + a\mathbf{v}_2,$$

# Vector space: $(V, +, \lambda \cdot)$

- Points = vectors.
- 2 operations:  $\mathbf{v}_1 + \mathbf{v}_2$ ,  $\lambda \mathbf{v}$  with the 'natural' properties

$$(\mathbf{v}_1 + \mathbf{v}_2) + \mathbf{v}_3 = \mathbf{v}_1 + (\mathbf{v}_2 + \mathbf{v}_3), \text{ (associativity)}$$

$$\mathbf{v}_1 + \mathbf{v}_2 = \mathbf{v}_2 + \mathbf{v}_1, \quad \text{ (commutativity)}$$

$$\exists \mathbf{0} : \mathbf{v} + \mathbf{0} = \mathbf{v},$$

$$\exists -\mathbf{v} : \mathbf{v} + (-\mathbf{v}) = \mathbf{0},$$

$$a(b\mathbf{v}) = (ab)\mathbf{v},$$

$$1\mathbf{v} = \mathbf{v},$$

$$a(\mathbf{v}_1 + \mathbf{v}_2) = a\mathbf{v}_1 + a\mathbf{v}_2,$$

$$(a + b)\mathbf{v} = a\mathbf{v} + b\mathbf{v} \quad \text{for } \forall \mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3, \mathbf{v} \in V, a, b \in \mathbb{R}.$$

# Vector space: examples

①  $(\mathbb{R}^d, +, \cdot)$  defined as

$$(x_1, \dots, x_d) + (y_1, \dots, y_d) := (x_1 + y_1, \dots, x_d + y_d),$$
$$\lambda \cdot (x_1, \dots, x_d) := (\lambda x_1, \dots, \lambda x_d).$$

# Vector space: examples

- ①  $(\mathbb{R}^d, +, \cdot)$  defined as

$$(x_1, \dots, x_d) + (y_1, \dots, y_d) := (x_1 + y_1, \dots, x_d + y_d),$$
$$\lambda \cdot (x_1, \dots, x_d) := (\lambda x_1, \dots, \lambda x_d).$$

- ② Polynomials:

$$\sum_i a_i x^i + \sum_i b_i x^i := \sum_i (a_i + b_i) x^i,$$
$$\lambda \cdot \left( \sum_i a_i x^i \right) := \sum_i (\lambda a_i) x^i.$$

# Vector space: examples

①  $(\mathbb{R}^d, +, \cdot)$  defined as

$$(x_1, \dots, x_d) + (y_1, \dots, y_d) := (x_1 + y_1, \dots, x_d + y_d),$$
$$\lambda \cdot (x_1, \dots, x_d) := (\lambda x_1, \dots, \lambda x_d).$$

② Polynomials:

$$\sum_i a_i x^i + \sum_i b_i x^i := \sum_i (a_i + b_i) x^i,$$

$$\lambda \cdot \left( \sum_i a_i x^i \right) := \sum_i (\lambda a_i) x^i.$$

③  $\mathcal{X} \rightarrow \mathbb{R}$ -valued functions:

$$(f + g)(x) := f(x) + g(x), \quad (\lambda \cdot f)(x) := \lambda f(x).$$

# Vector space: examples

- ①  $(\mathbb{R}^d, +, \cdot)$  defined as

$$(x_1, \dots, x_d) + (y_1, \dots, y_d) := (x_1 + y_1, \dots, x_d + y_d),$$
$$\lambda \cdot (x_1, \dots, x_d) := (\lambda x_1, \dots, \lambda x_d).$$

- ② Polynomials:

$$\sum_i a_i x^i + \sum_i b_i x^i := \sum_i (a_i + b_i) x^i,$$

$$\lambda \cdot \left( \sum_i a_i x^i \right) := \sum_i (\lambda a_i) x^i.$$

- ③  $\mathcal{X} \rightarrow \mathbb{R}$ -valued functions:

$$(f + g)(x) := f(x) + g(x), \quad (\lambda \cdot f)(x) := \lambda f(x).$$

Previously:  $\mathcal{X} = \{1, \dots, d\}$ ,  $\mathcal{X} = \mathbb{N}$ .

- ④  $\supset C[a, b]$ : continuous functions.

## Vector space: examples – continued

④  $\supset C[a, b]$ : continuous functions.

⑤  $(\mathbb{R}^d, \otimes, \cdot)$ :

$$(x_1, \dots, x_d) \otimes (y_1, \dots, y_d) := [x_i y_j]_{i,j=1}^d$$

$$\lambda \cdot (x_1, \dots, x_d) := (\lambda x_1, \dots, \lambda x_d).$$

## Vector space: examples – continued

④  $\supset C[a, b]$ : continuous functions.

⑤  $(\mathbb{R}^d, \otimes, \cdot)$ : NO!!!

$$(x_1, \dots, x_d) \otimes (y_1, \dots, y_d) := [x_i y_j]_{i,j=1}^d \in \mathbb{R}^{d \times d} \neq \mathbb{R}^d,$$

$$\lambda \cdot (x_1, \dots, x_d) := (\lambda x_1, \dots, \lambda x_d).$$

## Vector space: examples – continued

④  $\supset C[a, b]$ : continuous functions.

⑤  $(\mathbb{R}^d, \otimes, \cdot)$ : NO!!!

$$(x_1, \dots, x_d) \otimes (y_1, \dots, y_d) := [x_i y_j]_{i,j=1}^d \in \mathbb{R}^{d \times d} \neq \mathbb{R}^d,$$

$$\lambda \cdot (x_1, \dots, x_d) := (\lambda x_1, \dots, \lambda x_d).$$

⑥  $(\mathbb{R}^d, \langle \cdot, \cdot \rangle, \cdot)$ : similarly NO

$$\langle (x_1, \dots, x_d), (y_1, \dots, y_d) \rangle := \sum_{i=1}^d x_i y_i \in \mathbb{R} \neq \mathbb{R}^d,$$

$$\lambda \cdot (x_1, \dots, x_d) := (\lambda x_1, \dots, \lambda x_d).$$

## Vector space: examples – continued

④  $\supset C[a, b]$ : continuous functions.

⑤  $(\mathbb{R}^d, \otimes, \cdot)$ : NO!!!

$$(x_1, \dots, x_d) \otimes (y_1, \dots, y_d) := [x_i y_j]_{i,j=1}^d \in \mathbb{R}^{d \times d} \neq \mathbb{R}^d,$$

$$\lambda \cdot (x_1, \dots, x_d) := (\lambda x_1, \dots, \lambda x_d).$$

⑥  $(\mathbb{R}^d, \langle \cdot, \cdot \rangle, \cdot)$ : similarly NO

$$\langle (x_1, \dots, x_d), (y_1, \dots, y_d) \rangle := \sum_{i=1}^d x_i y_i \in \mathbb{R} \neq \mathbb{R}^d,$$

$$\lambda \cdot (x_1, \dots, x_d) := (\lambda x_1, \dots, \lambda x_d).$$

Now we put a notion of norm, inner product to vectors .

We define the 'length' of a vector.

$\mathcal{H}$ : vector space over  $\mathbb{R}$ .  $\|\cdot\| : \mathcal{H} \rightarrow [0, \infty)$  is norm on  $\mathcal{H}$ , if  
①  $\|f\| = 0$  iff.  $f = 0$  (norm separates points),

# Normed space

We define the 'length' of a vector.

$\mathcal{H}$ : vector space over  $\mathbb{R}$ .  $\|\cdot\| : \mathcal{H} \rightarrow [0, \infty)$  is norm on  $\mathcal{H}$ , if

- ①  $\|f\| = 0$  iff.  $f = 0$  (norm separates points),
- ②  $\|\lambda f\| = |\lambda| \|f\| \quad \forall \lambda \in \mathbb{R}, \forall f \in \mathcal{H}$  (positive homogeneity),

# Normed space

We define the 'length' of a vector.

$\mathcal{H}$ : vector space over  $\mathbb{R}$ .  $\|\cdot\| : \mathcal{H} \rightarrow [0, \infty)$  is norm on  $\mathcal{H}$ , if

- ①  $\|f\| = 0$  iff.  $f = 0$  (norm separates points),
- ②  $\|\lambda f\| = |\lambda| \|f\| \quad \forall \lambda \in \mathbb{R}, \forall f \in \mathcal{H}$  (positive homogeneity),
- ③  $\|f + g\| \leq \|f\| + \|g\| \quad \forall f, g \in \mathcal{H}$  (triangle inequality).

# Normed space

We define the 'length' of a vector.

$\mathcal{H}$ : vector space over  $\mathbb{R}$ .  $\|\cdot\| : \mathcal{H} \rightarrow [0, \infty)$  is norm on  $\mathcal{H}$ , if

- ①  $\|f\| = 0$  iff.  $f = 0$  (norm separates points),
- ②  $\|\lambda f\| = |\lambda| \|f\| \quad \forall \lambda \in \mathbb{R}, \forall f \in \mathcal{H}$  (positive homogeneity),
- ③  $\|f + g\| \leq \|f\| + \|g\| \quad \forall f, g \in \mathcal{H}$  (triangle inequality).

Note:

- norm  $\Rightarrow$  metric:  $\rho(f, g) = \|f - g\| \Rightarrow$
- study continuity, convergence.

- $(\mathbb{R}, |\cdot|)$ .

- $(\mathbb{R}, |\cdot|)$ .
- $\left(\mathbb{R}^d, \|\mathbf{x}\|_p = \left[\sum_i |x_i|^p\right]^{\frac{1}{p}}\right), p \in [0, \infty]$ .
  - $p = 2$ :  $\|\mathbf{x}\|_2 = \sqrt{\sum_i x_i^2}$  (**Euclidean**),

- $(\mathbb{R}, |\cdot|)$ .
- $(\mathbb{R}^d, \|\mathbf{x}\|_p = [\sum_i |x_i|^p]^{\frac{1}{p}})$ ,  $p \in [0, \infty]$ .
  - $p = 2$ :  $\|\mathbf{x}\|_2 = \sqrt{\sum_i x_i^2}$  (**Euclidean**),
  - $p = 1$ :  $\|\mathbf{x}\|_1 = \sum_i |x_i|$  (**Manhattan**),

- $(\mathbb{R}, |\cdot|)$ .
- $(\mathbb{R}^d, \|\mathbf{x}\|_p = [\sum_i |x_i|^p]^{\frac{1}{p}})$ ,  $p \in [0, \infty]$ .
  - $p = 2$ :  $\|\mathbf{x}\|_2 = \sqrt{\sum_i x_i^2}$  (**Euclidean**),
  - $p = 1$ :  $\|\mathbf{x}\|_1 = \sum_i |x_i|$  (**Manhattan**),
  - $p = \infty$ :  $\|\mathbf{x}\|_\infty = \max_i |x_i|$  (**maximum norm**).

- $(\mathbb{R}, |\cdot|)$ .
- $\left(\mathbb{R}^d, \|\mathbf{x}\|_p = \left[\sum_i |x_i|^p\right]^{\frac{1}{p}}\right), p \in [0, \infty]$ .
  - $p = 2$ :  $\|\mathbf{x}\|_2 = \sqrt{\sum_i x_i^2}$  (**Euclidean**),
  - $p = 1$ :  $\|\mathbf{x}\|_1 = \sum_i |x_i|$  (**Manhattan**),
  - $p = \infty$ :  $\|\mathbf{x}\|_\infty = \max_i |x_i|$  (**maximum norm**).
- $\left(C[a, b], \|f\|_p = \left[\int_a^b |f(x)|^p dx\right]^{\frac{1}{p}}\right), 1 \leq p < \infty$ .

- $(\mathbb{R}, |\cdot|)$ .
- $\left(\mathbb{R}^d, \|\mathbf{x}\|_p = \left[\sum_i |x_i|^p\right]^{\frac{1}{p}}\right), p \in [0, \infty]$ .
  - $p = 2$ :  $\|\mathbf{x}\|_2 = \sqrt{\sum_i x_i^2}$  (Euclidean),
  - $p = 1$ :  $\|\mathbf{x}\|_1 = \sum_i |x_i|$  (Manhattan),
  - $p = \infty$ :  $\|\mathbf{x}\|_\infty = \max_i |x_i|$  (maximum norm).
- $\left(C[a, b], \|f\|_p = \left[\int_a^b |f(x)|^p dx\right]^{\frac{1}{p}}\right), 1 \leq p < \infty$ .
- $(C[a, b], \|f\|_\infty = \max_{x \in [a, b]} |f(x)|), p = \infty$ .

# Inner product space (also called Euclidean space)

$\mathcal{H}$ : vector space over  $\mathbb{R}$ .  $\langle \cdot, \cdot \rangle : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$  is an **inner product** on  $\mathcal{H}$  if for  $\forall \alpha_i \in \mathbb{R}, f_i, f, g \in \mathcal{H}$

- ①  $\langle \alpha_1 f_1 + \alpha_2 f_2, g \rangle = \alpha_1 \langle f_1, g \rangle + \alpha_2 \langle f_2, g \rangle$  (**linearity**),

# Inner product space (also called Euclidean space)

$\mathcal{H}$ : vector space over  $\mathbb{R}$ .  $\langle \cdot, \cdot \rangle : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$  is an **inner product** on  $\mathcal{H}$  if for  $\forall \alpha_i \in \mathbb{R}, f_i, f, g \in \mathcal{H}$

- ①  $\langle \alpha_1 f_1 + \alpha_2 f_2, g \rangle = \alpha_1 \langle f_1, g \rangle + \alpha_2 \langle f_2, g \rangle$  (**linearity**),
- ②  $\langle f, g \rangle = \langle g, f \rangle$  (**symmetry**),

# Inner product space (also called Euclidean space)

$\mathcal{H}$ : vector space over  $\mathbb{R}$ .  $\langle \cdot, \cdot \rangle : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$  is an **inner product** on  $\mathcal{H}$  if for  $\forall \alpha_i \in \mathbb{R}, f_i, f, g \in \mathcal{H}$

- ①  $\langle \alpha_1 f_1 + \alpha_2 f_2, g \rangle = \alpha_1 \langle f_1, g \rangle + \alpha_2 \langle f_2, g \rangle$  (**linearity**),
- ②  $\langle f, g \rangle = \langle g, f \rangle$  (**symmetry**),
- ③  $\langle f, f \rangle \geq 0; \langle f, f \rangle = 0 \Leftrightarrow f = 0.$

# Inner product space (also called Euclidean space)

$\mathcal{H}$ : vector space over  $\mathbb{R}$ .  $\langle \cdot, \cdot \rangle : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$  is an **inner product** on  $\mathcal{H}$  if for  $\forall \alpha_i \in \mathbb{R}, f_i, f, g \in \mathcal{H}$

- ①  $\langle \alpha_1 f_1 + \alpha_2 f_2, g \rangle = \alpha_1 \langle f_1, g \rangle + \alpha_2 \langle f_2, g \rangle$  (**linearity**),
- ②  $\langle f, g \rangle = \langle g, f \rangle$  (**symmetry**),
- ③  $\langle f, f \rangle \geq 0; \langle f, f \rangle = 0 \Leftrightarrow f = 0.$

Notes:

- 1, 2  $\Rightarrow$  bilinearity.

# Inner product space (also called Euclidean space)

$\mathcal{H}$ : vector space over  $\mathbb{R}$ .  $\langle \cdot, \cdot \rangle : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$  is an **inner product** on  $\mathcal{H}$  if for  $\forall \alpha_i \in \mathbb{R}, f_i, f, g \in \mathcal{H}$

- ①  $\langle \alpha_1 f_1 + \alpha_2 f_2, g \rangle = \alpha_1 \langle f_1, g \rangle + \alpha_2 \langle f_2, g \rangle$  (**linearity**),
- ②  $\langle f, g \rangle = \langle g, f \rangle$  (**symmetry**),
- ③  $\langle f, f \rangle \geq 0; \langle f, f \rangle = 0 \Leftrightarrow f = 0.$

Notes:

- 1, 2  $\Rightarrow$  bilinearity. Inner product  $\Rightarrow$

$$\text{norm: } \|f\| = \sqrt{\langle f, f \rangle}, \quad \text{angle: } \cos(f, g) = \frac{\langle f, g \rangle}{\|f\| \|g\|}.$$

- $\left(\mathbb{R}^d, \langle \mathbf{x}, \mathbf{y} \rangle = \sum_{i=1}^d x_i y_i\right).$
- $\left(\mathbb{R}^{d_1 \times d_2}, \langle \mathbf{A}, \mathbf{B} \rangle_F = \text{tr}(\mathbf{A}^T \mathbf{B}) = \sum_{i=1}^{d_1} \sum_{j=1}^{d_2} A_{ij} B_{ij}\right).$
- $\left(C[a, b], \langle f, g \rangle = \int_a^b f(x)g(x)dx\right).$

# Norm vs inner product

Relations:

- $|\langle f, g \rangle| \leq \|f\| \cdot \|g\|$  (CBS),

# Norm vs inner product

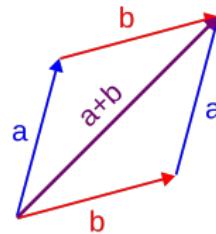
Relations:

- $|\langle f, g \rangle| \leq \|f\| \cdot \|g\|$  (CBS),
- $4\langle f, g \rangle = \|f + g\|^2 - \|f - g\|^2$  (polarization identity),

# Norm vs inner product

Relations:

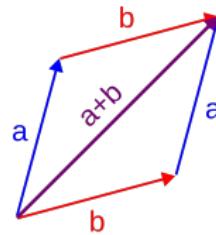
- $|\langle f, g \rangle| \leq \|f\| \cdot \|g\|$  (CBS),
- $4\langle f, g \rangle = \|f + g\|^2 - \|f - g\|^2$  (polarization identity),
- $\|a + b\|^2 + \|a - b\|^2 = (*)$  (parallelogram rule).



# Norm vs inner product

Relations:

- $|\langle f, g \rangle| \leq \|f\| \cdot \|g\|$  (CBS),
- $4\langle f, g \rangle = \|f + g\|^2 - \|f - g\|^2$  (polarization identity),
- $\|a + b\|^2 + \|a - b\|^2 = (*)$  (parallelogram rule).

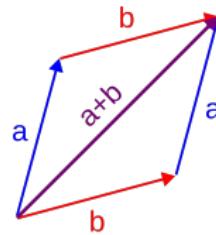


$$(*) = \langle a + b, a + b \rangle + \langle a - b, a - b \rangle$$

# Norm vs inner product

Relations:

- $|\langle f, g \rangle| \leq \|f\| \cdot \|g\|$  (CBS),
- $4\langle f, g \rangle = \|f + g\|^2 - \|f - g\|^2$  (polarization identity),
- $\|a + b\|^2 + \|a - b\|^2 = (*)$  (parallelogram rule).

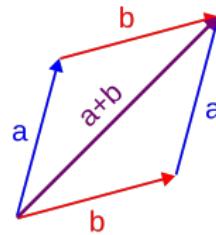


$$\begin{aligned} (*) &= \langle a + b, a + b \rangle + \langle a - b, a - b \rangle \\ &= 2(\|a\|^2 + \|b\|^2) \pm 2\langle a, b \rangle \end{aligned}$$

# Norm vs inner product

Relations:

- $|\langle f, g \rangle| \leq \|f\| \cdot \|g\|$  (CBS),
- $4\langle f, g \rangle = \|f + g\|^2 - \|f - g\|^2$  (polarization identity),
- $\|a + b\|^2 + \|a - b\|^2 = (*)$  (parallelogram rule).



$$\begin{aligned} (*) &= \langle a + b, a + b \rangle + \langle a - b, a - b \rangle \\ &= 2(\|a\|^2 + \|b\|^2) \pm 2\langle a, b \rangle = 2\|a\|^2 + 2\|b\|^2. \end{aligned}$$

The parallelogram rule holds in an inner product space.

## Example when the parallelogram rule fails

$C[0, 1]$  with  $\|f\|_\infty = \max_{x \in [0,1]} |f(x)|$ :

$$f(x) := 1 - x, \quad g(x) := x$$

## Example when the parallelogram rule fails

$C[0, 1]$  with  $\|f\|_\infty = \max_{x \in [0, 1]} |f(x)|$ :

$$f(x) := 1 - x, \quad g(x) := x,$$
$$\|f + g\|_\infty^2 + \|f - g\|_\infty^2 = \|1\|_\infty^2 + \underbrace{\|\underbrace{1 - 2x}_{\in [-1,1]}\|_\infty^2}_{\in [-1,1]} = 1^2 + 1^2 = 2,$$

## Example when the parallelogram rule fails

$C[0, 1]$  with  $\|f\|_\infty = \max_{x \in [0, 1]} |f(x)|$ :

$$f(x) := 1 - x, \quad g(x) := x,$$

$$\|f + g\|_\infty^2 + \|f - g\|_\infty^2 = \|1\|_\infty^2 + \underbrace{\|1 - 2x\|_\infty^2}_{\in [-1, 1]} = 1^2 + 1^2 = 2,$$

$$2 \left( \|f\|_\infty^2 + \|g\|_\infty^2 \right) = 2 \underbrace{\|1 - x\|_\infty^2}_{\in [0, 1]} + 2 \underbrace{\|x\|_\infty^2}_{\in [0, 1]} = 2 + 2 = 4.$$

## Example when the parallelogram rule fails

$C[0, 1]$  with  $\|f\|_\infty = \max_{x \in [0,1]} |f(x)|$ :

$$f(x) := 1 - x, \quad g(x) := x,$$

$$\|f + g\|_\infty^2 + \|f - g\|_\infty^2 = \|1\|_\infty^2 + \underbrace{\|1 - 2x\|_\infty^2}_{\in [-1,1]} = 1^2 + 1^2 = 2,$$

$$2 \left( \|f\|_\infty^2 + \|g\|_\infty^2 \right) = 2 \underbrace{\|1 - x\|_\infty^2}_{\in [0,1]} + 2 \underbrace{\|x\|_\infty^2}_{\in [0,1]} = 2 + 2 = 4.$$

## Characterization

A norm is induced by an inner product iff

$$\|f + g\|^2 + \|f - g\|^2 = 2 \left( \|f\|^2 + \|g\|^2 \right) \quad \forall f, g.$$

# Completeness: motivation

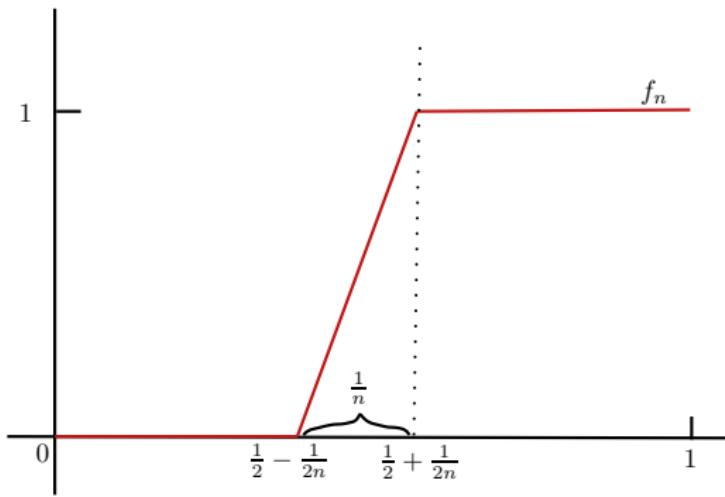
- $\{(f_n)\}_{n \in \mathbb{N}} \subset \mathbb{Q}$ : 1, 1.4, 1.41, 1.414, 1.4142, ...
  - $|f_n - f_m|$  can be arbitrary small for large enough  $(n, m) \leftrightarrow$  Cauchy seq ,

# Completeness: motivation

- $\{(f_n)\}_{n \in \mathbb{N}} \subset \mathbb{Q}$ : 1, 1.4, 1.41, 1.414, 1.4142, ...
  - $|f_n - f_m|$  can be arbitrary small for large enough  $(n, m) \leftrightarrow$  Cauchy seq ,
  - but  $(f_n)$  does not converge in  $\mathbb{Q} \notin \sqrt{2}$ .

# Completeness: motivation

- $\{(f_n)\}_{n \in \mathbb{N}} \subset \mathbb{Q}$ :  $1, 1.4, 1.41, 1.414, 1.4142, \dots$ 
  - $|f_n - f_m|$  can be arbitrary small for large enough  $(n, m) \leftrightarrow$  Cauchy seq ,
  - but  $(f_n)$  does not converge in  $\mathbb{Q} \notin \sqrt{2}$ .
- $(C[0, 1], \|\cdot\|_{L^2[0,1]})$ : similarly



# Hilbert space

- $\mathcal{H}$ ilbert space := complete Euclidean space. Prototype:

$$L^2(\mathcal{X}, \mu) := \left\{ f : \mathcal{X} \rightarrow \mathbb{R} : \|f\|_2 = \left[ \int_{\mathcal{X}} |f(x)|^2 d\mu(x) \right]^{1/2} < \infty \right\}.$$

# Hilbert space

- $\mathcal{H}$ ilbert space := complete Euclidean space. Prototype:

$$L^2(\mathcal{X}, \mu) := \left\{ f : \mathcal{X} \rightarrow \mathbb{R} : \|f\|_2 = \left[ \int_{\mathcal{X}} |f(x)|^2 d\mu(x) \right]^{1/2} < \infty \right\}.$$

Specifically:

$$(\mathbb{R}^d, \|\cdot\|_2), \text{ or } \ell^2(\mathbb{N}) = \left\{ (a_n)_{n \in \mathbb{N}} : \sqrt{\sum_{n=1}^{\infty} a_n^2} < \infty \right\}.$$

Recall:  $k(x, y) = \langle \varphi(x), \varphi(y) \rangle_{\mathcal{H}}$ .

# Hilbert space

- Hilbert space := complete Euclidean space. Prototype:

$$L^2(\mathcal{X}, \mu) := \left\{ f : \mathcal{X} \rightarrow \mathbb{R} : \|f\|_2 = \left[ \int_{\mathcal{X}} |f(x)|^2 d\mu(x) \right]^{1/2} < \infty \right\}.$$

Specifically:

$$(\mathbb{R}^d, \|\cdot\|_2), \text{ or } \ell^2(\mathbb{N}) = \left\{ (a_n)_{n \in \mathbb{N}} : \sqrt{\sum_{n=1}^{\infty} a_n^2} < \infty \right\}.$$

Recall:  $k(x, y) = \langle \varphi(x), \varphi(y) \rangle_{\mathcal{H}}$ .

- Banach space := complete normed space:

$$L^p(\mathcal{X}, \mu), \quad (\mathbb{R}^d, \|\cdot\|_p), \quad \ell^p(\mathbb{N}), \quad (C[a, b], \|\cdot\|_{\infty}).$$

## Kernels, RKHS: Definition-2

- Def-1 = feature space point of view,  $k(a, b) = \langle \varphi(a), \varphi(b) \rangle_{\mathcal{H}}$ .

# Kernels, RKHS: Definition-2

- Def-1 = feature space point of view,  $k(a, b) = \langle \varphi(a), \varphi(b) \rangle_{\mathcal{H}}$ .
- Def-2 = constructive.  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is a **reproducing kernel** of a  $\mathcal{H}$ (ilbert)  $\subset \mathbb{R}^{\mathcal{X}}$

$$\underbrace{k(\cdot, b)}_{\text{a blue bell curve}} \in \mathcal{H}, \quad \underbrace{f(b) = \langle f, k(\cdot, b) \rangle_{\mathcal{H}}}_{\text{reproducing property}}.$$

# Kernels, RKHS: Definition-2

- Def-1 = feature space point of view,  $k(a, b) = \langle \varphi(a), \varphi(b) \rangle_{\mathcal{H}}$ .
- Def-2 = constructive.  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is a **reproducing kernel** of a Hilbert space  $\mathcal{H} \subset \mathbb{R}^{\mathcal{X}}$

$$\underbrace{k(\cdot, b)}_{\text{represents a function}} \in \mathcal{H}, \quad f(b) = \underbrace{\langle f, k(\cdot, b) \rangle_{\mathcal{H}}}_{\text{reproducing property}}.$$


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

# Kernels, RKHS: Definition-2

- Def-1 = feature space point of view,  $k(a, b) = \langle \varphi(a), \varphi(b) \rangle_{\mathcal{H}}$ .
- Def-2 = constructive.  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is a **reproducing kernel** of a Hilbert space  $\mathcal{H} \subset \mathbb{R}^{\mathcal{X}}$

$$\underbrace{k(\cdot, b)}_{\text{a blue bell-shaped curve}} \in \mathcal{H}, \quad \underbrace{f(b) = \langle f, k(\cdot, b) \rangle_{\mathcal{H}}}_{\text{reproducing property}}.$$

$$\xrightarrow{\text{spec.}} k(a, b) = \langle k(\cdot, a), k(\cdot, b) \rangle_{\mathcal{H}}. \quad \mathcal{H}_k = \overline{\{\sum_{i=1}^n \alpha_i k(\cdot, x_i)\}}.$$

## Kernels: Definition-3

- Def-3: Gram matrix, optimization point of view.
- Intuition:  $\mathcal{X} := \mathbb{R}^d$ , data matrix  $\mathbf{X} = [x_1, \dots, x_n] \in \mathbb{R}^{d \times n}$ , then

$$\mathbf{G} := \mathbf{X}^T \mathbf{X} = [\langle x_i, x_j \rangle_2]_{i,j=1}^n \succeq 0.$$

## Kernels: Definition-3

- Def-3: Gram matrix, optimization point of view.
- Intuition:  $\mathcal{X} := \mathbb{R}^d$ , data matrix  $\mathbf{X} = [x_1, \dots, x_n] \in \mathbb{R}^{d \times n}$ , then

$$\mathbf{G} := \mathbf{X}^T \mathbf{X} = [\langle x_i, x_j \rangle_2]_{i,j=1}^n \geq 0.$$

i.e.

$$\mathbf{G}^T = \mathbf{G} \quad (\text{symmetry}),$$

## Kernels: Definition-3

- Def-3: Gram matrix, optimization point of view.
- Intuition:  $\mathcal{X} := \mathbb{R}^d$ , data matrix  $\mathbf{X} = [x_1, \dots, x_n] \in \mathbb{R}^{d \times n}$ , then

$$\mathbf{G} := \mathbf{X}^T \mathbf{X} = [\langle x_i, x_j \rangle_2]_{i,j=1}^n \geq 0.$$

i.e.

$$\begin{aligned}\mathbf{G}^T &= \mathbf{G} && \text{(symmetry),} \\ \mathbf{v}^T \mathbf{G} \mathbf{v} &= \mathbf{v}^T \mathbf{X}^T \mathbf{X} \mathbf{v} = \|\mathbf{X} \mathbf{v}\|_2^2 \geq 0 && (\forall \mathbf{v} \in \mathbb{R}^d).\end{aligned}$$

## Kernels: Definition-3

- Def-3: Gram matrix, optimization point of view.
- Intuition:  $\mathcal{X} := \mathbb{R}^d$ , data matrix  $\mathbf{X} = [x_1, \dots, x_n] \in \mathbb{R}^{d \times n}$ , then

$$\mathbf{G} := \mathbf{X}^T \mathbf{X} = [\langle x_i, x_j \rangle_2]_{i,j=1}^n \geq 0.$$

i.e.

$$\begin{aligned}\mathbf{G}^T &= \mathbf{G} && \text{(symmetry),} \\ \mathbf{v}^T \mathbf{G} \mathbf{v} &= \mathbf{v}^T \mathbf{X}^T \mathbf{X} \mathbf{v} = \|\mathbf{X} \mathbf{v}\|_2^2 \geq 0 && (\forall \mathbf{v} \in \mathbb{R}^d).\end{aligned}$$

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

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

- Def-4 intuition: We want

$$(f_n)_{n \in \mathbb{N}} \xrightarrow{\|\cdot\|} f \quad \Rightarrow \quad (f_n(x))_{n \in \mathbb{N}} \rightarrow f(x) \quad \forall x.$$

# Kernels: Definition-4 – motivation

- Def-4 intuition: We want

$$(f_n)_{n \in \mathbb{N}} \xrightarrow{\|\cdot\|} f \quad \Rightarrow \quad (f_n(x))_{n \in \mathbb{N}} \rightarrow f(x) \quad \forall x.$$

- Example-1: For  $\mathcal{H} := (C[0, 1], \|\cdot\|_\infty)$

$$|f_n(x) - f(x)| \leq \max_{x \in [0, 1]} |f_n(x) - f(x)| = \|f_n - f\|_\infty \xrightarrow{n \rightarrow \infty} 0,$$

# Kernels: Definition-4 – motivation

- Def-4 intuition: We want

$$(f_n)_{n \in \mathbb{N}} \xrightarrow{\|\cdot\|} f \Rightarrow (f_n(x))_{n \in \mathbb{N}} \rightarrow f(x) \quad \forall x.$$

- Example-1: For  $\mathcal{H} := (C[0, 1], \|\cdot\|_\infty)$

$$|f_n(x) - f(x)| \leq \max_{x \in [0, 1]} |f_n(x) - f(x)| = \|f_n - f\|_\infty \xrightarrow{n \rightarrow \infty} 0,$$

but no inner product in  $C[0, 1]$  (as we saw it – parallelograms).

## Kernels: Definition-4 – continued

Let us now try a Hilbert space:  $\mathcal{H} = L^2[0, 1] \ni f_n(x) = x^n$  (simple).

## Kernels: Definition-4 – continued

Let us now try a Hilbert space:  $\mathcal{H} = L^2[0, 1] \ni f_n(x) = x^n$  (simple).

①  $f_n \xrightarrow{n \rightarrow \infty} 0 := f^* \in \mathcal{H}$  since

$$\lim_{n \rightarrow \infty} \|f_n - 0\|_{L^2[0,1]} = \lim_{n \rightarrow \infty} \left( \int_0^1 x^{2n} dx \right)^{1/2}$$

## Kernels: Definition-4 – continued

Let us now try a Hilbert space:  $\mathcal{H} = L^2[0, 1] \ni f_n(x) = x^n$  (simple).

①  $f_n \xrightarrow{n \rightarrow \infty} 0 := f^* \in \mathcal{H}$  since

$$\lim_{n \rightarrow \infty} \|f_n - 0\|_{L^2[0,1]} = \lim_{n \rightarrow \infty} \left( \int_0^1 x^{2n} dx \right)^{1/2} = \lim_{n \rightarrow \infty} \sqrt{\left[ \frac{x^{2n+1}}{2n+1} \right]_{x=0}^{x=1}}$$

## Kernels: Definition-4 – continued

Let us now try a Hilbert space:  $\mathcal{H} = L^2[0, 1] \ni f_n(x) = x^n$  (simple).

①  $f_n \xrightarrow{n \rightarrow \infty} 0 := f^* \in \mathcal{H}$  since

$$\begin{aligned}\lim_{n \rightarrow \infty} \|f_n - 0\|_{L^2[0,1]} &= \lim_{n \rightarrow \infty} \left( \int_0^1 x^{2n} dx \right)^{1/2} = \lim_{n \rightarrow \infty} \sqrt{\left[ \frac{x^{2n+1}}{2n+1} \right]_{x=0}^{x=1}} \\ &= \lim_{n \rightarrow \infty} \frac{1}{\sqrt{2n+1}} = 0,\end{aligned}$$

## Kernels: Definition-4 – continued

Let us now try a Hilbert space:  $\mathcal{H} = L^2[0, 1] \ni f_n(x) = x^n$  (simple).

- ①  $f_n \xrightarrow{n \rightarrow \infty} 0 := f^* \in \mathcal{H}$  since

$$\begin{aligned}\lim_{n \rightarrow \infty} \|f_n - 0\|_{L^2[0,1]} &= \lim_{n \rightarrow \infty} \left( \int_0^1 x^{2n} dx \right)^{1/2} = \lim_{n \rightarrow \infty} \sqrt{\left[ \frac{x^{2n+1}}{2n+1} \right]_{x=0}^{x=1}} \\ &= \lim_{n \rightarrow \infty} \frac{1}{\sqrt{2n+1}} = 0,\end{aligned}$$

- ② but  $f_n(1) = 1 \Rightarrow f^*(1) = 0$ .

In  $L^2$ : norm convergence  $\Rightarrow$  pointwise convergence.

## Kernels: Definition-4

- Evaluation functional:  $\delta_x(f) := f(x)$  is linear

$$\delta_x(f + g)$$

## Kernels: Definition-4

- Evaluation functional:  $\delta_x(f) := f(x)$  is linear

$$\delta_x(f + g) = (f + g)(x)$$

## Kernels: Definition-4

- Evaluation functional:  $\delta_x(f) := f(x)$  is linear

$$\delta_x(f + g) = (f + g)(x) = f(x) + g(x)$$

## Kernels: Definition-4

- Evaluation functional:  $\delta_x(f) := f(x)$  is linear

$$\delta_x(f + g) = (f + g)(x) = f(x) + g(x) = \delta_x(f) + \delta_x(g),$$

## Kernels: Definition-4

- Evaluation functional:  $\delta_x(f) := f(x)$  is linear

$$\delta_x(f + g) = (f + g)(x) = f(x) + g(x) = \delta_x(f) + \delta_x(g),$$
$$\delta_x(\lambda f)$$

## Kernels: Definition-4

- Evaluation functional:  $\delta_x(f) := f(x)$  is linear

$$\begin{aligned}\delta_x(f + g) &= (f + g)(x) = f(x) + g(x) = \delta_x(f) + \delta_x(g), \\ \delta_x(\lambda f) &= (\lambda f)(x)\end{aligned}$$

## Kernels: Definition-4

- Evaluation functional:  $\delta_x(f) := f(x)$  is linear

$$\delta_x(f + g) = (f + g)(x) = f(x) + g(x) = \delta_x(f) + \delta_x(g),$$

$$\delta_x(\lambda f) = (\lambda f)(x) = \lambda f(x)$$

## Kernels: Definition-4

- Evaluation functional:  $\delta_x(f) := f(x)$  is linear

$$\begin{aligned}\delta_x(f + g) &= (f + g)(x) = f(x) + g(x) = \delta_x(f) + \delta_x(g), \\ \delta_x(\lambda f) &= (\lambda f)(x) = \lambda f(x) = \lambda \delta_x(f).\end{aligned}$$

## Kernels: Definition-4

- Evaluation functional:  $\delta_x(f) := f(x)$  is linear

$$\begin{aligned}\delta_x(f + g) &= (f + g)(x) = f(x) + g(x) = \delta_x(f) + \delta_x(g), \\ \delta_x(\lambda f) &= (\lambda f)(x) = \lambda f(x) = \lambda \delta_x(f).\end{aligned}$$

- Def-4 (evaluation point of view):  $\mathcal{H} \subset \mathbb{R}^{\mathcal{X}}$  Hilbert space,

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

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

## Relation of Definition 1-4

- Def-1 (feature space):

$$k(a, b) = \langle \varphi(a), \varphi(b) \rangle_{\mathcal{H}}.$$

- Def-2 (reproducing kernel, constructive):

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

- Def-3 (Gram matrix):  $\mathbf{G} = [k(x_i, x_j)] \succeq 0$ .
- Def-4 (evaluation):  $\delta_x(f) = f(x)$  is continuous for all  $x$ .

## Relation of Definition 1-4

- Def-1 (feature space):

$$k(a, b) = \langle \varphi(a), \varphi(b) \rangle_{\mathcal{H}}.$$

- Def-2 (reproducing kernel, constructive):

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

- Def-3 (Gram matrix):  $\mathbf{G} = [k(x_i, x_j)] \succeq 0$ .
- Def-4 (evaluation):  $\delta_x(f) = f(x)$  is continuous for all  $x$ .

- All these definitions are equivalent,  $k \overset{1:1}{\leftrightarrow} \mathcal{H}_k$ .

- Trickiest direction (Moore-Aronszajn theorem):

$k$  positive definite function  $\xrightarrow{\text{construction}}$  RKHS.

## Example: every kernel is positive definite

$$\mathbf{a}^T \mathbf{G} \mathbf{a} = \sum_{i=1}^n \sum_{j=1}^n a_i a_j k(x_i, x_j)$$

## Example: every kernel is positive definite

$$\mathbf{a}^T \mathbf{G} \mathbf{a} = \sum_{i=1}^n \sum_{j=1}^n a_i a_j k(x_i, x_j) \stackrel{(i)}{=} \left\langle \sum_{i=1}^n a_i \varphi(x_i), \sum_{j=1}^n a_j \varphi(x_j) \right\rangle_{\mathcal{H}}$$

(i):  $k$  definition,  $\langle \cdot, \cdot \rangle_{\mathcal{H}}$  linear

## Example: every kernel is positive definite

$$\begin{aligned}\mathbf{a}^T \mathbf{G} \mathbf{a} &= \sum_{i=1}^n \sum_{j=1}^n a_i a_j k(x_i, x_j) \stackrel{(i)}{=} \left\langle \sum_{i=1}^n a_i \varphi(x_i), \sum_{j=1}^n a_j \varphi(x_j) \right\rangle_{\mathcal{H}} \\ &\stackrel{(ii)}{=} \left\| \sum_{i=1}^n a_i \varphi(x_i) \right\|_{\mathcal{H}}^2 \geq 0.\end{aligned}$$

(i):  $k$  definition,  $\langle \cdot, \cdot \rangle_{\mathcal{H}}$  linear, (ii)  $\|\cdot\|_{\mathcal{H}} = \sqrt{\langle \cdot, \cdot \rangle_{\mathcal{H}}}$ .

# Kernels: further examples

- $\mathcal{X} = \mathbb{R}^d, \gamma > 0:$  
$$k(x, y) = \langle \varphi(x), \varphi(y) \rangle_{\mathcal{H}}$$
$$k_p(\mathbf{x}, \mathbf{y}) = (\langle \mathbf{x}, \mathbf{y} \rangle + \gamma)^p, \quad k_G(\mathbf{x}, \mathbf{y}) = e^{-\gamma \|\mathbf{x} - \mathbf{y}\|_2^2},$$
$$k_e(\mathbf{x}, \mathbf{y}) = e^{-\gamma \|\mathbf{x} - \mathbf{y}\|_2}, \quad k_C(\mathbf{x}, \mathbf{y}) = 1 + \frac{1}{\gamma \|\mathbf{x} - \mathbf{y}\|_2^2}.$$

# Kernels: further examples

- $\mathcal{X} = \mathbb{R}^d, \gamma > 0:$  
$$k(x, y) = \langle \varphi(x), \varphi(y) \rangle_{\mathcal{H}}$$
$$k_p(\mathbf{x}, \mathbf{y}) = (\langle \mathbf{x}, \mathbf{y} \rangle + \gamma)^p, \quad k_G(\mathbf{x}, \mathbf{y}) = e^{-\gamma \|\mathbf{x} - \mathbf{y}\|_2^2},$$
$$k_e(\mathbf{x}, \mathbf{y}) = e^{-\gamma \|\mathbf{x} - \mathbf{y}\|_2}, \quad k_C(\mathbf{x}, \mathbf{y}) = 1 + \frac{1}{\gamma \|\mathbf{x} - \mathbf{y}\|_2^2}.$$
- $\mathcal{X} = \text{strings:}$ 
  - $r$ -spectrum kernel: # of common  $\leq r$ -substrings.

# Kernels: further examples

- $\mathcal{X} = \mathbb{R}^d, \gamma > 0:$  
$$k(x, y) = \langle \varphi(x), \varphi(y) \rangle_{\mathcal{H}}$$
$$k_p(\mathbf{x}, \mathbf{y}) = (\langle \mathbf{x}, \mathbf{y} \rangle + \gamma)^p, \quad k_G(\mathbf{x}, \mathbf{y}) = e^{-\gamma \|\mathbf{x} - \mathbf{y}\|_2^2},$$
$$k_e(\mathbf{x}, \mathbf{y}) = e^{-\gamma \|\mathbf{x} - \mathbf{y}\|_2}, \quad k_C(\mathbf{x}, \mathbf{y}) = 1 + \frac{1}{\gamma \|\mathbf{x} - \mathbf{y}\|_2^2}.$$
- $\mathcal{X}$  = strings:
  - $r$ -spectrum kernel: # of common  $\leq r$ -substrings.
- $\mathcal{X}$  = time-series: dynamic time-warping.

# Kernel examples – continued

Matérn kernel: flexible family, well-suited for approximation (RFF)

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

where

- $K_\nu$ : modified Bessel function of the second kind of order  $\nu$ ,
- $\Gamma(t) = \int_0^\infty x^{t-1} e^{-x} dx$ : Gamma function ( $t > 0$ ).

## Kernel examples – continued

Matérn kernel: flexible family, well-suited for approximation (RFF)

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

where

- $K_\nu$ : modified Bessel function of the second kind of order  $\nu$ ,
- $\Gamma(t) = \int_0^\infty x^{t-1} e^{-x} dx$ : Gamma function ( $t > 0$ ).

# Kernel examples – continued

Matérn kernel: flexible family, well-suited for approximation (RFF)

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

where

- $K_v$ : modified Bessel function of the second kind of order  $v$ ,
- $\Gamma(t) = \int_0^\infty x^{t-1} e^{-x} dx$ : Gamma function ( $t > 0$ ).

Specific cases:

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

# Kernel puzzle

Let

$$\mathcal{X} = \{0, 1\},$$

$$k(x, x') = \begin{cases} 1, & \text{if } x \neq x' \\ -1, & \text{if } x = x' \end{cases}.$$

# Kernel puzzle

Let

$$\mathcal{X} = \{0, 1\},$$

$$k(x, x') = \begin{cases} 1, & \text{if } x \neq x' \\ -1, & \text{if } x = x' \end{cases}.$$

## Puzzle

Is  $k$  a kernel?

# Kernel puzzle

Let

$$\mathcal{X} = \{0, 1\},$$

$$k(x, x') = \begin{cases} 1, & \text{if } x \neq x' \\ -1, & \text{if } x = x' \end{cases}.$$

Puzzle

Is  $k$  a kernel?

No!

$$k(x, x) = \langle \varphi(x), \varphi(x) \rangle_{\mathcal{H}},$$

# Kernel puzzle

Let

$$\mathcal{X} = \{0, 1\},$$

$$k(x, x') = \begin{cases} 1, & \text{if } x \neq x' \\ -1, & \text{if } x = x' \end{cases}.$$

Puzzle

Is  $k$  a kernel?

No!

$$k(x, x) = \langle \varphi(x), \varphi(x) \rangle_{\mathcal{H}} = \|\varphi(x)\|_{\mathcal{H}}^2 \geq 0 \quad (\text{Gram with } n = 1),$$

# Kernel puzzle

Let

$$\mathcal{X} = \{0, 1\},$$

$$k(x, x') = \begin{cases} 1, & \text{if } x \neq x' \\ -1, & \text{if } x = x' \end{cases}.$$

## Puzzle

Is  $k$  a kernel?

No!

$$k(x, x) = \langle \varphi(x), \varphi(x) \rangle_{\mathcal{H}} = \|\varphi(x)\|_{\mathcal{H}}^2 \geq 0 \quad (\text{Gram with } n = 1),$$
$$k(0, 0) = k(1, 1) = -1 \quad (\text{in our case}).$$

# Kernel puzzle

Let

$$\mathcal{X} = \{0, 1\},$$

$$k(x, x') = \begin{cases} 1, & \text{if } x \neq x' \\ -1, & \text{if } x = x' \end{cases}.$$

## Puzzle

Is  $k$  a kernel?

No!

$$k(x, x) = \langle \varphi(x), \varphi(x) \rangle_{\mathcal{H}} = \|\varphi(x)\|_{\mathcal{H}}^2 \geq 0 \quad (\text{Gram with } n = 1),$$

$$k(0, 0) = k(1, 1) = -1 \quad (\text{in our case}).$$

Easy-to-check conditions for a  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  function to be kernel?

- We know:  $k(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y}$  is a kernel.

- We know:  $k(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y}$  is a kernel.
- A few **useful rules**:
  - ① **Non-negative shift**.  $k$ : kernel  $\Rightarrow k + \gamma$ : kernel ( $\gamma \in \mathbb{R}^{\geq 0}$ ). Why?

- We know:  $k(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y}$  is a kernel.
- A few **useful rules**:
  - ① **Non-negative shift**.  $k$ : kernel  $\Rightarrow k + \gamma$ : kernel ( $\gamma \in \mathbb{R}^{\geq 0}$ ). Why?  $\Leftarrow$  Gram.

- We know:  $k(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y}$  is a kernel.
- A few useful rules:

- ① Non-negative shift.  $k$ : kernel  $\Rightarrow k + \gamma$ : kernel ( $\gamma \in \mathbb{R}^{\geq 0}$ ). Why?  $\Leftarrow$  Gram.
- ② Cone. If  $k_m : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  kernel,  $\alpha_m \geq 0$  ( $m = 1, \dots, M$ ), then

$$\sum_{m=1}^M \alpha_m k_m(x, y) = ?$$

- We know:  $k(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y}$  is a kernel.
- A few useful rules:

- ① Non-negative shift.  $k$ : kernel  $\Rightarrow k + \gamma$ : kernel ( $\gamma \in \mathbb{R}^{\geq 0}$ ). Why?  $\Leftarrow$  Gram.
- ② Cone. If  $k_m : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  kernel,  $\alpha_m \geq 0$  ( $m = 1, \dots, M$ ), then

$$\sum_{m=1}^M \alpha_m k_m(x, y) = \sum_m \alpha_m \langle \varphi_m(x), \varphi_m(y) \rangle_{\mathcal{H}_m}$$

- We know:  $k(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y}$  is a kernel.
- A few useful rules:

- ① Non-negative shift.  $k$ : kernel  $\Rightarrow k + \gamma$ : kernel ( $\gamma \in \mathbb{R}^{\geq 0}$ ). Why?  $\Leftarrow$  Gram.
- ② Cone. If  $k_m : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  kernel,  $\alpha_m \geq 0$  ( $m = 1, \dots, M$ ), then

$$\begin{aligned}\sum_{m=1}^M \alpha_m k_m(x, y) &= \sum_m \alpha_m \langle \varphi_m(x), \varphi_m(y) \rangle_{\mathcal{H}_m} \\ &= \langle \varphi(x), \varphi(y) \rangle_{\mathcal{H}},\end{aligned}$$

- We know:  $k(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y}$  is a kernel.
- A few useful rules:

- ① Non-negative shift.  $k$ : kernel  $\Rightarrow k + \gamma$ : kernel ( $\gamma \in \mathbb{R}^{\geq 0}$ ). Why?  $\Leftarrow$  Gram.
- ② Cone. If  $k_m : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  kernel,  $\alpha_m \geq 0$  ( $m = 1, \dots, M$ ), then

$$\begin{aligned}\sum_{m=1}^M \alpha_m k_m(x, y) &= \sum_m \alpha_m \langle \varphi_m(x), \varphi_m(y) \rangle_{\mathcal{H}_m} \\ &= \langle \varphi(x), \varphi(y) \rangle_{\mathcal{H}}, \\ \varphi(x) &= (\sqrt{\alpha_1} \varphi_1(x), \dots, \sqrt{\alpha_M} \varphi_M(x)) \in \mathcal{H} := \bigoplus_{m=1}^M \mathcal{H}_m.\end{aligned}$$

- We know:  $k(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y}$  is a kernel.
- A few **useful rules**:

- ① **Non-negative shift.**  $k$ : kernel  $\Rightarrow k + \gamma$ : kernel ( $\gamma \in \mathbb{R}^{\geq 0}$ ). Why?  $\Leftarrow$  Gram.
- ② **Cone.** If  $k_m : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  kernel,  $\alpha_m \geq 0$  ( $m = 1, \dots, M$ ), then

$$\begin{aligned}\sum_{m=1}^M \alpha_m k_m(x, y) &= \sum_m \alpha_m \langle \varphi_m(x), \varphi_m(y) \rangle_{\mathcal{H}_m} \\ &= \langle \varphi(x), \varphi(y) \rangle_{\mathcal{H}},\end{aligned}$$

$$\varphi(x) = (\sqrt{\alpha_1} \varphi_1(x), \dots, \sqrt{\alpha_M} \varphi_M(x)) \in \mathcal{H} := \bigoplus_{m=1}^M \mathcal{H}_m.$$

Example:  $\bigoplus_{m=1}^M \mathbb{R} = \mathbb{R}^M$ .

- ④ **Product.** If  $(k_m)_{m=1}^M$  are kernels on  $\mathcal{X}_m$ , then

$$(\otimes_{m=1}^M k_m) \left( (x_1, \dots, x_M), (x'_1, \dots, x'_M) \right) = \prod_{m=1}^M k_m(x_m, x'_m).$$

- ④ **Product.** If  $(k_m)_{m=1}^M$  are kernels on  $\mathcal{X}_m$ , then

$$(\otimes_{m=1}^M k_m) \left( (x_1, \dots, x_M), (x'_1, \dots, x'_M) \right) = \prod_{m=1}^M k_m(x_m, x'_m).$$

- Thus,  $(k_m)_{m=1}^M : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  kernels  $\Rightarrow \prod_{m=1}^M k_m(x, x')$ : kernel on  $\mathcal{X}$ .

- ④ **Product.** If  $(k_m)_{m=1}^M$  are kernels on  $\mathcal{X}_m$ , then

$$(\otimes_{m=1}^M k_m) \left( (x_1, \dots, x_M), (x'_1, \dots, x'_M) \right) = \prod_{m=1}^M k_m(x_m, x'_m).$$

- Thus,  $(k_m)_{m=1}^M : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  kernels  $\Rightarrow \prod_{m=1}^M k_m(x, x')$ : kernel on  $\mathcal{X}$ .
- Recall:  $\otimes_{m=1}^M k_m$  will be in HSIC.

- ④ **Product.** If  $(k_m)_{m=1}^M$  are kernels on  $\mathcal{X}_m$ , then

$$(\otimes_{m=1}^M k_m) \left( (x_1, \dots, x_M), (x'_1, \dots, x'_M) \right) = \prod_{m=1}^M k_m(x_m, x'_m).$$

- Thus,  $(k_m)_{m=1}^M : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  kernels  $\Rightarrow \prod_{m=1}^M k_m(x, x')$ : kernel on  $\mathcal{X}$ .
- Recall:  $\otimes_{m=1}^M k_m$  will be in HSIC.
- Consequence ( $\gamma \geq 0$ ,  $p \in \mathbb{Z}^+$ ):

$$k(\mathbf{x}, \mathbf{y}) = (\langle \mathbf{x}, \mathbf{y} \rangle_2 + \gamma)^p$$

is a **kernel**.

Intuition for  $M = 2$  and assuming  $\varphi_m(x) \in \mathbb{R}^{d_m}$ :

$$(\textcolor{red}{k}_1 \otimes \textcolor{blue}{k}_2) ((x, y), (x', y')) = \textcolor{red}{k}_1(x, x') \textcolor{blue}{k}_2(y, y')$$

## Kernel factory: product indeed

Intuition for  $M = 2$  and assuming  $\varphi_m(x) \in \mathbb{R}^{d_m}$ :

$$\begin{aligned} (\textcolor{red}{k}_1 \otimes \textcolor{blue}{k}_2) ((x, y), (x', y')) &= \textcolor{red}{k}_1(x, x') \textcolor{blue}{k}_2(y, y') \\ &= \langle \varphi_1(x), \varphi_1(x') \rangle_{\mathcal{H}_1} \langle \varphi_2(x), \varphi_2(x') \rangle_{\mathcal{H}_2} \end{aligned}$$

## Kernel factory: product indeed

Intuition for  $M = 2$  and assuming  $\varphi_m(x) \in \mathbb{R}^{d_m}$ :

$$\begin{aligned} (\mathbf{k}_1 \otimes \mathbf{k}_2) ((x, y), (x', y')) &= \mathbf{k}_1(x, x') \mathbf{k}_2(y, y') \\ &= \langle \varphi_1(x), \varphi_1(x') \rangle_{\mathcal{H}_1} \langle \varphi_2(x), \varphi_2(x') \rangle_{\mathcal{H}_2} \\ &= \varphi_1(x)^T \varphi_1(x') \varphi_2(x)^T \varphi_2(x') \end{aligned}$$

## Kernel factory: product indeed

Intuition for  $M = 2$  and assuming  $\varphi_m(x) \in \mathbb{R}^{d_m}$ :

$$\begin{aligned} (\textcolor{red}{k}_1 \otimes \textcolor{blue}{k}_2)((x, y), (x', y')) &= \textcolor{red}{k}_1(x, x') \textcolor{blue}{k}_2(y, y') \\ &= \langle \varphi_1(x), \varphi_1(x') \rangle_{\mathcal{H}_1} \langle \varphi_2(x), \varphi_2(x') \rangle_{\mathcal{H}_2} \\ &= \varphi_1(x)^T \varphi_1(x') \varphi_2(x)^T \varphi_2(x') \\ &= \text{tr} \left( \varphi_1(x)^T \varphi_1(x') \varphi_2(x)^T \varphi_2(x') \right) \end{aligned}$$

## Kernel factory: product indeed

Intuition for  $M = 2$  and assuming  $\varphi_m(x) \in \mathbb{R}^{d_m}$ :

$$\begin{aligned} (\mathbf{k}_1 \otimes \mathbf{k}_2)((x, y), (x', y')) &= \mathbf{k}_1(x, x') \mathbf{k}_2(y, y') \\ &= \langle \varphi_1(x), \varphi_1(x') \rangle_{\mathcal{H}_1} \langle \varphi_2(x), \varphi_2(x') \rangle_{\mathcal{H}_2} \\ &= \varphi_1(x)^T \varphi_1(x') \varphi_2(x)^T \varphi_2(x') \\ &= \text{tr} \left( \varphi_1(x)^T \varphi_1(x') \varphi_2(x)^T \varphi_2(x') \right) \\ &= \text{tr} \left( \varphi_2(x') \varphi_1(x)^T \varphi_1(x') \varphi_2(x)^T \right) \end{aligned}$$

## Kernel factory: product indeed

Intuition for  $M = 2$  and assuming  $\varphi_m(x) \in \mathbb{R}^{d_m}$ :

$$\begin{aligned} (\mathbf{k}_1 \otimes \mathbf{k}_2)((x, y), (x', y')) &= \mathbf{k}_1(x, x') \mathbf{k}_2(y, y') \\ &= \langle \varphi_1(x), \varphi_1(x') \rangle_{\mathcal{H}_1} \langle \varphi_2(x), \varphi_2(x') \rangle_{\mathcal{H}_2} \\ &= \varphi_1(x)^T \varphi_1(x') \varphi_2(x)^T \varphi_2(x') \\ &= \text{tr} \left( \varphi_1(x)^T \varphi_1(x') \varphi_2(x)^T \varphi_2(x') \right) \\ &= \text{tr} \left( \varphi_2(x') \varphi_1(x)^T \varphi_1(x') \varphi_2(x)^T \right) \\ &= \underbrace{\left\| \varphi_1(x') \varphi_2(x)^T \right\|_F}_{\in \mathbb{R}^{d_1 \times d_2}}, \end{aligned}$$

where  $\|\mathbf{A}\|_F = \sqrt{\sum_{ij} \mathbf{A}_{ij}^2}$  is the Frobenius norm.

- ⑥ **Limit.** If  $(k_n)_{n \in \mathbb{N}}$  are kernels on  $\mathcal{X}$ , then

$$k(x, x') := \lim_{n \rightarrow \infty} k_n(x, x')$$

is a kernel. Why?

- ⑥ **Limit.** If  $(k_n)_{n \in \mathbb{N}}$  are kernels on  $\mathcal{X}$ , then

$$k(x, x') := \lim_{n \rightarrow \infty} k_n(x, x')$$

is a kernel. Why?  $\Leftarrow$  Gram.

- ⑥ **Limit.** If  $(k_n)_{n \in \mathbb{N}}$  are kernels on  $\mathcal{X}$ , then

$$k(x, x') := \lim_{n \rightarrow \infty} k_n(x, x')$$

is a kernel. Why?  $\Leftarrow$  Gram.

Example ( $\gamma > 0$ ):

$$k(\mathbf{x}, \mathbf{y}) = e^{\gamma \langle \mathbf{x}, \mathbf{y} \rangle_2} = \sum_{n \in \mathbb{N}} \frac{(\gamma \langle \mathbf{x}, \mathbf{y} \rangle_2)^n}{n!}$$

# Kernel factory

- ⑥ **Limit.** If  $(k_n)_{n \in \mathbb{N}}$  are kernels on  $\mathcal{X}$ , then

$$k(x, x') := \lim_{n \rightarrow \infty} k_n(x, x')$$

is a kernel. Why?  $\Leftarrow$  Gram.

Example ( $\gamma > 0$ ):

$$k(\mathbf{x}, \mathbf{y}) = e^{\gamma \langle \mathbf{x}, \mathbf{y} \rangle_2} = \sum_{n \in \mathbb{N}} \frac{(\gamma \langle \mathbf{x}, \mathbf{y} \rangle_2)^n}{n!}$$

Reason: polynomial kernel & limit rule.

## Kernel factory – continued

- ⑦ Pre-post multiplication.  $k$  kernel on  $\mathcal{X}$ ,  $f : \mathcal{X} \rightarrow \mathbb{R}$ , then

$$\tilde{k}(x, y) = f(x)k(x, y)f(y)$$

is a kernel. Check (feature view):

$$\tilde{k}(x, y) = f(x)\mathbf{k}(x, y)f(y)$$

- ⑦ Pre-post multiplication.  $k$  kernel on  $\mathcal{X}$ ,  $f : \mathcal{X} \rightarrow \mathbb{R}$ , then

$$\tilde{k}(x, y) = f(x)k(x, y)f(y)$$

is a kernel. Check (feature view):

$$\tilde{k}(x, y) = f(x)k(x, y)f(y) = f(x)\langle \varphi(x), \varphi(y) \rangle_{\mathcal{H}} f(y)$$

# Kernel factory – continued

- ⑦ Pre-post multiplication.  $k$  kernel on  $\mathcal{X}$ ,  $f : \mathcal{X} \rightarrow \mathbb{R}$ , then

$$\tilde{k}(x, y) = f(x)k(x, y)f(y)$$

is a kernel. Check (feature view):

$$\begin{aligned}\tilde{k}(x, y) &= f(x)k(x, y)f(y) = f(x)\langle\varphi(x), \varphi(y)\rangle_{\mathcal{H}}f(y) \\ &= \left\langle \underbrace{f(x)\varphi(x)}_{=: \tilde{\varphi}(x)}, f(y)\varphi(y) \right\rangle_{\mathcal{H}}.\end{aligned}$$

# Kernel factory – continued

- ⑦ Pre-post multiplication.  $k$  kernel on  $\mathcal{X}$ ,  $f : \mathcal{X} \rightarrow \mathbb{R}$ , then

$$\tilde{k}(x, y) = f(x)k(x, y)f(y)$$

is a kernel. Check (feature view):

$$\begin{aligned}\tilde{k}(x, y) &= f(x)\color{red}{k(x, y)}f(y) = f(x)\color{red}{\langle \varphi(x), \varphi(y) \rangle_{\mathcal{H}}}f(y) \\ &= \left\langle \underbrace{f(x)\varphi(x)}_{=: \tilde{\varphi}(x)}, f(y)\varphi(y) \right\rangle_{\mathcal{H}}.\end{aligned}$$

Example (Gaussian kernel,  $\gamma > 0$ ): previous example & new rule

$$k(\mathbf{x}, \mathbf{y}) = e^{-\gamma \|\mathbf{x}-\mathbf{y}\|_2^2}$$

by using  $\|\mathbf{x} - \mathbf{y}\|_2^2 = \|\mathbf{x}\|_2^2 + \|\mathbf{y}\|_2^2 - 2 \langle \mathbf{x}, \mathbf{y} \rangle$ .

## Kernel factory – continued

- 7 Pre-post multiplication.  $k$  kernel on  $\mathcal{X}$ ,  $f : \mathcal{X} \rightarrow \mathbb{R}$ , then

$$\tilde{k}(x, y) = f(x)k(x, y)f(y)$$

is a kernel. Check (feature view):

$$\begin{aligned}\tilde{k}(x, y) &= f(x)\mathbf{k}(x, y)f(y) = f(x)\langle\varphi(x), \varphi(y)\rangle_{\mathcal{H}}f(y) \\ &= \left\langle \underbrace{f(x)\varphi(x)}_{=: \tilde{\varphi}(x)}, f(y)\varphi(y) \right\rangle_{\mathcal{H}}.\end{aligned}$$

Example (Gaussian kernel,  $\gamma > 0$ ): previous example & new rule

$$k(\mathbf{x}, \mathbf{y}) = e^{-\gamma \|\mathbf{x}-\mathbf{y}\|_2^2} = e^{-\gamma \|\mathbf{x}\|^2} e^{2\gamma \langle \mathbf{x}, \mathbf{y} \rangle_2} e^{-\gamma \|\mathbf{x}\|^2}$$

by using  $\|\mathbf{x} - \mathbf{y}\|_2^2 = \|\mathbf{x}\|_2^2 + \|\mathbf{y}\|_2^2 - 2 \langle \mathbf{x}, \mathbf{y} \rangle$ .

# Kernel factory: $\mathbb{R}^d$ & Bochner theorem

We focus on continuous bounded shift-invariant kernels:

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

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

where  $\Lambda$  is a finite Borel measure (w.l.o.g. probability).

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

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

kernel name $k_0$	$\hat{k}_0(\omega)$
Gaussian	$e^{-\frac{x^2}{2\sigma^2}}$
Laplacian	$e^{-\sigma x }$
$B_{2n+1}$ -spline	$*^{2n+2}\chi_{[-\frac{1}{2}, \frac{1}{2}]}(x) \frac{4^{n+1}}{\sqrt{2\pi}} \frac{\sin^{2n+2}(\frac{\omega}{2})}{\omega^{2n+2}}$
Sinc	$\frac{\sin(\sigma x)}{x}$
Poisson	$\frac{1-\sigma^2}{\sigma^2 - 2\sigma \cos(x) + 1}$
Dirichlet	$\frac{\sin(\frac{(2n+1)x}{2})}{\sin(\frac{x}{2})}$
Fejér	$\frac{1}{n+1} \frac{\sin^2 \frac{(n+1)x}{2}}{\sin^2(\frac{x}{2})}$
Cosine	$\cos(\sigma x)$
	$\sigma e^{-\frac{\sigma^2 \omega^2}{2}}$
	$\sqrt{\frac{2}{\pi}} \frac{\sigma}{\sigma^2 + \omega^2}$
	$\sqrt{\frac{\pi}{2}} \chi_{[-\sigma, \sigma]}(\omega)$
	$\sqrt{2\pi} \sum_{j=-\infty}^{\infty} \sigma^{ j } \delta(\omega - j)$
	$\sqrt{2\pi} \sum_{j=-\infty}^{\infty} \delta(\omega - j)$
	$\sqrt{2\pi} \sum_{j=-n}^n \left(1 - \frac{ j }{n+1}\right) \delta(\omega - j)$
	$\sqrt{\frac{\pi}{2}} [\delta(\omega - \sigma) + \delta(\omega + \sigma)]$

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

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

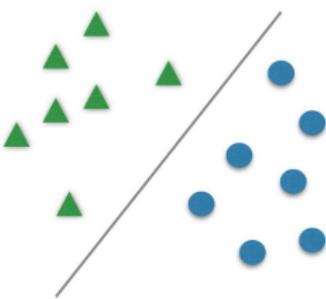
kernel name $k_0$	$\hat{k}_0(\omega)$
Gaussian	$e^{-\frac{x^2}{2\sigma^2}}$
Laplacian	$e^{-\sigma x }$
$B_{2n+1}$ -spline	$*^{2n+2}\chi_{[-\frac{1}{2}, \frac{1}{2}]}(x) \frac{4^{n+1}}{\sqrt{2\pi}} \frac{\sin^{2n+2}(\frac{\omega}{2})}{\omega^{2n+2}}$
Sinc	$\frac{\sin(\sigma x)}{x}$
Poisson	$\frac{1-\sigma^2}{\sigma^2 - 2\sigma \cos(x) + 1}$
Dirichlet	$\frac{\sin(\frac{(2n+1)x}{2})}{\sin(\frac{x}{2})}$
Fejér	$\frac{1}{n+1} \frac{\sin^2 \frac{(n+1)x}{2}}{\sin^2(\frac{x}{2})}$
Cosine	$\cos(\sigma x)$
	$\sigma e^{-\frac{\sigma^2 \omega^2}{2}}$
	$\sqrt{\frac{2}{\pi}} \frac{\sigma}{\sigma^2 + \omega^2}$
	$\sqrt{\frac{\pi}{2}} \chi_{[-\sigma, \sigma]}(\omega)$
	$\sqrt{2\pi} \sum_{j=-\infty}^{\infty} \sigma^{ j } \delta(\omega - j)$
	$\sqrt{2\pi} \sum_{j=-\infty}^{\infty} \delta(\omega - j)$
	$\sqrt{2\pi} \sum_{j=-n}^n \left(1 - \frac{ j }{n+1}\right) \delta(\omega - j)$
	$\sqrt{\frac{\pi}{2}} [\delta(\omega - \sigma) + \delta(\omega + \sigma)]$

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

# Kernels in action: classification, regression, dimensionality reduction

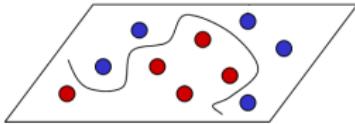
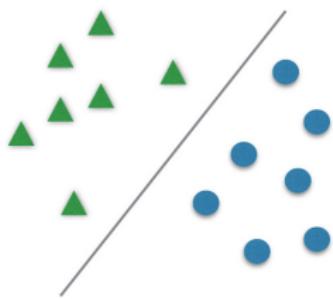
# Classification , regression

- Given:  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ ,  $y_i \in \{-1, 1\}$ .
- Goal: find an  $f$  classifier such that  $f(\mathbf{x}) \approx y$ .



# Classification, regression

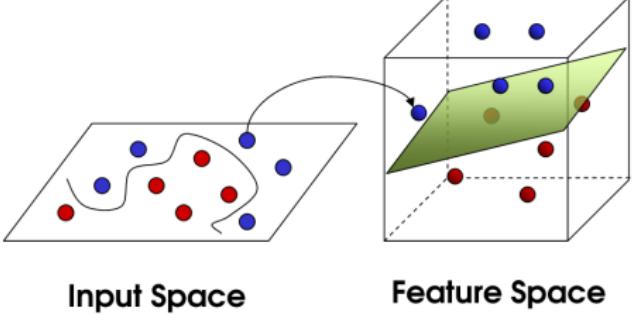
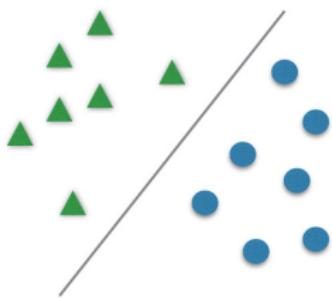
- Given:  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ ,  $y_i \in \{-1, 1\}$ .
- Goal: find an  $f$  classifier such that  $f(\mathbf{x}) \approx y$ .



Input Space

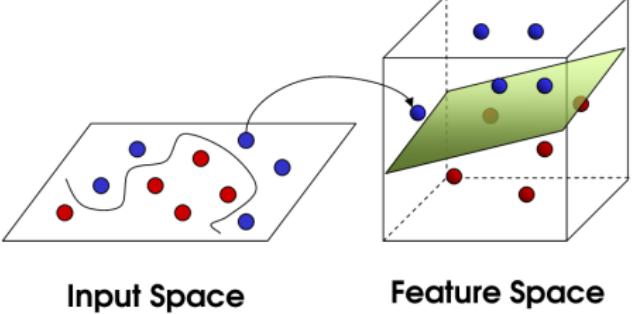
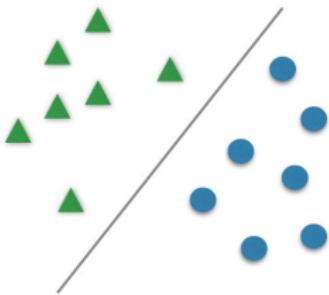
# Classification, regression

- Given:  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ ,  $y_i \in \{-1, 1\}$ .
- Goal: find an  $f$  classifier such that  $f(\mathbf{x}) \approx y$ .



# Classification, regression

- Given:  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ ,  $y_i \in \{-1, 1\}$ .
- Goal: find an  $f$  classifier such that  $f(\mathbf{x}) \approx y$ .
- Regression similarly:  $y_i \in \mathbb{R}$ .



# Dimensionality reduction: intuition

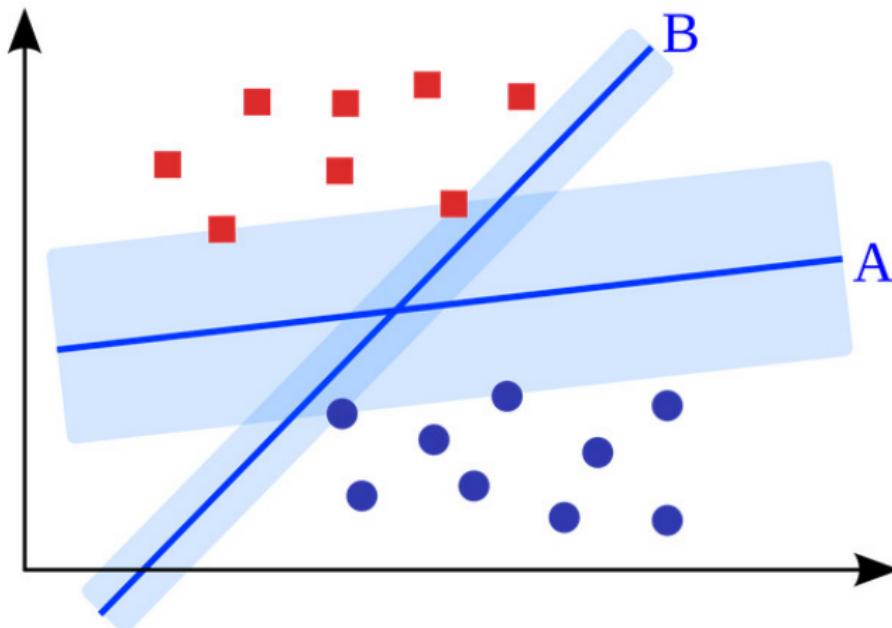
- Given: a set of observations  $X = \{\mathbf{x}_i\}_{i=1}^n \subset \mathbb{R}^D$ .
- Goal: find  $X' = \{\mathbf{x}'_i\}_{i=1}^n \subset \mathbb{R}^d$  'preserving' the geometry of  $X$ .
- $d \ll D$ : compression (images, music, ...).



# Classification: SVM

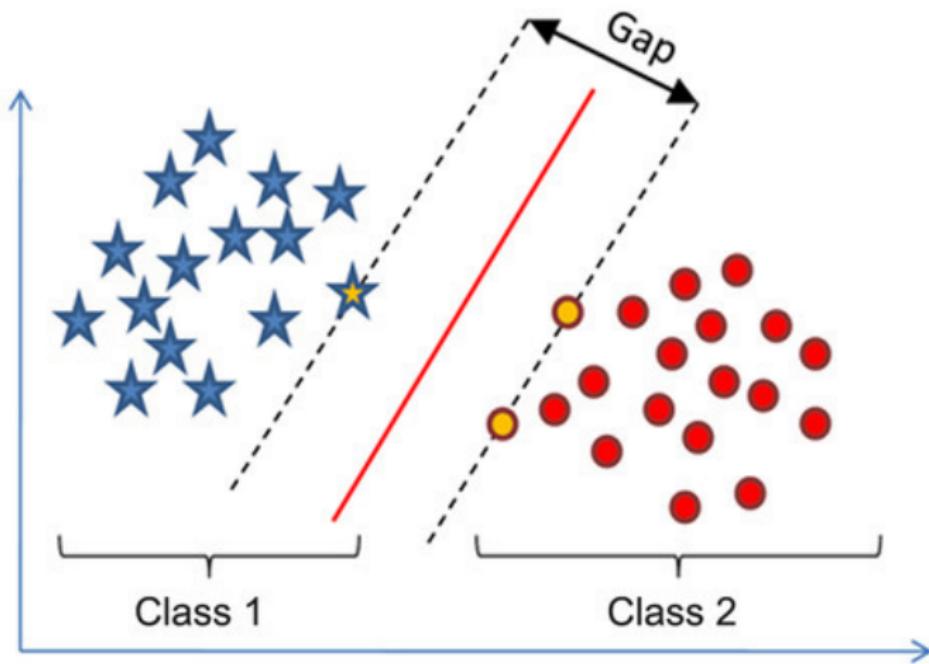
# Support vector machine (SVM) for classification

Which separating line is the 'best'?



# Support Vector Machine (SVM)

SVM answer: the one with the largest margin.



## SVM formulation: hard classification

- Hyperplane:  $f_{\mathbf{w}, b}(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + b$ .
  - $\mathbf{w}$ : normal vector,  $b$ : offset.

## SVM formulation: hard classification

- Hyperplane:  $f_{\mathbf{w}, b}(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + b$ .

- $\mathbf{w}$ : normal vector,  $b$ : offset.

- Goal:

$$\max_{\mathbf{w}, b} \underbrace{\frac{2}{\|\mathbf{w}\|_2}}_{\text{margin}} \Leftrightarrow \min \|\mathbf{w}\|_2^2, \text{ s.t. } \underbrace{\begin{cases} \langle \mathbf{w}, \mathbf{x}_i \rangle + b \geq 1 & \text{if } y_i = 1, \\ \langle \mathbf{w}, \mathbf{x}_i \rangle + b \leq -1 & \text{otherwise.} \end{cases}}_{\text{correct classification}}$$

## SVM formulation: hard classification

- Hyperplane:  $f_{\mathbf{w}, b}(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + b$ .

- $\mathbf{w}$ : normal vector,  $b$ : offset.

- Goal:

$$\max_{\mathbf{w}, b} \underbrace{\frac{2}{\|\mathbf{w}\|_2}}_{\text{margin}} \Leftrightarrow \min_{\mathbf{w}, b} \|\mathbf{w}\|_2^2, \text{ s.t. } \underbrace{\begin{cases} \langle \mathbf{w}, \mathbf{x}_i \rangle + b \geq 1 & \text{if } y_i = 1, \\ \langle \mathbf{w}, \mathbf{x}_i \rangle + b \leq -1 & \text{otherwise.} \end{cases}}_{\text{correct classification}}$$

- Shortly,

$$\min_{\mathbf{w}, b} \|\mathbf{w}\|_2^2, \text{ s.t. } y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1, \forall i.$$

## SVM formulation: hard classification

- Hyperplane:  $f_{\mathbf{w}, b}(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + b$ .

- $\mathbf{w}$ : normal vector,  $b$ : offset.

- Goal:

$$\max_{\mathbf{w}, b} \underbrace{\frac{2}{\|\mathbf{w}\|_2}}_{\text{margin}} \Leftrightarrow \min_{\mathbf{w}, b} \|\mathbf{w}\|_2^2, \text{ s.t. } \underbrace{\begin{cases} \langle \mathbf{w}, \mathbf{x}_i \rangle + b \geq 1 & \text{if } y_i = 1, \\ \langle \mathbf{w}, \mathbf{x}_i \rangle + b \leq -1 & \text{otherwise.} \end{cases}}_{\text{correct classification}}$$

- Shortly,

$$\min_{\mathbf{w}, b} \|\mathbf{w}\|_2^2, \text{ s.t. } y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1, \forall i.$$

- Decision:  $\hat{y} = \text{sign}(\langle \mathbf{w}, \mathbf{x} \rangle + b)$ .

## SVM formulation: soft classification

- Hard classification objective:

$$\min_{\mathbf{w}, b} \|\mathbf{w}\|_2^2 \text{ s.t. } y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1, \forall i.$$

There might not be solution! (non-linearly separable case)

# SVM formulation: soft classification

- Hard classification objective:

$$\min_{\mathbf{w}, b} \|\mathbf{w}\|_2^2 \text{ s.t. } y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1, \forall i.$$

There might not be solution! (non-linearly separable case)

- Soft classification objective:

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \xi_i \text{ s.t. } y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 - \xi_i, \xi_i \geq 0, \forall i.$$

Linear penalty on misclassification.

# SVM formulation: soft classification

Soft classification objective:

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \xi_i, \text{ s.t. } y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 - \xi_i, \xi_i \geq 0 \quad (\forall i).$$

Lagrangian function: with  $\alpha_i \geq 0, \beta_i \geq 0 \quad (\forall i)$

$L(\mathbf{w}, b, \xi; \alpha, \beta) = \text{objective} - \text{Lagrangian multipliers} \times \text{conditions}$

$$= \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i [y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) - 1 + \xi_i] - \sum_{i=1}^n \beta_i \xi_i.$$

Solving for  $\frac{\partial L}{\partial \text{primal}} = 0$ , we get ...

## SVM formulation: soft classification

$$L(\mathbf{w}, b, \xi; \alpha, \beta) =$$

$$= \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i [y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) - 1 + \xi_i] - \sum_{i=1}^n \beta_i \xi_i.$$

Optimality equations:

$$\mathbf{0} = \frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \quad (\mathbf{w} \leftrightarrow \alpha),$$

$$0 = \frac{\partial L}{\partial b} = \sum_{i=1}^n \alpha_i y_i,$$

$$0 = \frac{\partial L}{\partial \xi_i} = C - \alpha_i - \beta_i.$$

Plugging these equations back to  $L$ , we have . . .

# SVM formulation: soft classification

Dual form:

$$\max_{\alpha} \underbrace{\sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j}_{\text{quadratic in } \alpha}, \text{ s.t. } \underbrace{0 \leq \alpha_i \leq C, \sum_{i=1}^n \alpha_i y_i = 0}_{\text{linear in } \alpha}.$$

# SVM formulation: soft classification

Dual form:

$$\max_{\alpha} \underbrace{\sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j}_{\text{quadratic in } \alpha}, \text{ s.t. } 0 \leq \alpha_i \leq C, \underbrace{\sum_{i=1}^n \alpha_i y_i = 0}_{\text{linear in } \alpha}.$$

- $b \Leftarrow y_i(\mathbf{w}_i^T \mathbf{x}_i + b) = 1 \Leftarrow \alpha_i > 0$  [complementary slackness].

# SVM formulation: soft classification

Dual form:

$$\max_{\alpha} \underbrace{\sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j}_{\text{quadratic in } \alpha}, \text{ s.t. } 0 \leq \alpha_i \leq C, \underbrace{\sum_{i=1}^n \alpha_i y_i = 0}_{\text{linear in } \alpha}.$$

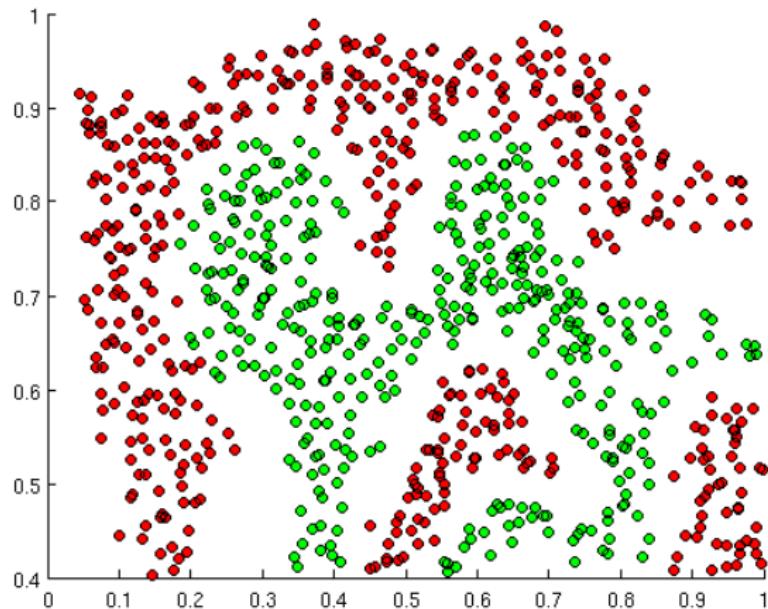
- $b \Leftarrow y_i(\mathbf{w}_i^T \mathbf{x}_i + b) = 1 \Leftarrow \alpha_i > 0$  [complementary slackness].
- QP: solvers are available.

# If linear separability does not hold

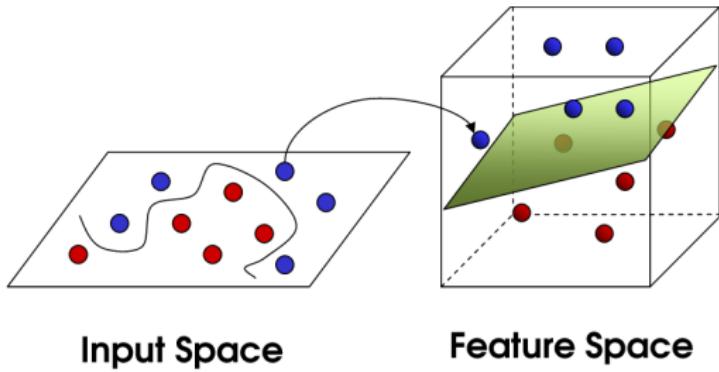
- Until this point:
  - (almost) **linearly separable** case.

# If linear separability does not hold

- Until this point:
  - (almost) **linearly separable** case.
- Now:



If linear separability does not hold: **kernel trick**



Input Space

Feature Space

- Linear SVM (dual):

$$\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j, \text{ s.t. } \sum_{i=1}^n \alpha_i y_i = 0, 0 \leq \alpha_i \leq C (\forall i).$$

# Nonlinear SVM

- Linear SVM (dual):

$$\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j, \text{ s.t. } \sum_{i=1}^n \alpha_i y_i = 0, 0 \leq \alpha_i \leq C (\forall i).$$

- Nonlinear SVM (dual):

$$\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j), \text{ s.t. } \sum_{i=1}^n \alpha_i y_i = 0, 0 \leq \alpha_i \leq C (\forall i).$$

# Nonlinear SVM

- Linear SVM (dual):

$$\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j, \text{ s.t. } \sum_{i=1}^n \alpha_i y_i = 0, 0 \leq \alpha_i \leq C (\forall i).$$

- Nonlinear SVM (dual):

$$\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{k}(\mathbf{x}_i, \mathbf{x}_j), \text{ s.t. } \sum_{i=1}^n \alpha_i y_i = 0, 0 \leq \alpha_i \leq C (\forall i).$$

- Nonlinear SVM (primal):

$$\min_{f \in \mathcal{H}_k, \xi} \frac{1}{2} \|f\|_{\mathcal{H}_k}^2 + C \sum_{i=1}^n \xi_i, \text{ s.t. } y_i f(x_i) \geq 1 - \xi_i, \xi_i \geq 0, \forall i.$$

# Kernel ridge regression

# Kernel ridge regression

- Given:  $\{(x_i, y_i)\}_{i=1}^n$ ,  $\mathcal{H} := \mathcal{H}_k$ ,  $y_i \in \mathbb{R}$ .
- Task ( $\lambda > 0$ ):

$$J(f) = \frac{1}{n} \sum_{i=1}^n [y_i - f(x_i)]^2 + \lambda \|f\|_{\mathcal{H}}^2 \rightarrow \min_{f \in \mathcal{H}}.$$

# Kernel ridge regression

- Given:  $\{(x_i, y_i)\}_{i=1}^n$ ,  $\mathcal{H} := \mathcal{H}_k$ ,  $y_i \in \mathbb{R}$ .
- Task ( $\lambda > 0$ ):

$$J(f) = \frac{1}{n} \sum_{i=1}^n [y_i - f(x_i)]^2 + \lambda \|f\|_{\mathcal{H}}^2 \rightarrow \min_{f \in \mathcal{H}}.$$

- Analytical solution (raised at distribution regression):

$$f(x) = [k(x_1, x), \dots, k(x_n, x)] (\mathbf{G} + \lambda n I)^{-1} [y_1; \dots; y_n],$$
$$\mathbf{G} = [k(x_i, x_j)]_{i,j=1}^n.$$

# Kernel ridge regression

- Given:  $\{(x_i, y_i)\}_{i=1}^n$ ,  $\mathcal{H} := \mathcal{H}_k$ ,  $y_i \in \mathbb{R}$ .
- Task ( $\lambda > 0$ ):

$$J(f) = \frac{1}{n} \sum_{i=1}^n [y_i - f(x_i)]^2 + \lambda \|f\|_{\mathcal{H}}^2 \rightarrow \min_{f \in \mathcal{H}}.$$

- Analytical solution (raised at distribution regression):

$$f(x) = [k(x_1, x), \dots, k(x_n, x)] (\mathbf{G} + \lambda n I)^{-1} [y_1; \dots; y_n],$$
$$\mathbf{G} = [k(x_i, x_j)]_{i,j=1}^n.$$

Question

How do we get this solution?

# Kernel ridge regression

By the representer theorem

$$f(x) = \sum_{i=1}^n a_i k(\cdot, x_i).$$

# Kernel ridge regression

By the representer theorem

$$f(x) = \sum_{i=1}^n a_i k(\cdot, x_i).$$

Multiplying the objective by  $n$ , using the reproducing property:

$$\begin{aligned} J(f) &= \sum_{j=1}^n [y_j - \langle f, k(\cdot, x_j) \rangle_{\mathcal{H}}]^2 + \lambda n \|f\|_{\mathcal{H}}^2 \\ &= \|\mathbf{y} - \mathbf{G}\mathbf{a}\|_2^2 + (\lambda n)\mathbf{a}^T \mathbf{G}\mathbf{a} \\ &= \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{G}\mathbf{a} + \mathbf{a}^T [\mathbf{G}^2 + (\lambda n)\mathbf{G}]\mathbf{a}. \end{aligned}$$

# Kernel ridge regression

By the representer theorem

$$f(x) = \sum_{i=1}^n a_i k(\cdot, x_i).$$

Multiplying the objective by  $n$ , using the reproducing property:

$$\begin{aligned} J(f) &= \sum_{j=1}^n [y_j - \langle f, k(\cdot, x_j) \rangle_{\mathcal{H}}]^2 + \lambda n \|f\|_{\mathcal{H}}^2 \\ &= \|\mathbf{y} - \mathbf{G}\mathbf{a}\|_2^2 + (\lambda n)\mathbf{a}^T \mathbf{G}\mathbf{a} \\ &= \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{G}\mathbf{a} + \mathbf{a}^T [\mathbf{G}^2 + (\lambda n)\mathbf{G}]\mathbf{a}. \end{aligned}$$

Solving  $\mathbf{0} = \frac{\partial J}{\partial \mathbf{a}}$ , one gets  $\mathbf{a}^* = (\mathbf{G} + \lambda n \mathbf{I})^{-1} \mathbf{y}$

# Kernel ridge regression

By the representer theorem

$$f(x) = \sum_{i=1}^n a_i k(\cdot, x_i).$$

Multiplying the objective by  $n$ , using the reproducing property:

$$\begin{aligned} J(f) &= \sum_{j=1}^n [y_j - \langle f, k(\cdot, x_j) \rangle_{\mathcal{H}}]^2 + \lambda n \|f\|_{\mathcal{H}}^2 \\ &= \|\mathbf{y} - \mathbf{G}\mathbf{a}\|_2^2 + (\lambda n)\mathbf{a}^T \mathbf{G}\mathbf{a} \\ &= \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{G}\mathbf{a} + \mathbf{a}^T [\mathbf{G}^T \mathbf{G} + (\lambda n)\mathbf{G}]\mathbf{a}. \end{aligned}$$

Solving  $\mathbf{0} = \frac{\partial J}{\partial \mathbf{a}}$ , one gets  $\mathbf{a}^* = (\mathbf{G} + \lambda n \mathbf{I})^{-1} \mathbf{y}$  by

$$\frac{\partial \mathbf{a}^T \mathbf{B} \mathbf{a}}{\partial \mathbf{a}} = (\mathbf{B} + \mathbf{B}^T) \mathbf{a}, \quad \frac{\partial \mathbf{c}^T \mathbf{a}}{\partial \mathbf{a}} = \mathbf{c}.$$

- Motivation: infoT objectives, hypothesis testing.
- Kernels, RKHS: definitions, construction.
- Kernel applications: classification, ridge regression.

# Notes

# Properties of $k$ control that of $\mathcal{H}_k$

[Steinwart and Christmann, 2008, Chapter 4]:

- $k$ : bounded  $[\sup_{x,y \in \mathcal{X}} k(x,y) \leq C] \Rightarrow \forall f \in \mathcal{H}_k$  is bounded

# Properties of $k$ control that of $\mathcal{H}_k$

[Steinwart and Christmann, 2008, Chapter 4]:

- $k$ : bounded  $[\sup_{x,y \in \mathcal{X}} k(x,y) \leq C] \Rightarrow \forall f \in \mathcal{H}_k$  is bounded:

$$|f(x)| = \left| \langle f, k(\cdot, x) \rangle_{\mathcal{H}_k} \right| \stackrel{\text{CBS}}{\leq} \|f\|_{\mathcal{H}_k} \underbrace{\|k(\cdot, x)\|_{\mathcal{H}_k}}_{\sqrt{k(x,x)}}.$$

## Properties of $k$ control that of $\mathcal{H}_k$

[Steinwart and Christmann, 2008, Chapter 4]:

- $k$ : bounded  $[\sup_{x,y \in \mathcal{X}} k(x,y) \leq C] \Rightarrow \forall f \in \mathcal{H}_k$  is bounded:

$$|f(x)| = \left| \langle f, k(\cdot, x) \rangle_{\mathcal{H}_k} \right| \stackrel{\text{CBS}}{\leq} \|f\|_{\mathcal{H}_k} \underbrace{\|k(\cdot, x)\|_{\mathcal{H}_k}}_{\sqrt{k(x,x)}}.$$

- $k$ : continuous  $\Rightarrow \mathcal{H}_k$ : separable  $[\ell^2(\mathbb{N})]$ .

# Properties of $k$ control that of $\mathcal{H}_k$

[Steinwart and Christmann, 2008, Chapter 4]:

- $k$ : bounded  $[\sup_{x,y \in \mathcal{X}} k(x,y) \leq C] \Rightarrow \forall f \in \mathcal{H}_k$  is bounded:

$$|f(x)| = \left| \langle f, k(\cdot, x) \rangle_{\mathcal{H}_k} \right| \stackrel{\text{CBS}}{\leq} \|f\|_{\mathcal{H}_k} \underbrace{\|k(\cdot, x)\|_{\mathcal{H}_k}}_{\sqrt{k(x,x)}}.$$

- $k$ : continuous  $\Rightarrow \mathcal{H}_k$ : separable  $[\ell^2(\mathbb{N})]$ .
- $k$ : bounded and continuous  $\Rightarrow \forall f \in \mathcal{H}_k$  is bounded & continuous.

# Properties of $k$ control that of $\mathcal{H}_k$

[Steinwart and Christmann, 2008, Chapter 4]:

- $k$ : bounded  $[\sup_{x,y \in \mathcal{X}} k(x,y) \leq C] \Rightarrow \forall f \in \mathcal{H}_k$  is bounded:

$$|f(x)| = \left| \langle f, k(\cdot, x) \rangle_{\mathcal{H}_k} \right| \stackrel{\text{CBS}}{\leq} \|f\|_{\mathcal{H}_k} \underbrace{\|k(\cdot, x)\|_{\mathcal{H}_k}}_{\sqrt{k(x,x)}}.$$

- $k$ : continuous  $\Rightarrow \mathcal{H}_k$ : separable  $[\ell^2(\mathbb{N})]$ .
- $k$ : bounded and continuous  $\Rightarrow \forall f \in \mathcal{H}_k$  is bounded & continuous.
- $k \in C^m \Rightarrow \forall f \in \mathcal{H}_k$  is  $m$ -times continuously differentiable.

# Properties of $k$ control that of $\mathcal{H}_k$

[Steinwart and Christmann, 2008, Chapter 4]:

- $k$ : bounded  $[\sup_{x,y \in \mathcal{X}} k(x,y) \leq C] \Rightarrow \forall f \in \mathcal{H}_k$  is bounded:

$$|f(x)| = \left| \langle f, k(\cdot, x) \rangle_{\mathcal{H}_k} \right| \stackrel{\text{CBS}}{\leq} \|f\|_{\mathcal{H}_k} \underbrace{\|k(\cdot, x)\|_{\mathcal{H}_k}}_{\sqrt{k(x,x)}}.$$

- $k$ : continuous  $\Rightarrow \mathcal{H}_k$ : separable  $[\ell^2(\mathbb{N})]$ .
- $k$ : bounded and continuous  $\Rightarrow \forall f \in \mathcal{H}_k$  is bounded & continuous.
- $k \in C^m \Rightarrow \forall f \in \mathcal{H}_k$  is  $m$ -times continuously differentiable.
- $k$ : analytic  $\Rightarrow \forall f \in \mathcal{H}_k$  is analytic.

# Hard vs soft-SVM classification

Recall:

- Hard SVM:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2, \text{ s.t. } y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1, \forall i.$$

- Soft SVM:

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \xi_i \text{ s.t. } y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 - \xi_i, \xi_i \geq 0, \forall i$$

# Hard vs soft-SVM classification

Recall:

- Hard SVM:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2, \text{ s.t. } y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1, \forall i.$$

- Soft SVM:

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \xi_i \quad \text{s.t. } y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 - \xi_i, \xi_i \geq 0, \forall i. \Leftrightarrow$$

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \max(1 - y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b), 0),$$

# Hard vs soft-SVM classification

Recall:

- Hard SVM:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2, \text{ s.t. } y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1, \forall i.$$

- Soft SVM:

$$\begin{aligned} & \min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \xi_i \quad \text{s.t. } y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 - \xi_i, \xi_i \geq 0, \forall i. \Leftrightarrow \\ & \min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \underbrace{\max \left( 1 - y_i \underbrace{(\langle \mathbf{w}, \mathbf{x}_i \rangle + b)}_{=f(\mathbf{x}_i)}, 0 \right)}_{=:h(y_i f(\mathbf{x}_i))}, \end{aligned}$$

where  $h(u) = \max(1 - u, 0)$  is the hinge loss.

## Hard vs soft-SVM classification – continued

The hinge loss is the convex envelope of the zero-one loss :

$$\textcolor{red}{z}(u) = \mathbb{I}_{u < 0},$$

$$u = y_i f(x_i),$$

$$\textcolor{blue}{h}(u) = \max(1 - u, 0).$$

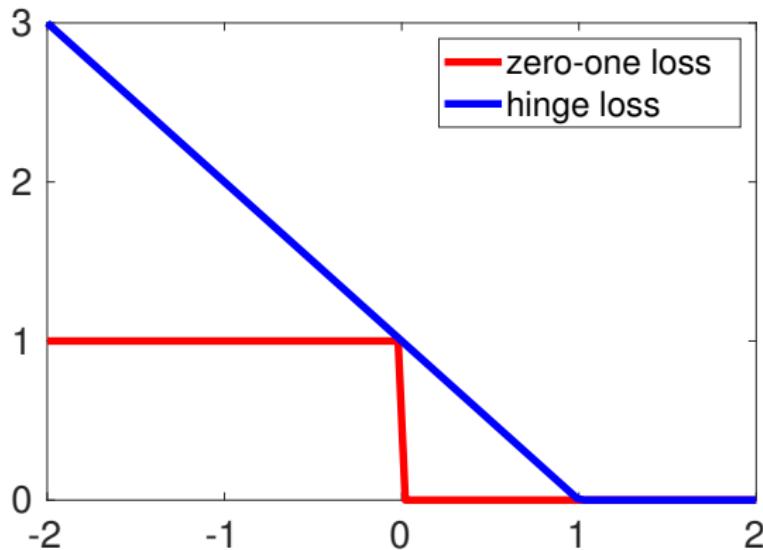
## Hard vs soft-SVM classification – continued

The hinge loss is the convex envelope of the zero-one loss:

$$z(u) = \mathbb{I}_{u < 0},$$

$$u = y_i f(x_i),$$

$$h(u) = \max(1 - u, 0).$$



## Representer theorem

[Schölkopf et al., 2001, Yu et al., 2013]

- Given:  $\{(x_i, y_i)\}_{i=1}^n$ , say classification/regression.
- Goal:

$$J(f) = V(x_1, y_1, f(x_1), \dots, x_n, y_n, f(x_n)) + r\left(\|f\|_{\mathcal{H}_k}^2\right) \rightarrow \min_{\mathcal{H}_k},$$

$r$  : monotonically increasing.

## Representer theorem

[Schölkopf et al., 2001, Yu et al., 2013]

- Given:  $\{(x_i, y_i)\}_{i=1}^n$ , say classification/regression.
- Goal:

$$J(f) = V(x_1, y_1, f(x_1), \dots, x_n, y_n, f(x_n)) + r\left(\|f\|_{\mathcal{H}_k}^2\right) \rightarrow \min_{\mathcal{H}_k},$$

$r$  : monotonically increasing.

- Example:

$$V(\dots) = \frac{1}{n} \sum_{i=1}^n \max(1 - y_i f(x_i), 0) \quad (\text{soft classification}),$$

$$V(\dots) = \frac{1}{n} \sum_{i=1}^n [f(x_i) - y_i]^2 \quad (\text{regression}).$$

# Representer theorem – continued

. . . then

- $\exists$  solution in the form:

$$f^* = \sum_{i=1}^n \alpha_i k(\cdot, x_i), \quad \alpha_i \in \mathbb{R}.$$

- $r$ : strictly increasing  $\Rightarrow \forall$  solution is of this form.
- Example:  $r(z) = \lambda z$ ,  $\lambda > 0$ .

# Representer theorem – proof

Objective

$$J(f) = V(x_1, y_1, f(x_1), \dots, x_n, y_n, f(x_n)) + r(\|f\|_{\mathcal{H}_k}^2) \rightarrow \min_{\mathcal{H}_k} .$$

Decompose & Pythagorean theorem:

$$S = \text{span}(k(\cdot, x_i), i = 1, \dots, n),$$

$$f = f_S + f_{\perp},$$

$$\|f\|_{\mathcal{H}_k}^2 = \|f_S\|_{\mathcal{H}_k}^2 + \underbrace{\|f_{\perp}\|_{\mathcal{H}_k}^2}_{\geq 0} \geq \|f_S\|_{\mathcal{H}_k}^2.$$

# Representer theorem – proof

## Objective

$$J(f) = V(x_1, y_1, f(x_1), \dots, x_n, y_n, f(x_n)) + r(\|f\|_{\mathcal{H}_k}^2) \rightarrow \min_{\mathcal{H}_k}.$$

Decompose & Pythagorean theorem:

$$S = \text{span}(k(\cdot, x_i), i = 1, \dots, n),$$

$$f = f_S + f_{\perp},$$

$$\|f\|_{\mathcal{H}_k}^2 = \|f_S\|_{\mathcal{H}_k}^2 + \underbrace{\|f_{\perp}\|_{\mathcal{H}_k}^2}_{\geq 0} \geq \|f_S\|_{\mathcal{H}_k}^2.$$

In  $J$

- 1st term: depends on  $f_S$  only,  $f(x_i) = \langle f, k(\cdot, x_i) \rangle$ .

# Representer theorem – proof

## Objective

$$J(f) = V(x_1, y_1, f(x_1), \dots, x_n, y_n, f(x_n)) + r \left( \|f\|_{\mathcal{H}_k}^2 \right) \rightarrow \min_{\mathcal{H}_k} .$$

Decompose & Pythagorean theorem:

$$S = \text{span}(k(\cdot, x_i), i = 1, \dots, n),$$

$$f = f_S + f_{\perp},$$

$$\|f\|_{\mathcal{H}_k}^2 = \|f_S\|_{\mathcal{H}_k}^2 + \underbrace{\|f_{\perp}\|_{\mathcal{H}_k}^2}_{\geq 0} \geq \|f_S\|_{\mathcal{H}_k}^2.$$

In  $J$

- 1st term: depends on  $f_S$  only,  $f(x_i) = \langle f, k(\cdot, x_i) \rangle$ .
- 2nd term: can only decrease by neglecting  $f_{\perp}$  ( $r \nearrow$ ).

$M$ -fold cross-validation [ $\theta := (C, \sigma)$ ]:

① Split data:

- training set ( $X_{tr}, Y_{tr}$ ):  $X_{val,i}, Y_{val,i}, i = 1, \dots, M$ .
- test set:  $X_{te}, Y_{te}$ .

$M$ -fold cross-validation [ $\theta := (C, \sigma)$ ]:

① Split data:

- training set  $(X_{tr}, Y_{tr})$ :  $X_{val,i}, Y_{val,i}$ ,  $i = 1, \dots, M$ .
- test set:  $X_{te}, Y_{te}$ .

② For fixed  $\theta$ : evaluate the average error while

- trained on:  $X_{tr} \setminus X_{val,i}, Y_{tr} \setminus Y_{val,i}$ ,
- tested on:  $X_{val,i}, Y_{val,i}$ .

$M$ -fold cross-validation  $[\theta := (C, \sigma)]$ :

- ① Split data:
  - training set  $(X_{tr}, Y_{tr})$ :  $X_{val,i}, Y_{val,i}$ ,  $i = 1, \dots, M$ .
  - test set:  $X_{te}, Y_{te}$ .
- ② For fixed  $\theta$ : evaluate the average error while
  - trained on:  $X_{tr} \setminus X_{val,i}, Y_{tr} \setminus Y_{val,i}$ ,
  - tested on:  $X_{val,i}, Y_{val,i}$ .
- ③  $\theta^* :=$  minimizer of CV error.

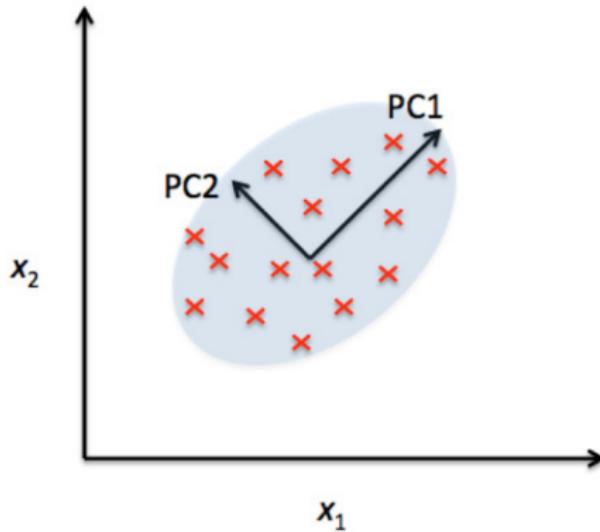
$M$ -fold cross-validation [ $\theta := (C, \sigma)$ ]:

- ① Split data:
  - training set  $(X_{tr}, Y_{tr})$ :  $X_{val,i}, Y_{val,i}$ ,  $i = 1, \dots, M$ .
  - test set:  $X_{te}, Y_{te}$ .
- ② For fixed  $\theta$ : evaluate the average error while
  - trained on:  $X_{tr} \setminus X_{val,i}, Y_{tr} \setminus Y_{val,i}$ ,
  - tested on:  $X_{val,i}, Y_{val,i}$ .
- ③  $\theta^* :=$  minimizer of CV error.
- ④ Report: performance of  $\theta^*$  on  $X_{te}, Y_{te}$ .

# PCA and its kernelized version

# PCA: intuition

Task: find the best  $d$ -dimensional subspace approximating  $\{\mathbf{x}_i\}_{i=1}^n \subset \mathbb{R}^D$ .



# PCA example: 100%

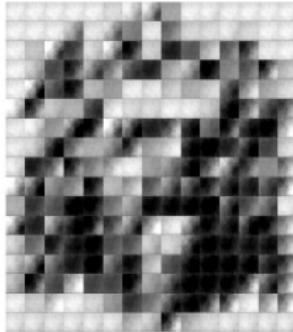


(A)

# PCA example: 100% → 1%



(A)

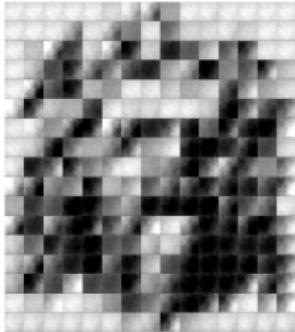


(B)

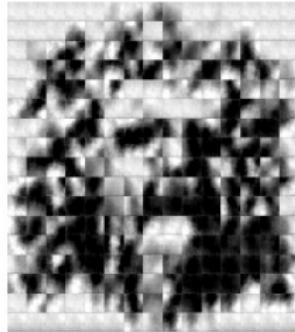
# PCA example: 100% → 2%



(A)



(B)

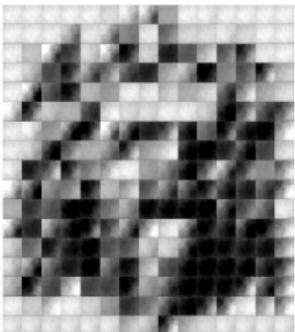


(C)

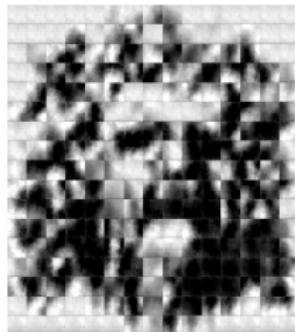
# PCA example: 100% → 5%



(A)



(B)



(C)

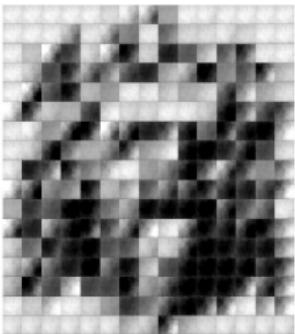


(D)

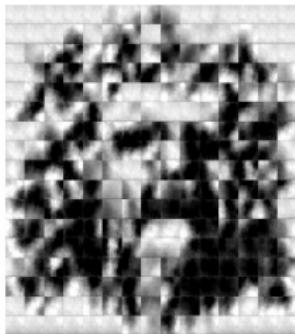
# PCA example: 100% → 10%



(A)



(B)



(C)



(D)

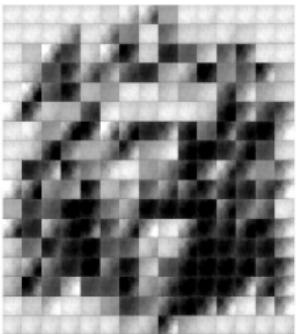


(E)

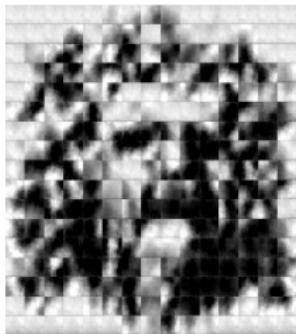
# PCA example: 100% → 20%



(A)



(B)



(C)



(D)



(E)



(F)

# PCA formulation: $d = 1$

- We are looking for the best one-dimensional projection.



- $\mathbb{E}$ := empirical/population expectation:  $\mathbb{E}\mathbf{x} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$ .
- Assumption:  $\mathbb{E}\mathbf{x} = \mathbf{0}$ .

# PCA formulation: $d = 1$

- We are looking for the best one-dimensional projection.



- $\mathbb{E}$ := empirical/population expectation:  $\mathbb{E}\mathbf{x} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$ .
- Assumption:  $\mathbb{E}\mathbf{x} = \mathbf{0}$ .
  - centering:  $\mathbf{x} \rightarrow \mathbf{x} - \mathbb{E}\mathbf{x}$ .

# PCA: projection

Projection ( $\|\mathbf{w}\|_2 = 1$ ):

- $\hat{\mathbf{x}} = \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w}$ .
- zero mean:  $\mathbf{0} \stackrel{?}{=} \mathbb{E}\hat{\mathbf{x}} = \mathbb{E}[\langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w}]$

# PCA: projection

Projection ( $\|\mathbf{w}\|_2 = 1$ ):

- $\hat{\mathbf{x}} = \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w}$ .
- zero mean:  $\mathbf{0} \stackrel{?}{=} \mathbb{E}\hat{\mathbf{x}} = \mathbb{E} [\langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w}] = \langle \mathbf{w}, \underbrace{\mathbb{E}\mathbf{x}}_{=\mathbf{0}} \rangle \mathbf{w}$ .

# PCA: min residual $\Leftrightarrow$ max squared projection

- Goal:  $\mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 \rightarrow \min_{\mathbf{w}}$ .

# PCA: min residual $\Leftrightarrow$ max squared projection

- Goal:  $\mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 \rightarrow \min_{\mathbf{w}}$ .
- Residual  $\Rightarrow$  objective:

$$\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 = \|\mathbf{x} - \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w}\|_2^2$$

# PCA: min residual $\Leftrightarrow$ max squared projection

- Goal:  $\mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 \rightarrow \min_{\mathbf{w}}$ .
- Residual  $\Rightarrow$  objective:

$$\begin{aligned}\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 &= \|\mathbf{x} - \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w}\|_2^2 \\ \|\mathbf{w}\|_2^2 = 1 &\quad \|\mathbf{x}\|_2^2 - \langle \mathbf{w}, \mathbf{x} \rangle^2 \Rightarrow\end{aligned}$$

# PCA: min residual $\Leftrightarrow$ max squared projection

- Goal:  $\mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 \rightarrow \min_{\mathbf{w}}$ .
- Residual  $\Rightarrow$  objective:

$$\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 = \|\mathbf{x} - \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w}\|_2^2$$

$$\|\mathbf{w}\|_2^2 = 1 \quad \|\mathbf{x}\|_2^2 - \langle \mathbf{w}, \mathbf{x} \rangle^2 \Rightarrow$$

$$\mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 = \mathbb{E} \left[ \|\mathbf{x}\|_2^2 - \langle \mathbf{w}, \mathbf{x} \rangle^2 \right] = \underbrace{\mathbb{E} \|\mathbf{x}\|_2^2}_{\text{independent of } \mathbf{w}} - \mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle^2 \Leftrightarrow$$

# PCA: min residual $\Leftrightarrow$ max squared projection

- Goal:  $\mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 \rightarrow \min_{\mathbf{w}}$ .
- Residual  $\Rightarrow$  objective:

$$\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 = \|\mathbf{x} - \langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w}\|_2^2$$

$$\stackrel{\|\mathbf{w}\|_2^2=1}{=} \|\mathbf{x}\|_2^2 - \langle \mathbf{w}, \mathbf{x} \rangle^2 \Rightarrow$$

$$\mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 = \mathbb{E} \left[ \|\mathbf{x}\|_2^2 - \langle \mathbf{w}, \mathbf{x} \rangle^2 \right] = \underbrace{\mathbb{E} \|\mathbf{x}\|_2^2}_{\text{independent of } \mathbf{w}} - \mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle^2 \Leftrightarrow$$

## Solution

maximizes the mean squared projection.

PCA: max squared projection  $\Leftrightarrow$  max variance of projection

By using  $\mathbb{E}y^2 = (\mathbb{E}y)^2 + \text{var}(y)$ :

$$\max_{\mathbf{w}} \leftarrow \mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle^2 = \left( \underbrace{\mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle}_{=0} \right)^2 + \text{var}(\langle \mathbf{w}, \mathbf{x} \rangle).$$

PCA: max squared projection  $\Leftrightarrow$  max variance of projection

By using  $\mathbb{E}y^2 = (\mathbb{E}y)^2 + \text{var}(y)$ :

$$\max_{\mathbf{w}} \leftarrow \mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle^2 = \left( \underbrace{\mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle}_{=0} \right)^2 + \text{var}(\langle \mathbf{w}, \mathbf{x} \rangle).$$

To sum up:

Minimize MSE of the residual :  $\min_{\mathbf{w}} \mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 \Leftrightarrow$

PCA: max squared projection  $\Leftrightarrow$  max variance of projection

By using  $\mathbb{E}y^2 = (\mathbb{E}y)^2 + \text{var}(y)$ :

$$\max_{\mathbf{w}} \leftarrow \mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle^2 = \left( \underbrace{\mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle}_{=0} \right)^2 + \text{var}(\langle \mathbf{w}, \mathbf{x} \rangle).$$

To sum up:

Minimize MSE of the residual :  $\min_{\mathbf{w}} \mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 \Leftrightarrow$

Maximize mean squared projection :  $\max_{\mathbf{w}} \mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle^2 \Leftrightarrow$

PCA: max squared projection  $\Leftrightarrow$  max variance of projection

By using  $\mathbb{E}y^2 = (\mathbb{E}y)^2 + \text{var}(y)$ :

$$\max_{\mathbf{w}} \leftarrow \mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle^2 = \underbrace{\left( \mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle \right)^2}_{=0} + \text{var}(\langle \mathbf{w}, \mathbf{x} \rangle).$$

To sum up:

Minimize MSE of the residual :  $\min_{\mathbf{w}} \mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 \Leftrightarrow$

Maximize mean squared projection :  $\max_{\mathbf{w}} \mathbb{E} \langle \mathbf{w}, \mathbf{x} \rangle^2 \Leftrightarrow$

Maximize variance of the projection :  $\max_{\mathbf{w}} \text{var}(\langle \mathbf{w}, \mathbf{x} \rangle)$ .

# PCA: Optimization

By the bilinearity of cov:

$$\text{var}(\langle \mathbf{w}, \mathbf{x} \rangle) = \text{cov}(\mathbf{w}^T \mathbf{x}, \mathbf{w}^T \mathbf{x})$$

# PCA: Optimization

By the bilinearity of cov:

$$\text{var}(\langle \mathbf{w}, \mathbf{x} \rangle) = \text{cov}(\mathbf{w}^T \mathbf{x}, \mathbf{w}^T \mathbf{x}) = \mathbf{w}^T \text{cov}(\mathbf{x}) \mathbf{w} =: \mathbf{w}^T \boldsymbol{\Sigma} \mathbf{w} \rightarrow \max_{\|\mathbf{w}\|_2=1} .$$

# PCA: Optimization

By the bilinearity of cov:

$$\text{var}(\langle \mathbf{w}, \mathbf{x} \rangle) = \text{cov}(\mathbf{w}^T \mathbf{x}, \mathbf{w}^T \mathbf{x}) = \mathbf{w}^T \text{cov}(\mathbf{x}) \mathbf{w} =: \mathbf{w}^T \boldsymbol{\Sigma} \mathbf{w} \rightarrow \max_{\|\mathbf{w}\|_2=1} .$$

Lagrange function, solving for 'derivatives = 0':

# PCA: Optimization

By the bilinearity of cov:

$$\text{var}(\langle \mathbf{w}, \mathbf{x} \rangle) = \text{cov}(\mathbf{w}^T \mathbf{x}, \mathbf{w}^T \mathbf{x}) = \mathbf{w}^T \text{cov}(\mathbf{x}) \mathbf{w} =: \mathbf{w}^T \Sigma \mathbf{w} \rightarrow \max_{\|\mathbf{w}\|_2=1} .$$

Lagrange function, solving for 'derivatives = 0':

$$L(\mathbf{w}, \lambda) = \underbrace{\mathbf{w}^T \Sigma \mathbf{w}}_{=\text{objective}} - \lambda (\underbrace{\mathbf{w}^T \mathbf{w} - 1}_{=\text{condition}}) \Rightarrow$$

# PCA: Optimization

By the bilinearity of cov:

$$\text{var}(\langle \mathbf{w}, \mathbf{x} \rangle) = \text{cov}(\mathbf{w}^T \mathbf{x}, \mathbf{w}^T \mathbf{x}) = \mathbf{w}^T \text{cov}(\mathbf{x}) \mathbf{w} =: \mathbf{w}^T \boldsymbol{\Sigma} \mathbf{w} \rightarrow \max_{\|\mathbf{w}\|_2=1} .$$

Lagrange function, solving for 'derivatives = 0':

$$L(\mathbf{w}, \lambda) = \underbrace{\mathbf{w}^T \boldsymbol{\Sigma} \mathbf{w}}_{\text{objective}} - \lambda (\underbrace{\mathbf{w}^T \mathbf{w} - 1}_{\text{condition}}) \Rightarrow$$
$$0 = \frac{\partial L(\mathbf{w}, \lambda)}{\partial \lambda} = -(\mathbf{w}^T \mathbf{w} - 1),$$

# PCA: Optimization

By the bilinearity of cov:

$$\text{var}(\langle \mathbf{w}, \mathbf{x} \rangle) = \text{cov}(\mathbf{w}^T \mathbf{x}, \mathbf{w}^T \mathbf{x}) = \mathbf{w}^T \text{cov}(\mathbf{x}) \mathbf{w} =: \mathbf{w}^T \boldsymbol{\Sigma} \mathbf{w} \rightarrow \max_{\|\mathbf{w}\|_2=1} .$$

Lagrange function, solving for 'derivatives = 0':

$$L(\mathbf{w}, \lambda) = \underbrace{\mathbf{w}^T \boldsymbol{\Sigma} \mathbf{w}}_{=\text{objective}} - \lambda (\underbrace{\mathbf{w}^T \mathbf{w} - 1}_{=\text{condition}}) \Rightarrow$$

$$0 = \frac{\partial L(\mathbf{w}, \lambda)}{\partial \lambda} = -(\mathbf{w}^T \mathbf{w} - 1),$$

$$\mathbf{0} = \frac{\partial L(\mathbf{w}, \lambda)}{\partial \mathbf{w}} = 2\boldsymbol{\Sigma} \mathbf{w} - 2\lambda \mathbf{w} \Rightarrow$$

# PCA: Optimization

By the bilinearity of cov:

$$\text{var}(\langle \mathbf{w}, \mathbf{x} \rangle) = \text{cov}(\mathbf{w}^T \mathbf{x}, \mathbf{w}^T \mathbf{x}) = \mathbf{w}^T \text{cov}(\mathbf{x}) \mathbf{w} =: \mathbf{w}^T \boldsymbol{\Sigma} \mathbf{w} \rightarrow \max_{\|\mathbf{w}\|_2=1} .$$

Lagrange function, solving for 'derivatives = 0':

$$L(\mathbf{w}, \lambda) = \underbrace{\mathbf{w}^T \boldsymbol{\Sigma} \mathbf{w}}_{\text{objective}} - \lambda (\underbrace{\mathbf{w}^T \mathbf{w} - 1}_{\text{condition}}) \Rightarrow$$
$$0 = \frac{\partial L(\mathbf{w}, \lambda)}{\partial \lambda} = -(\mathbf{w}^T \mathbf{w} - 1),$$
$$\mathbf{0} = \frac{\partial L(\mathbf{w}, \lambda)}{\partial \mathbf{w}} = 2\boldsymbol{\Sigma} \mathbf{w} - 2\lambda \mathbf{w} \Rightarrow$$

## Solution

$\mathbf{w}^*$ : eigenvector associated to  $\lambda_{\max}(\boldsymbol{\Sigma})$ .

PCA:  $d \geq 1$

# PCA ( $d \geq 1$ ): basis, approximation

- Goal: approximate with a  $d$ -dimensional subspace.
- ONB in the subspace ( $\mathbf{W}^T \mathbf{W} = \mathbf{I}$ ):

$$\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_d] \in \mathbb{R}^{D \times d},$$

- Approximation:

$$\hat{\mathbf{x}} = \sum_{i=1}^d \langle \mathbf{w}_i, \mathbf{x} \rangle \mathbf{w}_i = \mathbf{W} \mathbf{W}^T \mathbf{x}.$$

PCA ( $d \geq 1$ ): min residual  $\Leftrightarrow$  max squared projection

$$\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 = \left\| \mathbf{x} - \mathbf{W}\mathbf{W}^T \mathbf{x} \right\|_2^2$$

# PCA ( $d \geq 1$ ): min residual $\Leftrightarrow$ max squared projection

$$\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 = \left\| \mathbf{x} - \mathbf{W}\mathbf{W}^T \mathbf{x} \right\|_2^2 = \mathbf{x}^T \underbrace{\left( \mathbf{I} - \mathbf{W}\mathbf{W}^T \right) \left( \mathbf{I} - \mathbf{W}\mathbf{W}^T \right)}_{= \mathbf{I} - 2\mathbf{W}\mathbf{W}^T + \mathbf{W}\mathbf{W}^T \mathbf{W}\mathbf{W}^T = \mathbf{I} - \mathbf{W}\mathbf{W}^T} \mathbf{x}$$

Using  $\mathbf{W}^T \mathbf{W} = \mathbf{I}$

$$\begin{aligned}\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 &= \left\| \mathbf{x} - \mathbf{W} \mathbf{W}^T \mathbf{x} \right\|_2^2 = \mathbf{x}^T \underbrace{\left( \mathbf{I} - \mathbf{W} \mathbf{W}^T \right) \left( \mathbf{I} - \mathbf{W} \mathbf{W}^T \right)}_{= \mathbf{I} - 2\mathbf{W} \mathbf{W}^T + \mathbf{W} \mathbf{W}^T \mathbf{W} \mathbf{W}^T = \mathbf{I} - \mathbf{W} \mathbf{W}^T} \mathbf{x} \\ &= \|\mathbf{x}\|_2^2 - \left\| \mathbf{W}^T \mathbf{x} \right\|_2^2,\end{aligned}$$

Using  $\mathbf{W}^T \mathbf{W} = \mathbf{I}$

$$\begin{aligned}\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 &= \left\| \mathbf{x} - \mathbf{W} \mathbf{W}^T \mathbf{x} \right\|_2^2 = \mathbf{x}^T \underbrace{\left( \mathbf{I} - \mathbf{W} \mathbf{W}^T \right) \left( \mathbf{I} - \mathbf{W} \mathbf{W}^T \right)}_{= \mathbf{I} - 2\mathbf{W} \mathbf{W}^T + \mathbf{W} \mathbf{W}^T \mathbf{W} \mathbf{W}^T = \mathbf{I} - \mathbf{W} \mathbf{W}^T} \mathbf{x} \\ &= \|\mathbf{x}\|_2^2 - \left\| \mathbf{W}^T \mathbf{x} \right\|_2^2,\end{aligned}$$

$$\mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 = \underbrace{\mathbb{E} \|\mathbf{x}\|_2^2}_{\text{independent of } \mathbf{W}} - \mathbb{E} \left\| \mathbf{W}^T \mathbf{x} \right\|_2^2.$$

# PCA ( $d \geq 1$ ): min residual $\Leftrightarrow$ max squared projection

Using  $\mathbf{W}^T \mathbf{W} = \mathbf{I}$

$$\begin{aligned}\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 &= \left\| \mathbf{x} - \mathbf{W} \mathbf{W}^T \mathbf{x} \right\|_2^2 = \mathbf{x}^T \underbrace{\left( \mathbf{I} - \mathbf{W} \mathbf{W}^T \right) \left( \mathbf{I} - \mathbf{W} \mathbf{W}^T \right)}_{= \mathbf{I} - 2\mathbf{W} \mathbf{W}^T + \mathbf{W} \mathbf{W}^T \mathbf{W} \mathbf{W}^T = \mathbf{I} - \mathbf{W} \mathbf{W}^T} \mathbf{x} \\ &= \|\mathbf{x}\|_2^2 - \left\| \mathbf{W}^T \mathbf{x} \right\|_2^2,\end{aligned}$$

$$\mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 = \underbrace{\mathbb{E} \|\mathbf{x}\|_2^2}_{\text{independent of } \mathbf{W}} - \mathbb{E} \left\| \mathbf{W}^T \mathbf{x} \right\|_2^2.$$

Thus  $\min_w \mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 \Leftrightarrow \max_w \mathbb{E} \left\| \mathbf{W}^T \mathbf{x} \right\|_2^2$ .

PCA ( $d \geq 1$ ): max squared projection  $\Leftrightarrow$  max variance of projection

Let  $\mathbf{y} = \mathbf{W}^T \mathbf{x}$ :

$$\mathbb{E} \|\mathbf{y}\|_2^2 - \|\mathbb{E} \mathbf{y}\|_2^2 = \text{var}(\mathbf{y})?$$

PCA ( $d \geq 1$ ): max squared projection  $\Leftrightarrow$  max variance of projection

Let  $\mathbf{y} = \mathbf{W}^T \mathbf{x}$ :

$$\mathbb{E} \|\mathbf{y}\|_2^2 - \|\mathbb{E} \mathbf{y}\|_2^2 = \text{var}(\mathbf{y})?$$

$$= \mathbb{E} \left[ \sum_i y_i^2 \right] - \sum_i (\mathbb{E} y_i)^2 = \sum_i \text{var}(y_i) \Rightarrow$$

PCA ( $d \geq 1$ ): max squared projection  $\Leftrightarrow$  max variance of projection

Let  $\mathbf{y} = \mathbf{W}^T \mathbf{x}$ :

$$\mathbb{E} \|\mathbf{y}\|_2^2 - \|\mathbb{E}\mathbf{y}\|_2^2 = \text{var}(\mathbf{y})?$$

$$= \mathbb{E} \left[ \sum_i y_i^2 \right] - \sum_i (\mathbb{E} y_i)^2 = \sum_i \text{var}(y_i) \Rightarrow$$

$$\begin{aligned} \mathbb{E} \left\| \mathbf{W}^T \mathbf{x} \right\|_2^2 - \left\| \underbrace{\mathbb{E}[\mathbf{W}^T \mathbf{x}]}_{=\mathbf{W}^T \mathbb{E}\mathbf{x}=\mathbf{0}} \right\|_2^2 &= \sum_i \text{var} \left( \left( \mathbf{W}^T \mathbf{x} \right)_i \right) \rightarrow \max_{\mathbf{W}}. \end{aligned}$$

- The  $d$  principal components:

$$\{\mathbf{w}_i\}_{i=1}^d = \text{top } d \text{ eigenvectors of } \boldsymbol{\Sigma} = \text{cov}(\mathbf{x}).$$

- The  $d$  principal components:

$\{\mathbf{w}_i\}_{i=1}^d = \text{top } d \text{ eigenvectors of } \boldsymbol{\Sigma} = \text{cov}(\mathbf{x}).$

- $\boldsymbol{\Sigma}$ : symmetric, positive semi-definite  $\Rightarrow \{\mathbf{w}_i\}$ : ONS,  $\lambda_i \geq 0$ .

- The  $d$  principal components:

$$\{\mathbf{w}_i\}_{i=1}^d = \text{top } d \text{ eigenvectors of } \boldsymbol{\Sigma} = \text{cov}(\mathbf{x}).$$

- $\boldsymbol{\Sigma}$ : symmetric, positive semi-definite  $\Rightarrow \{\mathbf{w}_i\}$ : ONS,  $\lambda_i \geq 0$ .
- Variance decomposition:  $\text{cov}(\mathbf{x}) = \sum_{i=1}^D \lambda_i \mathbf{w}_i \mathbf{w}_i^T$ .

# PCA: $d \geq 1$

- The  $d$  principal components:

$$\{\mathbf{w}_i\}_{i=1}^d = \text{top } d \text{ eigenvectors of } \boldsymbol{\Sigma} = \text{cov}(\mathbf{x}).$$

- $\boldsymbol{\Sigma}$ : symmetric, positive semi-definite  $\Rightarrow \{\mathbf{w}_i\}$ : ONS,  $\lambda_i \geq 0$ .
- Variance decomposition:  $\text{cov}(\mathbf{x}) = \sum_{i=1}^D \lambda_i \mathbf{w}_i \mathbf{w}_i^T$ .
- Energy preserved using  $d$  components:  $\sum_{i=1}^d \lambda_i \Rightarrow$

$$R^2 = R^2(d) := \frac{\sum_{i=1}^d \lambda_i}{\sum_{i=1}^D \lambda_i} \in [0, 1].$$

# PCA: $d \geq 1$

- The  $d$  principal components:

$$\{\mathbf{w}_i\}_{i=1}^d = \text{top } d \text{ eigenvectors of } \boldsymbol{\Sigma} = \text{cov}(\mathbf{x}).$$

- $\boldsymbol{\Sigma}$ : symmetric, positive semi-definite  $\Rightarrow \{\mathbf{w}_i\}$ : ONS,  $\lambda_i \geq 0$ .
- Variance decomposition:  $\text{cov}(\mathbf{x}) = \sum_{i=1}^D \lambda_i \mathbf{w}_i \mathbf{w}_i^T$ .
- Energy preserved using  $d$  components:  $\sum_{i=1}^d \lambda_i \Rightarrow$

$$R^2 = R^2(d) := \frac{\sum_{i=1}^d \lambda_i}{\sum_{i=1}^D \lambda_i} \in [0, 1].$$

- In practice: choose  $d$  such that  $R^2 \approx 0.8 - 0.9$ .

# Kernel PCA: idea for ' $d = 1$ ' $\leftrightarrow f$

Let  $\mathcal{H} = \mathcal{H}_k$ .

- Objective function:

$$J(f) = \frac{1}{n} \sum_{i=1}^n \left\langle f, \underbrace{\varphi(x_i) - \frac{1}{n} \sum_{j=1}^n \varphi(x_j)}_{=: \tilde{\varphi}(x_i)} \right\rangle^2 = \text{var}(f) \rightarrow \max_{f: \|f\|_{\mathcal{H}} \leq 1} .$$

# Kernel PCA: idea for ' $d = 1$ ' $\leftrightarrow f$

Let  $\mathcal{H} = \mathcal{H}_k$ .

- Objective function:

$$J(f) = \frac{1}{n} \sum_{i=1}^n \left\langle f, \underbrace{\varphi(x_i) - \frac{1}{n} \sum_{j=1}^n \varphi(x_j)}_{=: \tilde{\varphi}(x_i)} \right\rangle^2 = \text{var}(f) \rightarrow \max_{f: \|f\|_{\mathcal{H}} \leq 1} .$$

- The solution can be searched in the form ( $\mathcal{H} \ni f \leftrightarrow \mathbf{a} \in \mathbb{R}^n$ ):

$$f = \sum_{i=1}^n a_i \tilde{\varphi}(x_i)$$

since component  $\perp \text{span}(\{\tilde{\varphi}(x_i)\}_{i=1}^n)$  has no contribution.

# Kernel PCA: idea for ' $d = 1$ ' $\leftrightarrow f$

Let  $\mathcal{H} = \mathcal{H}_k$ .

- Objective function:

$$J(f) = \frac{1}{n} \sum_{i=1}^n \left\langle f, \underbrace{\varphi(x_i) - \frac{1}{n} \sum_{j=1}^n \varphi(x_j)}_{=: \tilde{\varphi}(x_i)} \right\rangle^2 = \text{var}(f) \rightarrow \max_{f: \|f\|_{\mathcal{H}} \leq 1} .$$

- The solution can be searched in the form ( $\mathcal{H} \ni f \leftrightarrow \mathbf{a} \in \mathbb{R}^n$ ):

$$f = \sum_{i=1}^n a_i \tilde{\varphi}(x_i)$$

since component  $\perp \text{span}(\{\tilde{\varphi}(x_i)\}_{i=1}^n)$  has no contribution.

- We will get an eigenvalue problem for  $\mathbf{a}$ .

## (Empirical) covariance operator

$$C := \frac{1}{n} \sum_{i=1}^n \tilde{\varphi}(x_i) \otimes \tilde{\varphi}(x_i).$$

$c \otimes d$  is the analogue of  $cd^T$ :

$$(c \otimes d)(e) = c \langle d, e \rangle_{\mathcal{H}}.$$

## (Empirical) covariance operator

$$C := \frac{1}{n} \sum_{i=1}^n \tilde{\varphi}(x_i) \otimes \tilde{\varphi}(x_i).$$

$c \otimes d$  is the analogue of  $cd^T$ :

$$(c \otimes d)(e) = c \langle d, e \rangle_{\mathcal{H}}.$$

Similarly to the finite-dimensional case:

$$Cf_j = \lambda_j f_j.$$

### Challenge

How do we solve this eigenvalue problem?

# Computation of $Cf_j$

Assume  $j$  is fixed ( $Cf = \lambda f$ ):

$$Cf = \left[ \frac{1}{n} \sum_{i=1}^n \tilde{\varphi}(x_i) \otimes \tilde{\varphi}(x_i) \right] \textcolor{blue}{f}$$

# Computation of $Cf_j$

Assume  $j$  is fixed ( $Cf = \lambda f$ ):

$$\begin{aligned} Cf &= \left[ \frac{1}{n} \sum_{i=1}^n \tilde{\varphi}(x_i) \otimes \tilde{\varphi}(x_i) \right] \textcolor{blue}{f} \\ &\stackrel{\otimes \text{ def}}{=} \frac{1}{n} \sum_{i=1}^n \tilde{\varphi}(x_i) \left\langle \tilde{\varphi}(x_i), \sum_{j=1}^n \textcolor{blue}{a}_j \tilde{\varphi}(x_j) \right\rangle_{\mathcal{H}} \end{aligned}$$

# Computation of $Cf_j$

Assume  $j$  is fixed ( $Cf = \lambda f$ ):

$$Cf = \left[ \frac{1}{n} \sum_{i=1}^n \tilde{\varphi}(x_i) \otimes \tilde{\varphi}(x_i) \right] \textcolor{blue}{f}$$
$$\stackrel{\otimes \text{def}}{=} \frac{1}{n} \sum_{i=1}^n \tilde{\varphi}(x_i) \left\langle \tilde{\varphi}(x_i), \sum_{j=1}^n \textcolor{blue}{a}_j \tilde{\varphi}(x_j) \right\rangle_{\mathcal{H}} = \frac{1}{n} \sum_{i=1}^n \tilde{\varphi}(x_i) \sum_{j=1}^n a_j \tilde{k}(x_i, x_j)$$

with  $\tilde{\mathbf{G}} = \mathbf{HGH} = \left[ \tilde{k}(x_i, x_j) \right]_{i,j=1}^n$ ,  $\mathbf{H} = \mathbf{I}_n - \frac{\mathbf{E}_n}{n}$ .

# Eigenvalue problem

- We want to solve  $Cf = \lambda f$ ,  $\textcolor{red}{C}\textcolor{blue}{f} = \frac{1}{n} \sum_{i=1}^n \tilde{\varphi}(x_i) \sum_{j=1}^n a_j \tilde{k}(x_i, x_j)$ .
- Idea: multiple by  $\tilde{\varphi}(x_r)$

$$\langle \tilde{\varphi}(x_r), \lambda \textcolor{blue}{f} \rangle_{\mathcal{H}} = \left\langle \tilde{\varphi}(x_r), \lambda \sum_{j=1}^n a_j \tilde{\varphi}(x_j) \right\rangle_{\mathcal{H}}$$

# Eigenvalue problem

- We want to solve  $Cf = \lambda f$ ,  $\textcolor{red}{C}\textcolor{blue}{f} = \frac{1}{n} \sum_{i=1}^n \tilde{\varphi}(x_i) \sum_{j=1}^n a_j \tilde{k}(x_i, x_j)$ .
- Idea: multiple by  $\tilde{\varphi}(x_r)$

$$\langle \tilde{\varphi}(x_r), \lambda \textcolor{blue}{f} \rangle_{\mathcal{H}} = \left\langle \tilde{\varphi}(x_r), \lambda \sum_{j=1}^n a_j \tilde{\varphi}(x_j) \right\rangle_{\mathcal{H}} = \lambda \underbrace{\sum_{j=1}^n a_j \tilde{G}_{rj}}_{(\tilde{\mathbf{G}}\mathbf{a})_r},$$

# Eigenvalue problem

- We want to solve  $Cf = \lambda f$ ,  $\textcolor{red}{Cf} = \frac{1}{n} \sum_{i=1}^n \tilde{\varphi}(x_i) \sum_{j=1}^n a_j \tilde{k}(x_i, x_j)$ .
- Idea: multiple by  $\tilde{\varphi}(x_r)$

$$\langle \tilde{\varphi}(x_r), \lambda \textcolor{blue}{f} \rangle_{\mathcal{H}} = \left\langle \tilde{\varphi}(x_r), \lambda \sum_{j=1}^n a_j \tilde{\varphi}(x_j) \right\rangle_{\mathcal{H}} = \lambda \underbrace{\sum_{j=1}^n a_j \tilde{G}_{rj}}_{(\tilde{\mathbf{G}}\mathbf{a})_r},$$

$$\langle \tilde{\varphi}(x_r), \textcolor{red}{Cf} \rangle_{\mathcal{H}} = \left\langle \tilde{\varphi}(x_r), \frac{1}{n} \sum_{i=1}^n \tilde{\varphi}(x_i) \sum_{j=1}^n a_j \tilde{k}(x_i, x_j) \right\rangle_{\mathcal{H}}$$

# Eigenvalue problem

- We want to solve  $Cf = \lambda f$ ,  $\textcolor{red}{Cf} = \frac{1}{n} \sum_{i=1}^n \tilde{\varphi}(x_i) \sum_{j=1}^n a_j \tilde{k}(x_i, x_j)$ .
- Idea: multiple by  $\tilde{\varphi}(x_r)$

$$\langle \tilde{\varphi}(x_r), \lambda \textcolor{blue}{f} \rangle_{\mathcal{H}} = \left\langle \tilde{\varphi}(x_r), \lambda \sum_{j=1}^n a_j \tilde{\varphi}(x_j) \right\rangle_{\mathcal{H}} = \lambda \underbrace{\sum_{j=1}^n a_j \tilde{G}_{rj}}_{(\tilde{\mathbf{G}}\mathbf{a})_r},$$

$$\begin{aligned}\langle \tilde{\varphi}(x_r), \textcolor{red}{Cf} \rangle_{\mathcal{H}} &= \left\langle \tilde{\varphi}(x_r), \frac{1}{n} \sum_{i=1}^n \tilde{\varphi}(x_i) \sum_{j=1}^n a_j \tilde{k}(x_i, x_j) \right\rangle_{\mathcal{H}} \\ &= \frac{1}{n} \sum_{i=1}^n \tilde{G}_{ri} \underbrace{\sum_{j=1}^n a_j \tilde{G}_{ij}}_{(\tilde{\mathbf{G}}\mathbf{a})_i}\end{aligned}$$

# Eigenvalue problem

- We want to solve  $Cf = \lambda f$ ,  $\textcolor{red}{Cf} = \frac{1}{n} \sum_{i=1}^n \tilde{\varphi}(x_i) \sum_{j=1}^n a_j \tilde{k}(x_i, x_j)$ .
- Idea: multiple by  $\tilde{\varphi}(x_r)$

$$\langle \tilde{\varphi}(x_r), \lambda \textcolor{blue}{f} \rangle_{\mathcal{H}} = \left\langle \tilde{\varphi}(x_r), \lambda \sum_{j=1}^n a_j \tilde{\varphi}(x_j) \right\rangle_{\mathcal{H}} = \lambda \underbrace{\sum_{j=1}^n a_j \tilde{G}_{rj}}_{(\tilde{\mathbf{G}}\mathbf{a})_r},$$

$$\begin{aligned}\langle \tilde{\varphi}(x_r), \textcolor{red}{Cf} \rangle_{\mathcal{H}} &= \left\langle \tilde{\varphi}(x_r), \frac{1}{n} \sum_{i=1}^n \tilde{\varphi}(x_i) \sum_{j=1}^n a_j \tilde{k}(x_i, x_j) \right\rangle_{\mathcal{H}} \\ &= \frac{1}{n} \sum_{i=1}^n \tilde{G}_{ri} \underbrace{\sum_{j=1}^n a_j \tilde{G}_{ij}}_{(\tilde{\mathbf{G}}\mathbf{a})_i} = \frac{1}{n} (\tilde{\mathbf{G}}^2 \mathbf{a})_r.\end{aligned}$$

# Eigenvalue problem

- We want to solve  $Cf = \lambda f$ ,  $\textcolor{red}{Cf} = \frac{1}{n} \sum_{i=1}^n \tilde{\varphi}(x_i) \sum_{j=1}^n a_j \tilde{k}(x_i, x_j)$ .
- Idea: multiple by  $\tilde{\varphi}(x_r)$

$$\langle \tilde{\varphi}(x_r), \lambda \textcolor{blue}{f} \rangle_{\mathcal{H}} = \left\langle \tilde{\varphi}(x_r), \lambda \sum_{j=1}^n a_j \tilde{\varphi}(x_j) \right\rangle_{\mathcal{H}} = \lambda \underbrace{\sum_{j=1}^n a_j \tilde{G}_{rj}}_{(\tilde{\mathbf{G}}\mathbf{a})_r},$$

$$\begin{aligned}\langle \tilde{\varphi}(x_r), \textcolor{red}{Cf} \rangle_{\mathcal{H}} &= \left\langle \tilde{\varphi}(x_r), \frac{1}{n} \sum_{i=1}^n \tilde{\varphi}(x_i) \sum_{j=1}^n a_j \tilde{k}(x_i, x_j) \right\rangle_{\mathcal{H}} \\ &= \frac{1}{n} \sum_{i=1}^n \tilde{G}_{ri} \underbrace{\sum_{j=1}^n a_j \tilde{G}_{ij}}_{(\tilde{\mathbf{G}}\mathbf{a})_i} = \frac{1}{n} (\tilde{\mathbf{G}}^2 \mathbf{a})_r.\end{aligned}$$

- Eigenvalue problem:  $\tilde{\mathbf{G}}^2 \mathbf{a} = n\lambda \tilde{\mathbf{G}}\mathbf{a}$ , i.e.  $\tilde{\mathbf{G}}\mathbf{a} = (n\lambda)\mathbf{a}$ .

# Orthogonal eigenvectors in kernel PCA

Taking two (eigenvector, eigenvalue) pairs:

$$\mathbf{f}_1 = \sum_{i=1}^n a_{1i} \tilde{\varphi}(x_i), \quad \tilde{\mathbf{G}}\mathbf{a}_1 = \lambda_1 \mathbf{a}_1,$$

$$\mathbf{f}_2 = \sum_{j=1}^n a_{2j} \tilde{\varphi}(x_j), \quad \tilde{\mathbf{G}}\mathbf{a}_2 = \lambda_2 \mathbf{a}_2.$$

one has

$$0 \stackrel{?}{=} \langle f_1, f_2 \rangle_{\mathcal{H}}$$

# Orthogonal eigenvectors in kernel PCA

Taking two (eigenvector, eigenvalue) pairs:

$$f_1 = \sum_{i=1}^n a_{1i} \tilde{\varphi}(x_i), \quad \tilde{\mathbf{G}}\mathbf{a}_1 = \lambda_1 \mathbf{a}_1,$$

$$f_2 = \sum_{j=1}^n a_{2j} \tilde{\varphi}(x_j), \quad \tilde{\mathbf{G}}\mathbf{a}_2 = \lambda_2 \mathbf{a}_2.$$

one has

$$0 \stackrel{?}{=} \langle f_1, f_2 \rangle_{\mathcal{H}} = \left\langle \sum_{i=1}^n a_{1i} \tilde{\varphi}(x_i), \sum_{j=1}^n a_{2j} \tilde{\varphi}(x_j) \right\rangle_{\mathcal{H}}$$

# Orthogonal eigenvectors in kernel PCA

Taking two (eigenvector, eigenvalue) pairs:

$$f_1 = \sum_{i=1}^n a_{1i} \tilde{\varphi}(x_i), \quad \tilde{\mathbf{G}}\mathbf{a}_1 = \lambda_1 \mathbf{a}_1,$$

$$f_2 = \sum_{j=1}^n a_{2j} \tilde{\varphi}(x_j), \quad \tilde{\mathbf{G}}\mathbf{a}_2 = \lambda_2 \mathbf{a}_2.$$

one has

$$0 \stackrel{?}{=} \langle f_1, f_2 \rangle_{\mathcal{H}} = \left\langle \sum_{i=1}^n a_{1i} \tilde{\varphi}(x_i), \sum_{j=1}^n a_{2j} \tilde{\varphi}(x_j) \right\rangle_{\mathcal{H}} = \mathbf{a}_1^T \tilde{\mathbf{G}} \mathbf{a}_2$$

# Orthogonal eigenvectors in kernel PCA

Taking two (eigenvector, eigenvalue) pairs:

$$f_1 = \sum_{i=1}^n a_{1i} \tilde{\varphi}(x_i), \quad \tilde{\mathbf{G}}\mathbf{a}_1 = \lambda_1 \mathbf{a}_1,$$

$$f_2 = \sum_{j=1}^n a_{2j} \tilde{\varphi}(x_j), \quad \tilde{\mathbf{G}}\mathbf{a}_2 = \lambda_2 \mathbf{a}_2.$$

one has

$$0 \stackrel{?}{=} \langle f_1, f_2 \rangle_{\mathcal{H}} = \left\langle \sum_{i=1}^n a_{1i} \tilde{\varphi}(x_i), \sum_{j=1}^n a_{2j} \tilde{\varphi}(x_j) \right\rangle_{\mathcal{H}} = \mathbf{a}_1^T \tilde{\mathbf{G}} \mathbf{a}_2 = \mathbf{a}_1^T \lambda_2 \mathbf{a}_2.$$

# Orthogonality $\Rightarrow$ projection is easy

Projection of a new  $x^*$  to the first  $d$ -PCs:

$$\Pi [\tilde{\varphi}(x^*)] = \sum_{j=1}^d \langle \tilde{\varphi}(x^*), f_j \rangle_{\mathcal{H}} f_j.$$

# Orthogonality $\Rightarrow$ projection is easy

Projection of a new  $x^*$  to the first  $d$ -PCs:

$$\Pi [\tilde{\varphi}(x^*)] = \sum_{j=1}^d \langle \tilde{\varphi}(x^*), f_j \rangle_{\mathcal{H}} f_j.$$

For fixed  $f = f_j$ , using  $f = \sum_{i=1}^n a_i \tilde{\varphi}(x_i)$ :

$$\langle \tilde{\varphi}(x^*), f \rangle_{\mathcal{H}} f = \sum_i a_i \tilde{k}(x_i, x^*) f = \sum_{i,j=1}^n a_i a_j \tilde{k}(x_i, x^*) \tilde{\varphi}(x_j).$$

# In denoising application: PCA vs kernel PCA

The pre-image problem to solve:  $\widehat{x^*} = \arg \min_{x \in \mathcal{X}} \|\tilde{\varphi}(x) - \Pi[\tilde{\varphi}(x^*)]\|_{\mathcal{H}}^2$ .

		Gaussian noise									
orig.	noisy	0	1	2	3	4	5	6	7	8	9
n = 1	0	1	2	3	4	5	6	7	8	9	
4	0	1	2	3	4	5	6	7	8	9	
16	0	1	2	3	4	5	6	7	8	9	
64	0	1	2	3	4	5	6	7	8	9	
256	0	1	2	3	4	5	6	7	8	9	
n = 1	0	1	2	3	4	5	6	7	8	9	
4	0	1	2	3	4	5	6	7	8	9	
16	0	1	2	3	4	5	6	7	8	9	
64	0	1	2	3	4	5	6	7	8	9	
256	0	1	2	3	4	5	6	7	8	9	

# Kernel-based Divergence & Independence Measures

- Mean embedding:

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

# KL Divergence and Mutual Information Alternatives

- Mean embedding:

$$\mu_k(\mathbb{P}) = \int_{\mathcal{X}} \underbrace{\varphi(x)}_{k(\cdot, x)} d\mathbb{P}(x) \in \mathcal{H}_k = \overline{\text{span}}(k(\cdot, x) : x \in \mathcal{X}).$$



- Maximum mean discrepancy:

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

# KL Divergence and Mutual Information Alternatives

- Mean embedding:

$$\mu_k(\mathbb{P}) = \int_{\mathcal{X}} \underbrace{\varphi(x)}_{k(\cdot, x)} d\mathbb{P}(x) \in \mathcal{H}_k = \overline{\text{span}}(k(\cdot, x) : x \in \mathcal{X}).$$



- Maximum mean discrepancy:

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

- Hilbert-Schmidt independence criterion,  $k = k_1 \otimes k_2$ :

$$\begin{aligned}\text{HSIC}_k(\mathbb{P}) &= \text{MMD}_k(\mathbb{P}, \mathbb{P}_1 \otimes \mathbb{P}_2), \\ (k_1 \otimes k_2)((x, y), (x', y')) &= k_1(x, x')k_2(y, y').\end{aligned}$$

# KL Divergence and Mutual Information Alternatives

- Mean embedding:

$$\mu_k(\mathbb{P}) = \int_{\mathcal{X}} \underbrace{\varphi(x)}_{k(\cdot, x)} d\mathbb{P}(x) \in \mathcal{H}_k = \overline{\text{span}}(k(\cdot, x) : x \in \mathcal{X}).$$



- Maximum mean discrepancy:

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

- Hilbert-Schmidt independence criterion,  $k = k_1 \otimes k_2$ :

$$\begin{aligned}\text{HSIC}_k(\mathbb{P}) &= \text{MMD}_k(\mathbb{P}, \mathbb{P}_1 \otimes \mathbb{P}_2), \\ (k_1 \otimes k_2)((x, y), (x', y')) &= k_1(x, x')k_2(y, y').\end{aligned}$$

- Kernel Canonical Correlation Analysis:

$$\text{KCCA}(\mathbb{P}_{xy}) = \sup_{f \in \mathcal{H}_k, g \in \mathcal{H}_\ell} \text{corr}(f(x), g(y)).$$

# Independence measures – History of KCCA

- 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 – History of KCCA

- 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:

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

satisfies 1-4.

- He showed:

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

satisfies 1-4.

- Too ambitious:

- computationally intractable.
- many functions.

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

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

- Independence measure,
- distance,
- inner product

measures/estimates on probability distributions

without density estimation!

# Kernel Canonical Correlation Analysis (KCCA)

# 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  $\mathcal{H}_k$ ,  $\mathcal{H}_\ell$ .

# 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  $\mathcal{H}_k$ ,  $\mathcal{H}_\ell$ .
- 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**.
- $k, \ell$  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**.
- $k, \ell$  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

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

# KCCA: empirical estimate

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

Similarly:

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

# 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:  $g$  as  $\langle g, \tilde{\psi}(y_n) \rangle_{\mathcal{H}_\ell}$ ].  $\Rightarrow$

## KCCA: empirical estimate

- $f$ : appears only as  $\langle f, \tilde{\varphi}(x_n) \rangle_{\mathcal{H}_k}$  [similarly:  $g$  as  $\langle g, \tilde{\psi}(y_n) \rangle_{\mathcal{H}_\ell}$ ].  $\Rightarrow$
- $\forall$  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:  $g$  as  $\langle g, \tilde{\psi}(y_n) \rangle_{\mathcal{H}_\ell}$ ].  $\Rightarrow$
- $\forall$  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  $a, b \in \mathbb{R} \setminus \{0\}$ .
- 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.

## KCCA: symmetry, other form

For a

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

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

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

## KCCA: symmetry, other form

For a

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

$([\mathbf{c}, \mathbf{d}], \lambda)$  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

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

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

# 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

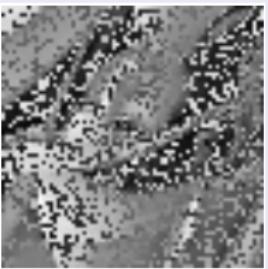
In short

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

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

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

Recall: outlier-robust image registration (it was KCCA)



KCCA: finished.

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

# Mean embedding: 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  $k$ ).

# Mean embedding: 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  $k$ ).

Intuition of MMD and HSIC estimation follows

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

$$\text{HSIC}_k(\mathbb{P}) = \text{MMD}_k(\mathbb{P}, \mathbb{P}_1 \otimes \mathbb{P}_2).$$

# Maximum Mean Discrepancy (MMD)

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

---

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

---

# MMD estimator: intuition

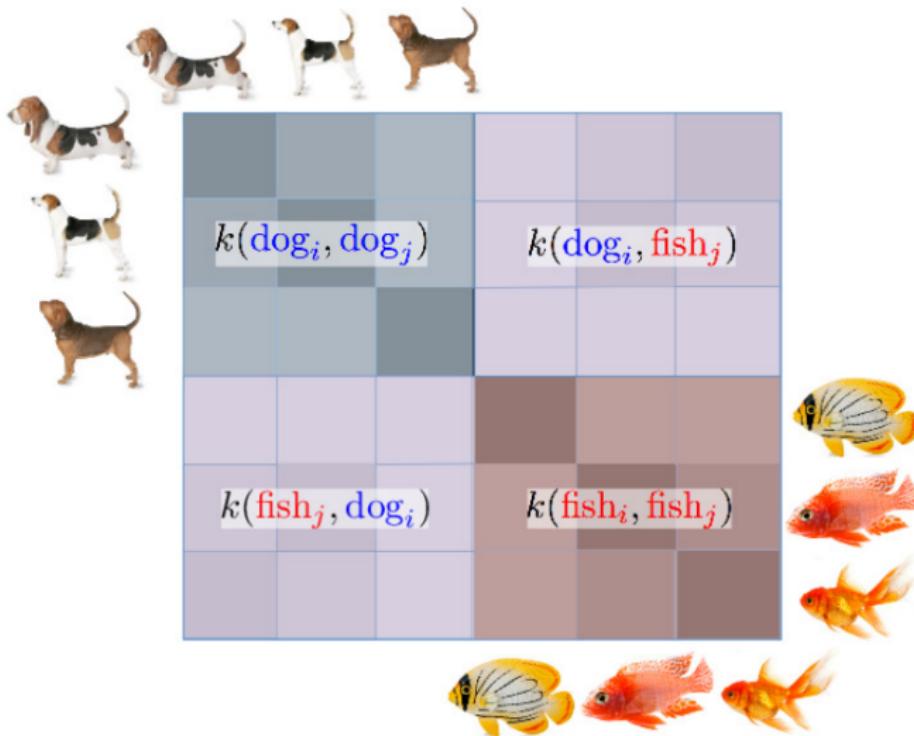


$\sim P$

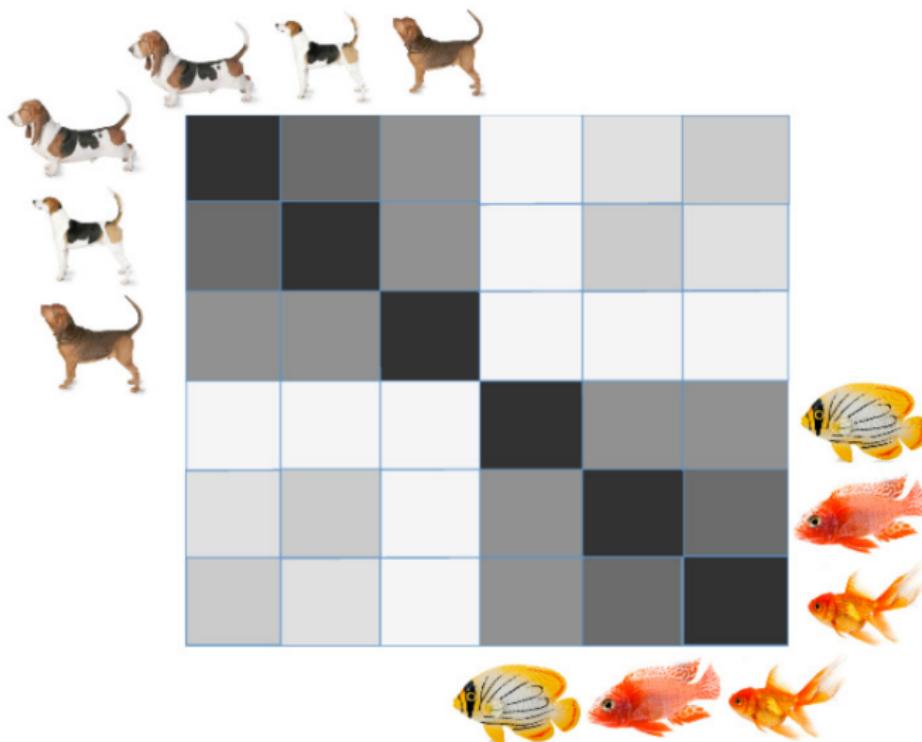


$\sim Q$

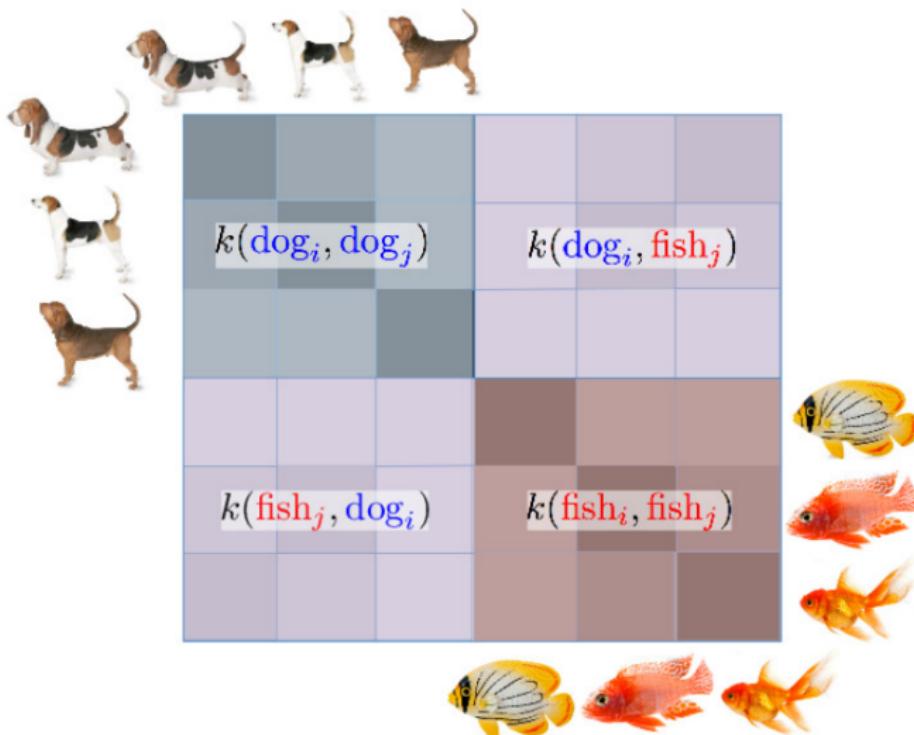
# MMD estimator: intuition



# MMD estimator: intuition



# MMD estimator: intuition



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

<sup>†</sup>  $\widehat{\text{MMD}}$  &  $\widehat{\text{HSIC}}$  illustration credit: Arthur Gretton

- Feature of a distribution:  $\mu_{\mathbb{P}} = \mathbb{E}_{x \sim \mathbb{P}} \varphi(x)$ .

- Feature of a distribution:  $\mu_{\mathbb{P}} = \mathbb{E}_{x \sim \mathbb{P}} \varphi(x)$ .
- MMD = difference between feature means:

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

- Feature of a distribution:  $\mu_{\mathbb{P}} = \mathbb{E}_{x \sim \mathbb{P}} \varphi(x)$ .
- MMD = difference between feature means:

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

- Feature of a distribution:  $\mu_{\mathbb{P}} = \mathbb{E}_{x \sim \mathbb{P}} \varphi(x)$ .
- MMD = difference between feature means:

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

## MMD estimator: mean of kernel values

- Feature of a distribution:  $\mu_{\mathbb{P}} = \mathbb{E}_{x \sim \mathbb{P}} \varphi(x)$ .
- MMD = difference between feature means:

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

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

using  $\{x_i\}_{i=1}^m \sim \mathbb{P}$ ,  $\{y_j\}_{j=1}^n \sim \mathbb{Q}$  samples.

- Computational complexity:  $\mathcal{O}((m+n)^2)$ , quadratic.

# Hilbert-Schmidt Independence Criterion (HSIC)

# 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.



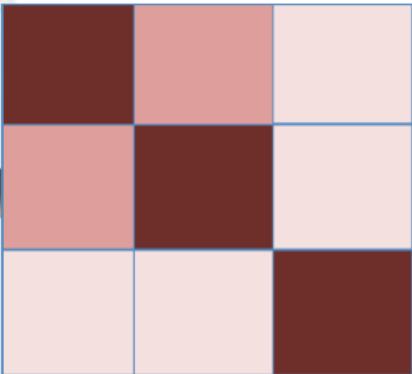
A large animal who slings slobber, exudes a distinctive houndy odor, and wants nothing more than to follow his nose. They need a significant amount of exercise and mental stimulation.



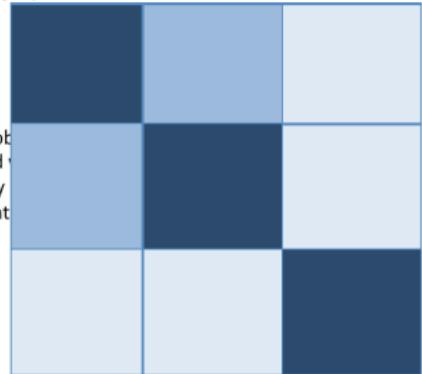
Known for their curiosity, intelligence, and excellent communication skills, the Javanese breed is perfect if you want a responsive, interactive pet, one that will blow in your ear and follow you everywhere.

Text from [dogtime.com](http://dogtime.com) and [petfinder.com](http://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

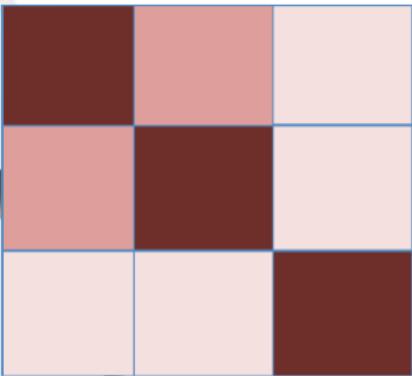


Known for their curiosity, intelligence, and excellent communication skills, the Javanese breed is perfect if you want a responsive, interactive pet, one that will blow in your ear and follow you everywhere.

# HSIC intuition: Gram matrices

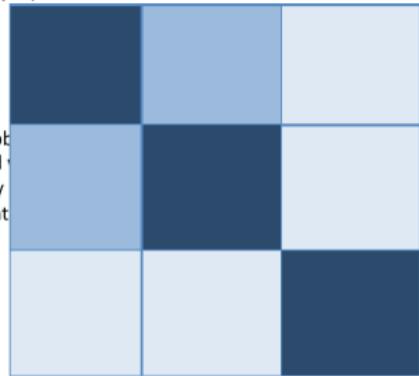


$\tilde{\mathbf{G}}_x$



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.

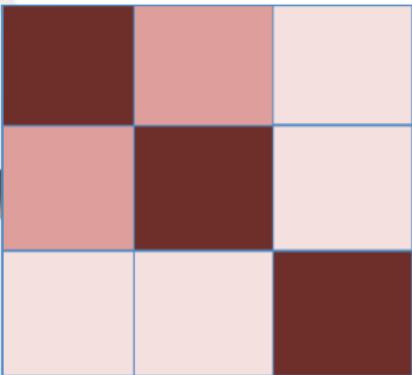
Empirical estimate:

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

# HSIC intuition: Gram matrices

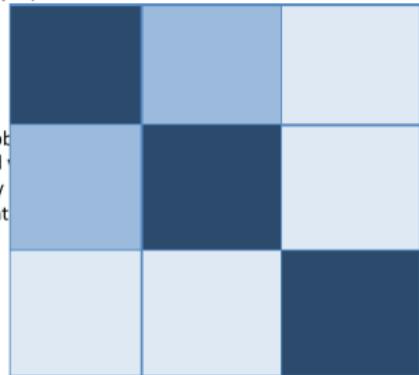


$\tilde{\mathbf{G}}_x$



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:

$$\widehat{\text{HSIC}}^2 = \frac{1}{n^2} \left\langle \tilde{\mathbf{G}}_x, \tilde{\mathbf{G}}_y \right\rangle_F. \quad \text{HSIC}(\mathbb{P}_{xy}) = \text{MMD}(\mathbb{P}_{xy}, \mathbb{P}_x \otimes \mathbb{P}_y).$$

# Idea of the HSIC estimator

MMD in terms of kernel evaluations:

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

# Idea of the HSIC estimator

MMD in terms of kernel evaluations:

$$\text{MMD}^2(\mathbb{P}, \mathbb{Q}) = \mathbb{E}_{x \sim \mathbb{P}, x' \sim \mathbb{P}} k(x, x') + \mathbb{E}_{y \sim \mathbb{Q}, y' \sim \mathbb{Q}} k(y, y') \\ - 2\mathbb{E}_{x \sim \mathbb{P}, y \sim \mathbb{Q}} k(x, y).$$

## Question

Can we rewrite HSIC in terms of expected kernel values ?

# Idea of the HSIC estimator

MMD in terms of kernel evaluations:

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

## Question

Can we rewrite HSIC in terms of expected kernel values ?

$$\begin{aligned} \text{HSIC}^2(x, y) &= \mathbb{E}_{xy} \mathbb{E}_{x'y'} k(x, x') \ell(y, y') + \mathbb{E}_{xx'} k(x, x') \mathbb{E}_{yy'} \ell(y, y') \\ &\quad - 2\mathbb{E}_{xy} [\mathbb{E}_{x'} k(x, x') \mathbb{E}_{y'} \ell(y, y')]. \end{aligned}$$

# Idea of the HSIC estimator

MMD in terms of kernel evaluations:

$$\text{MMD}^2(\mathbb{P}, \mathbb{Q}) = \mathbb{E}_{x \sim \mathbb{P}, x' \sim \mathbb{P}} k(x, x') + \mathbb{E}_{y \sim \mathbb{Q}, y' \sim \mathbb{Q}} k(y, y') \\ - 2\mathbb{E}_{x \sim \mathbb{P}, y \sim \mathbb{Q}} k(x, y).$$

## Question

Can we rewrite HSIC in terms of expected kernel values ?

$$\text{HSIC}^2(x, y) = \mathbb{E}_{xy} \mathbb{E}_{x'y'} k(x, x') \ell(y, y') + \mathbb{E}_{xx'} k(x, x') \mathbb{E}_{yy'} \ell(y, y') \\ - 2\mathbb{E}_{xy} [\mathbb{E}_{x'} k(x, x') \mathbb{E}_{y'} \ell(y, y')].$$

Empirical estimation results in

$$\widehat{\text{HSIC}}_b^2(x, y) = \frac{1}{n^2} \left\langle \tilde{\mathbf{G}}_x, \tilde{\mathbf{G}}_x \right\rangle_F.$$

# Cocktail party: HSIC demo



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

where  $\mathbf{s}^m$ -s are non-Gaussian & independent.

- Goal:  $\{\mathbf{x}_t\}_{t=1}^T \rightarrow \mathbf{W} = \mathbf{A}^{-1}, \{\mathbf{s}_t\}_{t=1}^T$ ,

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

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 ( $s$ ):



- Observation ( $x$ ):



- Estimated sources ( $\hat{s}$ ):



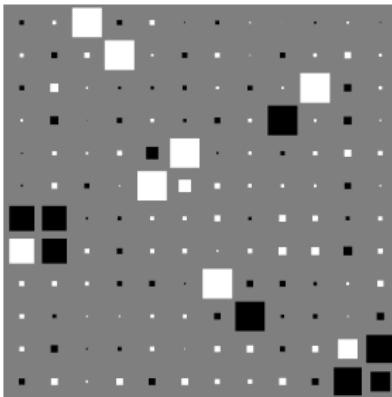
The image displays the words "BROADWAY" in a bold, sans-serif font. Each letter is constructed from numerous small, dark gray or black dots, giving it a granular, point-based appearance. The letters are slightly overlapping, and the overall effect is a dense, textured representation of the text.

# ISA: estimated sources using HSIC, ambiguity

- Estimated sources ( $\hat{s}$ ):



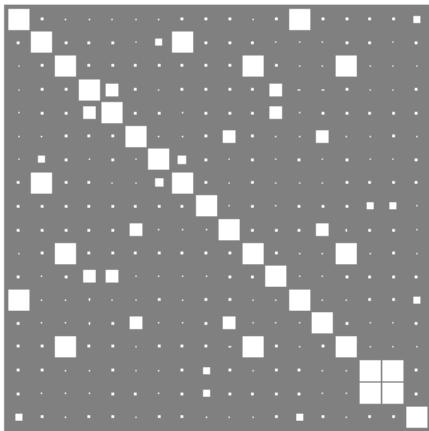
- Performance ( $\hat{W}A$ ), ambiguity:



- $\text{ISA} = \text{ICA} + \text{permutation.}$

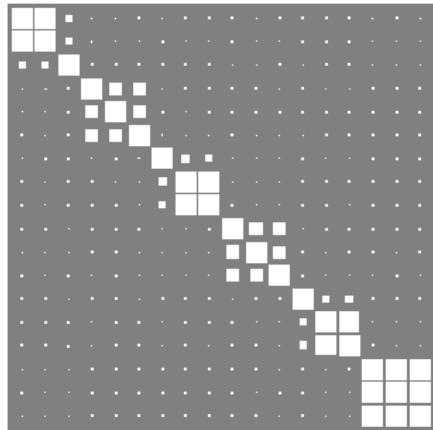
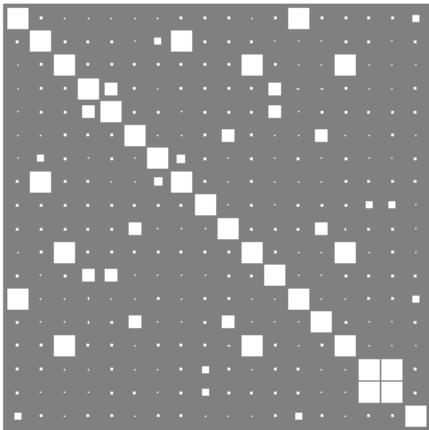
## Conjecture: ISA separation theorem [Cardoso, 1998]

- ISA = ICA + permutation.  $\widehat{\text{HSIC}}(\hat{s}_i, \hat{s}_j)$ . Here:  $\dim(\mathbf{s}^m) = 3$ .



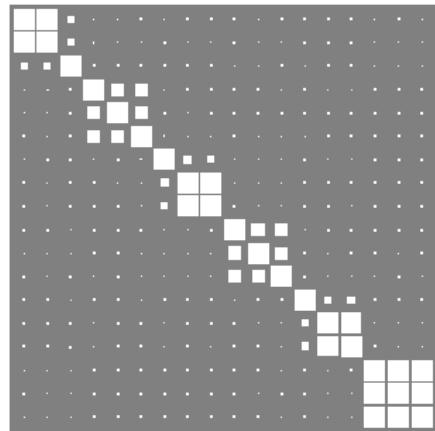
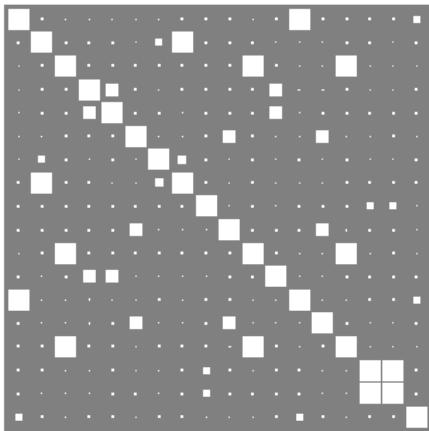
## Conjecture: ISA separation theorem [Cardoso, 1998]

- ISA = ICA + permutation.  $\widehat{\text{HSIC}}(\hat{s}_i, \hat{s}_j)$ . Here:  $\dim(\mathbf{s}^m) = 3$ .



## Conjecture: ISA separation theorem [Cardoso, 1998]

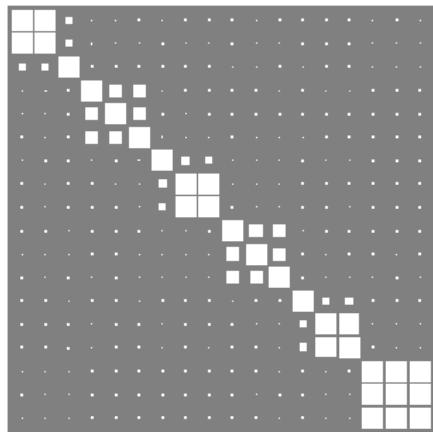
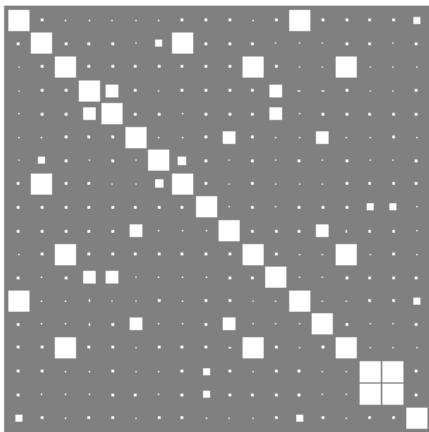
- ISA = ICA + permutation.  $\widehat{\text{HSIC}}(\hat{s}_i, \hat{s}_j)$ . Here:  $\dim(\mathbf{s}^m) = 3$ .



- Basis of the state-of-the-art ISA solvers.

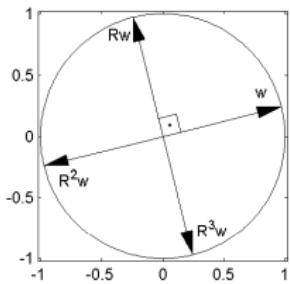
## Conjecture: ISA separation theorem [Cardoso, 1998]

- ISA = ICA + permutation.  $\widehat{\text{HSIC}}(\hat{s}_i, \hat{s}_j)$ . Here:  $\dim(\mathbf{s}^m) = 3$ .



- Basis of the state-of-the-art ISA solvers.
- Sufficient conditions [Szabó et al., 2012]:
  - $\mathbf{s}^m$ : spherical.

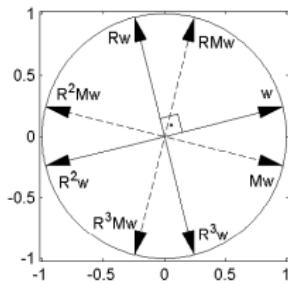
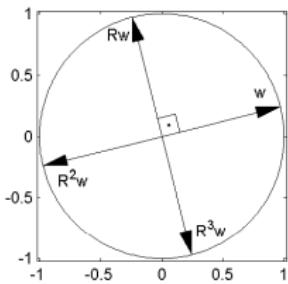
# ISA separation theorem



Invariance to

- $90^\circ$  rotation:  $f(u_1, u_2) = f(-u_2, u_1) = f(-u_1, -u_2) = f(u_2, -u_1)$ .

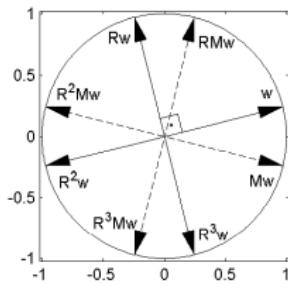
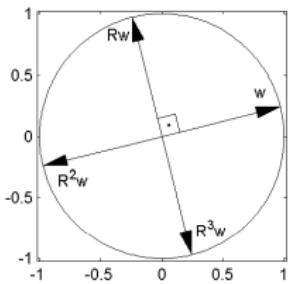
# ISA separation theorem



Invariance to

- $90^\circ$  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



Invariance to

- $90^\circ$  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:  $f(u_1, u_2) = h(\sum_i |u_i|^p)$  ( $p > 0$ ).

Universal kernel (see KCCA)

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

# Universality

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  $\Rightarrow$

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

# 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.

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

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

is universal.

# Universal Taylor kernels

[Steinwart, 2001, Steinwart and Christmann, 2008]

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

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

# Universal Taylor kernels

[Steinwart, 2001, Steinwart and Christmann, 2008]

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

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

- 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, $\alpha > 0$

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

## Universal kernels, $\alpha > 0$

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

# Universal kernels, $\alpha > 0$

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

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

# Universality: notes

- $k$ : universal  $\Leftrightarrow \mathbb{P} \mapsto \mu_{\mathbb{P}}$  is injective on finite signed measures (Hahn-Banach).

- $k$ : universal  $\Leftrightarrow \mathbb{P} \mapsto \mu_{\mathbb{P}}$  is injective on finite signed measures (Hahn-Banach).
- Thus, universal  $\Rightarrow$  characteristic.

- $k$ : universal  $\Leftrightarrow \mathbb{P} \mapsto \mu_{\mathbb{P}}$  is injective on finite signed measures (Hahn-Banach).
- Thus, universal  $\Rightarrow$  characteristic.
- 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, 2016].

Characteristic property, i.e. when MMD is a metric?

[Sriperumbudur et al., 2010b]:

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

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

# Polynomial kernels are not characteristic

[Sriperumbudur et al., 2010b]:

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

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

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

$$\text{MMD}_k^2(\mathbb{P}, \mathbb{Q}) = 2 \|\mathbf{m}_{\mathbb{P}} - \mathbf{m}_{\mathbb{Q}}\|^2 + \left\| \boldsymbol{\Sigma}_{\mathbb{P}} - \boldsymbol{\Sigma}_{\mathbb{Q}} + \mathbf{m}_{\mathbb{P}} \mathbf{m}_{\mathbb{P}}^T - \mathbf{m}_{\mathbb{Q}} \mathbf{m}_{\mathbb{Q}}^T \right\|_F^2,$$

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

# MMD in terms of characteristic functions

Using Bochner's theorem:

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

# MMD in terms of characteristic functions

Using Bochner's theorem:

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

# MMD in terms of characteristic functions

Using Bochner's theorem:

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

# MMD in terms of characteristic functions

Using Bochner's theorem:

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

# MMD in terms of characteristic functions

Using Bochner's theorem:

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

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

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

$$k(\mathbf{x}, \mathbf{x}') = \int_{\mathbb{R}^d} e^{-i\langle \mathbf{x}-\mathbf{x}', \boldsymbol{\omega} \rangle} d\Lambda(\boldsymbol{\omega}).$$

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

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

$$k(\mathbf{x}, \mathbf{x}') = k_0(\mathbf{x} - \mathbf{x}') = \int_{\mathbb{R}^d} e^{-i\langle \mathbf{x}-\mathbf{x}', \boldsymbol{\omega} \rangle} d\Lambda(\boldsymbol{\omega}).$$

Example on  $\mathbb{R}$ :

kernel name	$k_0$	$\hat{k}_0(\omega)$	$\text{supp}(\hat{k}_0)$
Gaussian	$e^{-\frac{x^2}{2\sigma^2}}$	$\sigma e^{-\frac{\sigma^2\omega^2}{2}}$	$\mathbb{R}$
Laplacian	$e^{-\sigma x }$	$\sqrt{\frac{2}{\pi}} \frac{\sigma}{\sigma^2 + \omega^2}$	$\mathbb{R}$
Sinc	$\frac{\sin(\sigma x)}{x}$	$\sqrt{\frac{\pi}{2}} \chi_{[-\sigma, \sigma]}(\omega)$	$[-\sigma, \sigma]$

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

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

$$k(\mathbf{x}, \mathbf{x}') = k_0(\mathbf{x} - \mathbf{x}') = \int_{\mathbb{R}^d} e^{-i\langle \mathbf{x}-\mathbf{x}', \boldsymbol{\omega} \rangle} d\Lambda(\boldsymbol{\omega}).$$

Example on  $\mathbb{R}$ :

kernel name	$k_0$	$\hat{k}_0(\omega)$	$\text{supp}(\hat{k}_0)$
Gaussian	$e^{-\frac{x^2}{2\sigma^2}}$	$\sigma e^{-\frac{\sigma^2\omega^2}{2}}$	$\mathbb{R}$
Laplacian	$e^{-\sigma x }$	$\sqrt{\frac{2}{\pi}} \frac{\sigma}{\sigma^2 + \omega^2}$	$\mathbb{R}$
Sinc	$\frac{\sin(\sigma x)}{x}$	$\sqrt{\frac{\pi}{2}} \chi_{[-\sigma, \sigma]}(\omega)$	$[-\sigma, \sigma]$

Similar characterization  $\exists$  on 'Bochner domains'  
 [Berg et al., 1984, Fukumizu et al., 2009].

## 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  $\mathcal{H}_k$ .

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

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

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

## 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  $\mathcal{H}_k$ .

$$\begin{aligned}\text{MMD}(\mathbb{P}, \mathbb{Q}) &= \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k} \\ &= \sup_{f \in \mathcal{F}} \langle f, \mu_{\mathbb{P}} - \mu_{\mathbb{Q}} \rangle_{\mathcal{H}_k} \\ &= \sup_{f \in \mathcal{F}} (\mathbb{P}f - \mathbb{Q}f), \quad \mathbb{P}f := \int_{\mathcal{X}} f(x) d\mathbb{P}(x).\end{aligned}$$

## 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  $\mathcal{H}_k$ .

$$\begin{aligned}\text{MMD}(\mathbb{P}, \mathbb{Q}) &= \|\mu_{\mathbb{P}} - \mu_{\mathbb{Q}}\|_{\mathcal{H}_k} \\ &= \sup_{f \in \mathcal{F}} \langle f, \mu_{\mathbb{P}} - \mu_{\mathbb{Q}} \rangle_{\mathcal{H}_k} \\ &= \sup_{f \in \mathcal{F}} (\mathbb{P}f - \mathbb{Q}f), \quad \mathbb{P}f := \int_{\mathcal{X}} f(x) d\mathbb{P}(x).\end{aligned}$$

- 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]:

$$\text{MMD}_k(\mathbb{P}, \mathbb{Q}) \leq \sup_{x \in \mathcal{X}} \sqrt{k(x, x)} \text{TV}(\mathbb{P}, \mathbb{Q}).$$

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

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

[Sriperumbudur et al., 2012]:

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

$\mathcal{I}$ -characteristic property, i.e. when HSIC is  
an independence measure?

## Central in applications: characteristic property

- HSIC,  $k = \otimes_{m=1}^M k_m$ ,  $x = (x_m)_{m=1}^M$ :

$$\text{HSIC}_k(\mathbb{P}) := \text{MMD}_{\textcolor{green}{k}} \left( \mathbb{P}, \otimes_{m=1}^M \mathbb{P}_m \right), \quad k(x, x') := \prod_{m=1}^M k_m(x_m, x'_m).$$

# Central in applications: characteristic property

- HSIC,  $k = \otimes_{m=1}^M k_m$ ,  $x = (x_m)_{m=1}^M$ :

$$\text{HSIC}_k(\mathbb{P}) := \text{MMD}_{\textcolor{violet}{k}} \left( \mathbb{P}, \otimes_{m=1}^M \mathbb{P}_m \right), \quad k(x, x') := \prod_{m=1}^M k_m(x_m, x'_m).$$

$k = \otimes_{m=1}^M k_m$  will be called **I-characteristic** if

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

# 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,  $k = \otimes_{m=1}^M k_m$ ,  $x = (x_m)_{m=1}^M$ :

$$\text{HSIC}_k(\mathbb{P}) := \text{MMD}_k\left(\mathbb{P}, \otimes_{m=1}^M \mathbb{P}_m\right), \quad k(x, x') := \prod_{m=1}^M k_m(x_m, x'_m).$$

$k = \otimes_{m=1}^M k_m$  will be called  **$\mathcal{I}$ -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.

### Question

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.

### Question

Extension to  $M \geq 2$ ?

### Main Challenge

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

# Results [Szabó and Sriperumbudur, 2018]

## Proposition (characteristic property)

- $\otimes_{m=1}^M k_m$ : characteristic  $\Rightarrow (k_m)_{m=1}^M$  are characteristic.
- $\Leftarrow [|\mathcal{X}_m| = 2, k_m(x, x') = 2\delta_{x,x'} - 1]$

# Results [Szabó and Sriperumbudur, 2018]

## Proposition (characteristic property)

- $\otimes_{m=1}^M k_m$ : characteristic  $\Rightarrow (k_m)_{m=1}^M$  are characteristic.
- $\Leftarrow$   $[\mathcal{X}_m] = 2, k_m(x, x') = 2\delta_{x,x'} - 1]$

## 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  $\not\Rightarrow \otimes_{m=1}^3 k_m$ :  $\mathcal{I}$ -characteristic [Ex].

# Results [Szabó and Sriperumbudur, 2018]

## Proposition (characteristic property)

- $\otimes_{m=1}^M k_m$ : characteristic  $\Rightarrow (k_m)_{m=1}^M$  are characteristic.
- $\Leftarrow \exists [|\mathcal{X}_m| = 2, k_m(x, x') = 2\delta_{x,x'} - 1]$

## 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  $\Rightarrow \otimes_{m=1}^3 k_m$ :  $\mathcal{I}$ -characteristic [Ex].

## 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 [Szabó and Sriperumbudur, 2018]

## Proposition (characteristic property)

- $\otimes_{m=1}^M k_m$ : characteristic  $\Rightarrow (k_m)_{m=1}^M$  are characteristic.
- $\Leftarrow \exists |X_m| = 2, k_m(x, x') = 2\delta_{x,x'} - 1$

## Proposition ( $\mathcal{I}$ -characteristic property)

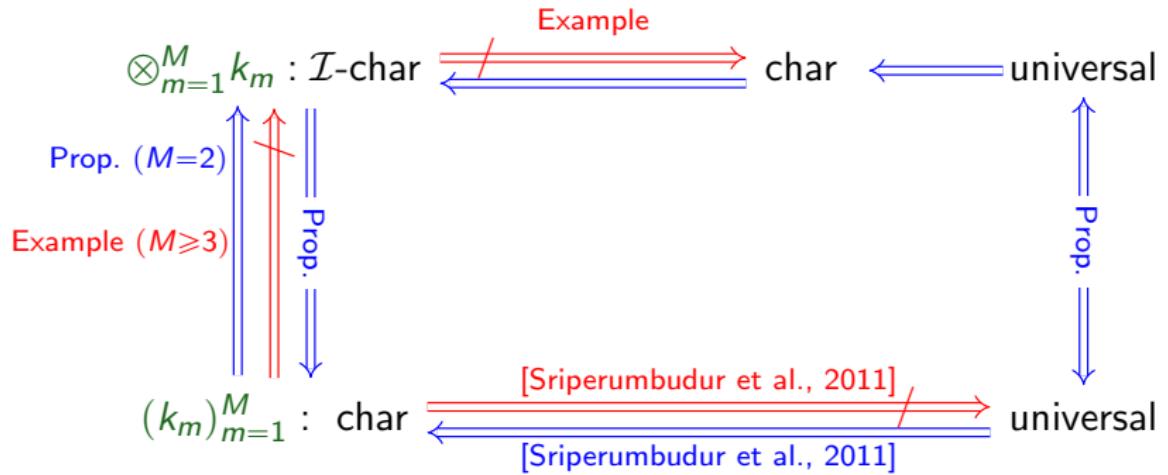
- $k_1, k_2$ : characteristic  $\Rightarrow k_1 \otimes k_2$ :  $\mathcal{I}$ -characteristic.
- $\Leftarrow$ : for  $\forall M \geq 2$ .
- $k_1, k_2, k_3$ : characteristic  $\Rightarrow \otimes_{m=1}^3 k_m$ :  $\mathcal{I}$ -characteristic [Ex].

## Proposition ( $X_m = \mathbb{R}^{d_m}$ , $k_m$ : 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.



# Hypothesis Testing

## Two-sample testing: recall

- Given:
  - $X = \{\mathbf{x}_i\}_{i=1}^n \sim \mathbb{P}$ ,  $Y = \{\mathbf{y}_j\}_{j=1}^n \sim \mathbb{Q}$ .
  - Example:  $\mathbf{x}_i = i^{th}$  happy face,  $\mathbf{y}_j = j^{th}$  sad face.



## Two-sample testing: recall

- Given:

- $X = \{\mathbf{x}_i\}_{i=1}^n \sim \mathbb{P}$ ,  $Y = \{\mathbf{y}_j\}_{j=1}^n \sim \mathbb{Q}$ .
- Example:  $\mathbf{x}_i = i^{th}$  happy face,  $\mathbf{y}_j = j^{th}$  sad face.



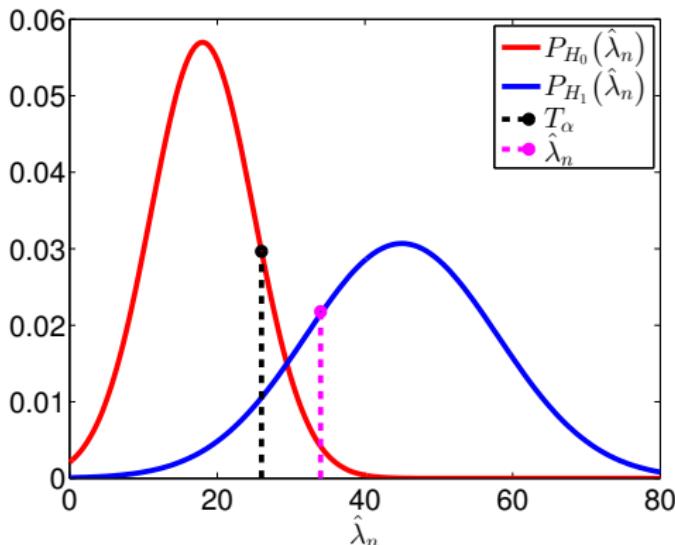
- Problem: using  $X, Y$  test

$$H_0 : \mathbb{P} = \mathbb{Q}, \text{ vs }$$

$$H_1 : \mathbb{P} \neq \mathbb{Q}.$$

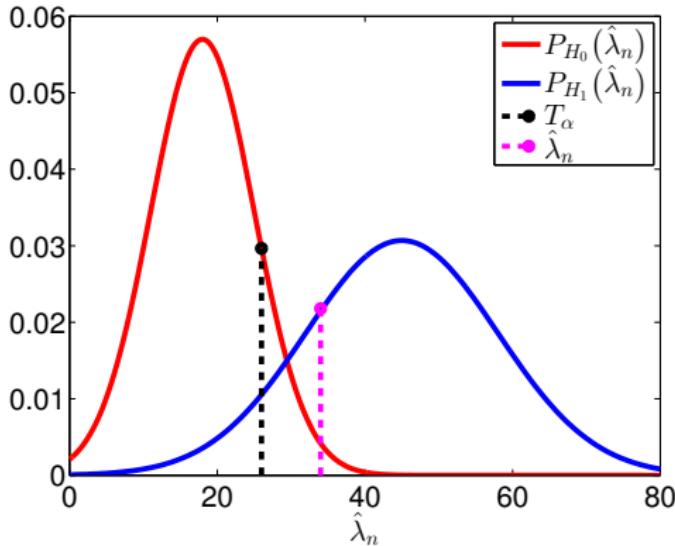
# Ingredients of two-sample test

- Test statistic:  $\hat{\lambda}_n = \hat{\lambda}_n(X, Y)$ , random.
- Significance level:  $\alpha = 0.01$ .
- Under  $H_0$ :  $P_{H_0}(\underbrace{\hat{\lambda}_n \leq T_\alpha}_{\text{correctly accepting } H_0}) = 1 - \alpha$ .



# Ingredients of two-sample test

- 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  $H_1$ :  $P_{H_1}(T_\alpha < \hat{\lambda}_n) = P(\text{correctly rejecting } H_0) =: \text{power}$ .



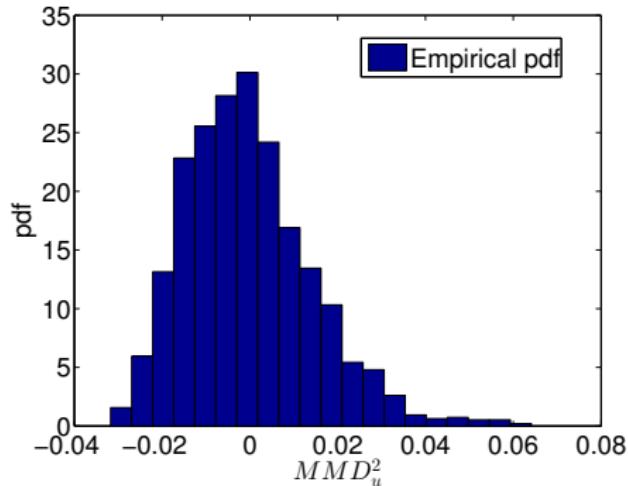
## Two-sample test using MMD asymptotics: $H_0$

Under  $H_0$  [Gretton et al., 2007, Gretton et al., 2012]  $\xrightarrow{\text{U-statistics}}$

$$\widehat{n\text{MMD}_u^2}(\mathbb{P}, \mathbb{P}) \sim \sum_{i=1}^{\infty} \lambda_i (z_i^2 - 2),$$

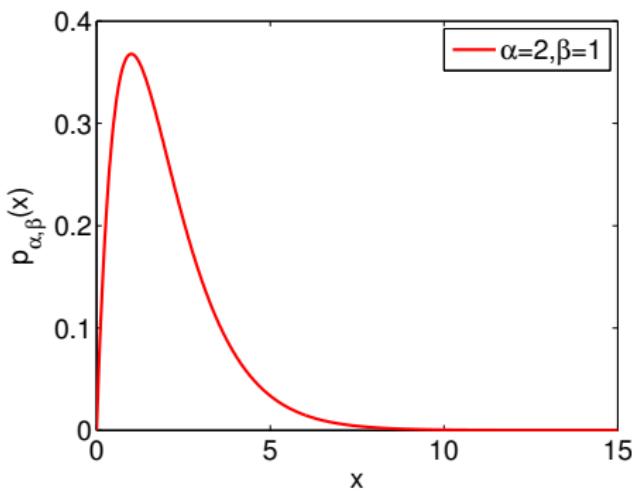
where  $z_i \sim N(0, 2)$  i.i.d.,

$$\int_{\mathcal{X}} \tilde{k}(x, x') v_i(x) d\mathbb{P}(x) = \lambda_i v_i(x'), \quad \tilde{k}(x, x') = \langle \varphi(x) - \mu_{\mathbb{P}}, \varphi(x') - \mu_{\mathbb{P}} \rangle_{\mathcal{H}_k}.$$



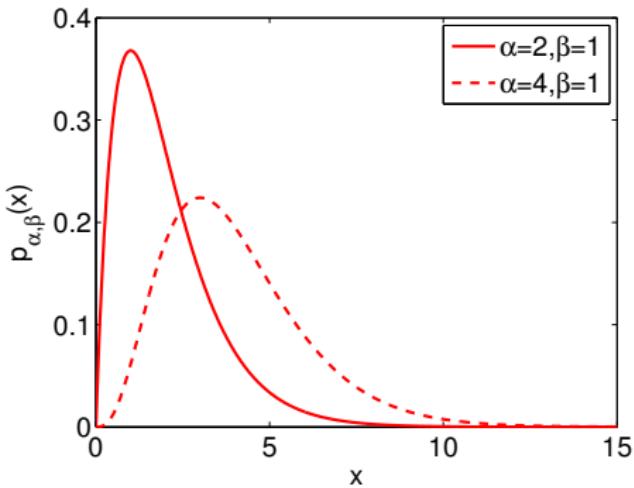
# Null approximations; test statistics: quadratic in time

- Small sample size: permutation test.
- Medium sample size:
  - gamma approximation:



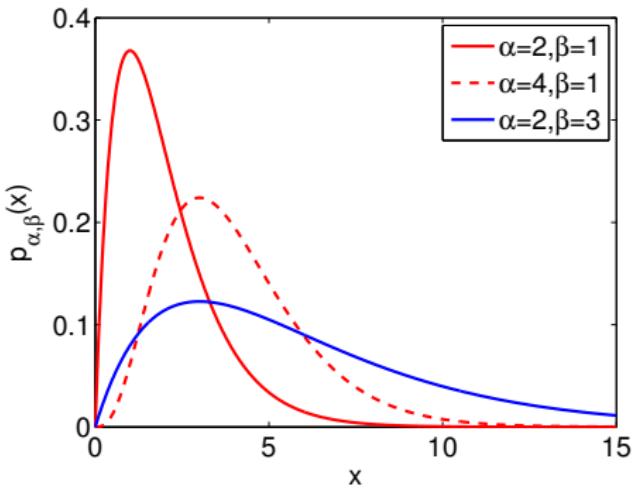
# Null approximations; test statistics: quadratic in time

- Small sample size: permutation test.
- Medium sample size:
  - gamma approximation:



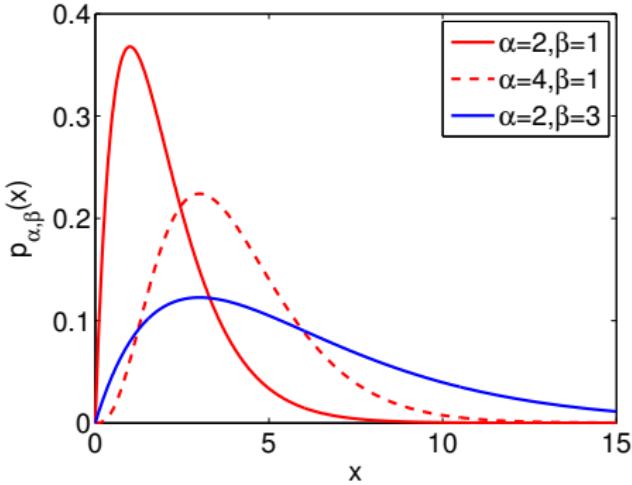
# Null approximations; test statistics: quadratic in time

- Small sample size: permutation test.
- Medium sample size:
  - gamma approximation:



# Null approximations; test statistics: quadratic in time

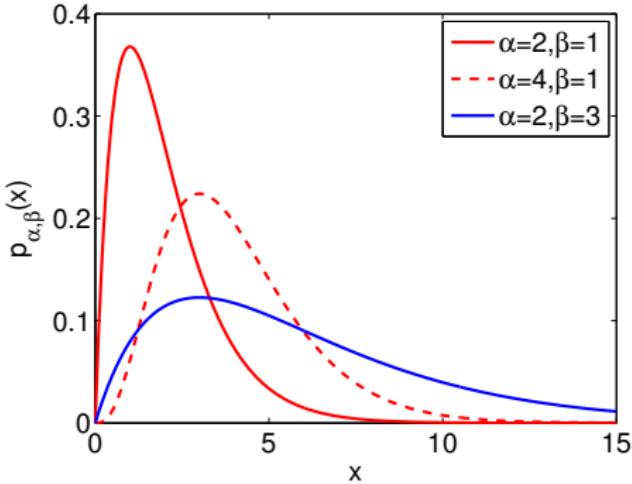
- Small sample size: permutation test.
- Medium sample size:
  - gamma approximation:



- truncated expansion [Gretton et al., 2009].

# Null approximations; test statistics: quadratic in time

- 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] (large var),
  - recent linear methods (soon).

# Independence testing with HSIC

Similary story [Gretton et al., 2008, Pfister et al., 2017]:

- Null asymptotics:

$$\sum_{i=1}^{\infty} \lambda_i z_i^2, \quad z_i \sim N(0, 1).$$

- In practice: permutation-test/gamma-approximation.

## Related work

- 2-sample testing: block-MMD [Zaremba et al., 2013]
  - var ↘

## Related work

- 2-sample testing: [block-MMD](#) [Zaremba et al., 2013]
  - var ↘
- 3-variable [interaction](#) [Sejdinovic et al., 2013a].

## Related work

- 2-sample testing: **block-MMD** [Zaremba et al., 2013]
  - var ↘
- 3-variable **interaction** [Sejdinovic et al., 2013a].
- **Goodness-of-fit** [Chwialkowski et al., 2016].

## Related work

- 2-sample testing: **block-MMD** [Zaremba et al., 2013]
  - var ↘
- 3-variable **interaction** [Sejdinovic et al., 2013a].
- **Goodness-of-fit** [Chwialkowski et al., 2016].
- **Time-series:**
  - independence (stationary → shift) [Chwialkowski and Gretton, 2014],
  - wild bootstrap: [Chwialkowski et al., 2014, Rubenstein et al., 2016].

## Related work

- 2-sample testing: **block-MMD** [Zaremba et al., 2013]
  - var ↘
- 3-variable **interaction** [Sejdinovic et al., 2013a].
- **Goodness-of-fit** [Chwialkowski et al., 2016].
- **Time-series:**
  - independence (stationary → shift) [Chwialkowski and Gretton, 2014],
  - wild bootstrap: [Chwialkowski et al., 2014, Rubenstein et al., 2016].
- **block-HSIC** [Zhang et al., 2017]:
  - RFF acceleration.

## Related work

- 2-sample testing: **block-MMD** [Zaremba et al., 2013]
  - var ↘
- 3-variable **interaction** [Sejdinovic et al., 2013a].
- **Goodness-of-fit** [Chwialkowski et al., 2016].
- **Time-series:**
  - independence (stationary → shift) [Chwialkowski and Gretton, 2014],
  - wild bootstrap: [Chwialkowski et al., 2014, Rubenstein et al., 2016].
- **block-HSIC** [Zhang et al., 2017]:
  - RFF acceleration.
- **Conditional independence** & RFF [Strobl et al., 2017].

# Linear-time Tests

# Linear-time 'MMD'

Idea [Chwialkowski et al., 2015]

Replace  $\|\cdot\|_{\mathcal{H}_k}$  in MMD with  $\|\cdot\|_{L^2(\mathcal{V})}$ . Metric a.s. for analytic & characteristic  $k = k_\sigma$ .

$$\rho(\mathbb{P}, \mathbb{Q}) = \sqrt{\frac{1}{J} \sum_{j=1}^J [\mu_{\mathbb{P}}(\mathbf{v}_j) - \mu_{\mathbb{Q}}(\mathbf{v}_j)]^2}, \quad \mathcal{V} = \{\mathbf{v}_j\}_{j=1}^J,$$

# Linear-time 'MMD'

Idea [Chwialkowski et al., 2015]

Replace  $\|\cdot\|_{\mathcal{H}_k}$  in MMD with  $\|\cdot\|_{L^2(\mathcal{V})}$ . Metric a.s. for analytic & characteristic  $k = k_\sigma$ .

Plug-in estimate:  $\mathcal{O}(n)$ -time.

$$\rho(\mathbb{P}, \mathbb{Q}) = \sqrt{\frac{1}{J} \sum_{j=1}^J [\mu_{\mathbb{P}}(\mathbf{v}_j) - \mu_{\mathbb{Q}}(\mathbf{v}_j)]^2}, \quad \mathcal{V} = \{\mathbf{v}_j\}_{j=1}^J,$$

$$\hat{\rho}(\mathbb{P}, \mathbb{Q}) = \frac{\bar{\mathbf{z}}_n^T \bar{\mathbf{z}}_n}{J},$$

$$\bar{\mathbf{z}}_n = \frac{1}{n} \underbrace{\sum_{i=1}^n \left[ k(\mathbf{x}_i, \mathbf{v}_j) - k(\mathbf{y}_i, \mathbf{v}_j) \right]_{j=1}^J}_{=: \mathbf{z}(\mathbf{x}_i, \mathbf{y}_i)}$$

# Linear-time 'MMD'

Idea [Chwialkowski et al., 2015]

Replace  $\|\cdot\|_{\mathcal{H}_k}$  in MMD with  $\|\cdot\|_{L^2(\mathcal{V})}$ . Metric a.s. for analytic & characteristic  $k = k_\sigma$ .

Plug-in estimate:  $\mathcal{O}(n)$ -time. Whitened test statistic:  $\chi_J^2$  null.

$$\rho(\mathbb{P}, \mathbb{Q}) = \sqrt{\frac{1}{J} \sum_{j=1}^J [\mu_{\mathbb{P}}(\mathbf{v}_j) - \mu_{\mathbb{Q}}(\mathbf{v}_j)]^2}, \quad \mathcal{V} = \{\mathbf{v}_j\}_{j=1}^J,$$

$$\hat{\rho}(\mathbb{P}, \mathbb{Q}) = \frac{\bar{\mathbf{z}}_n^T \bar{\mathbf{z}}_n}{J},$$

$$\bar{\mathbf{z}}_n = \frac{1}{n} \underbrace{\sum_{i=1}^n [k(\mathbf{x}_i, \mathbf{v}_j) - k(\mathbf{y}_i, \mathbf{v}_j)]_{j=1}^J}_{=: \mathbf{z}(\mathbf{x}_i, \mathbf{y}_i)}$$

$$\hat{\lambda}_n = n \bar{\mathbf{z}}_n^T \Sigma_n^{-1} \bar{\mathbf{z}}_n,$$

$$\Sigma_n = \widehat{\text{cov}} (\{\mathbf{z}(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n),$$

# Linear-time 'MMD'

Idea [Chwialkowski et al., 2015], [Jitkrittum et al., 2016]

Replace  $\|\cdot\|_{\mathcal{H}_k}$  in MMD with  $\|\cdot\|_{L^2(\mathcal{V})}$ . Metric a.s. for analytic & characteristic  $k = k_\sigma$ .

Plug-in estimate:  $\mathcal{O}(n)$ -time. Whitened test statistic:  $\chi_J^2$  null. Power opt.

$$\rho(\mathbb{P}, \mathbb{Q}) = \sqrt{\frac{1}{J} \sum_{j=1}^J [\mu_{\mathbb{P}}(\mathbf{v}_j) - \mu_{\mathbb{Q}}(\mathbf{v}_j)]^2}, \quad \mathcal{V} = \{\mathbf{v}_j\}_{j=1}^J,$$

$$\hat{\rho}(\mathbb{P}, \mathbb{Q}) = \frac{\bar{\mathbf{z}}_n^T \bar{\mathbf{z}}_n}{J},$$

$$\bar{\mathbf{z}}_n = \frac{1}{n} \underbrace{\sum_{i=1}^n [k(\mathbf{x}_i, \mathbf{v}_j) - k(\mathbf{y}_i, \mathbf{v}_j)]_{j=1}^J}_{=: \mathbf{z}(\mathbf{x}_i, \mathbf{y}_i)}$$

$$\hat{\lambda}_n = n \bar{\mathbf{z}}_n^T \Sigma_n^{-1} \bar{\mathbf{z}}_n,$$

$$\Sigma_n = \widehat{\text{cov}} (\{\mathbf{z}(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n),$$

$$(\sigma^*, \mathcal{V}^*) = \arg \max_{\sigma, \mathcal{V}} \lambda,$$

$$\lambda = n \mathbf{m}^T \Sigma_n^{-1} \mathbf{m}.$$

# Linear-time 'HSIC' [Jitkrittum et al., 2017]

Use different norm of the witness function ( $u$ ):

$$\text{HSIC}(x, y) = \|\mu_{xy} - \mu_x \otimes \mu_y\|_{\mathcal{H}_{k_1 \otimes k_2}}, \quad u(\mathbf{v}, \mathbf{w}) = \mu_{xy}(\mathbf{v}, \mathbf{w}) - \mu_x(\mathbf{v})\mu_y(\mathbf{w}),$$

# Linear-time 'HSIC' [Jitkrittum et al., 2017]

Use different norm of the witness function ( $u$ ):

$$\text{HSIC}(x, y) = \|\mu_{xy} - \mu_x \otimes \mu_y\|_{\mathcal{H}_{k_1 \otimes k_2}}, \quad u(\mathbf{v}, \mathbf{w}) = \mu_{xy}(\mathbf{v}, \mathbf{w}) - \mu_x(\mathbf{v})\mu_y(\mathbf{w}),$$

$$\text{FSIC}(x, y) = \sqrt{\frac{1}{J} \sum_{j=1}^J u^2(\mathbf{v}_j, \mathbf{w}_j)}, \quad \mathcal{V} = \{(\mathbf{v}_j, \mathbf{w}_j)\}_{j=1}^J,$$

# Linear-time 'HSIC' [Jitkrittum et al., 2017]

Use different norm of the witness function ( $u$ ):

$$\text{HSIC}(x, y) = \|\mu_{xy} - \mu_x \otimes \mu_y\|_{\mathcal{H}_{k_1 \otimes k_2}}, \quad u(\mathbf{v}, \mathbf{w}) = \mu_{xy}(\mathbf{v}, \mathbf{w}) - \mu_x(\mathbf{v})\mu_y(\mathbf{w}),$$

$$\begin{aligned}\text{FSIC}(x, y) &= \sqrt{\frac{1}{J} \sum_{j=1}^J u^2(\mathbf{v}_j, \mathbf{w}_j)}, & \mathcal{V} &= \{(\mathbf{v}_j, \mathbf{w}_j)\}_{j=1}^J, \\ &= \|u\|_{L^2(\mathcal{V})}.\end{aligned}$$

# Linear-time 'HSIC' [Jitkrittum et al., 2017]

Use different norm of the witness function ( $u$ ):

$$\text{HSIC}(x, y) = \|\mu_{xy} - \mu_x \otimes \mu_y\|_{\mathcal{H}_{k_1 \otimes k_2}}, \quad u(\mathbf{v}, \mathbf{w}) = \mu_{xy}(\mathbf{v}, \mathbf{w}) - \mu_x(\mathbf{v})\mu_y(\mathbf{w}),$$

$$\begin{aligned}\text{FSIC}(x, y) &= \sqrt{\frac{1}{J} \sum_{j=1}^J u^2(\mathbf{v}_j, \mathbf{w}_j)}, & \mathcal{V} &= \{(\mathbf{v}_j, \mathbf{w}_j)\}_{j=1}^J, \\ &= \|u\|_{L^2(\mathcal{V})}.\end{aligned}$$

- Whitening  $\Rightarrow \chi_J^2$  null. Computation:  $\mathcal{O}(n)$ . Power optimization.

# Linear-time 'HSIC' [Jitkrittum et al., 2017]

Use different norm of the witness function ( $u$ ):

$$\text{HSIC}(x, y) = \|\mu_{xy} - \mu_x \otimes \mu_y\|_{\mathcal{H}_{k_1 \otimes k_2}}, \quad u(\mathbf{v}, \mathbf{w}) = \mu_{xy}(\mathbf{v}, \mathbf{w}) - \mu_x(\mathbf{v})\mu_y(\mathbf{w}),$$

$$\begin{aligned}\text{FSIC}(x, y) &= \sqrt{\frac{1}{J} \sum_{j=1}^J u^2(\mathbf{v}_j, \mathbf{w}_j)}, & \mathcal{V} &= \{(\mathbf{v}_j, \mathbf{w}_j)\}_{j=1}^J, \\ &= \|u\|_{L^2(\mathcal{V})}.\end{aligned}$$

- Whitening  $\Rightarrow \chi_J^2$  null. Computation:  $\mathcal{O}(n)$ . Power optimization.
- Alternative view:  $u(\mathbf{v}, \mathbf{w}) = \text{cov}_{\mathbf{xy}}(k_1(\mathbf{x}, \mathbf{v}), k_2(\mathbf{y}, \mathbf{w})) = (\mathbf{v}, \mathbf{w})^{th}$  entry of

$$C_{xy} = \mathbb{E}_{xy} [\varphi_1(x) \otimes \varphi_2(y)] - \mu_x \otimes \mu_y.$$

We

- assumed analytic, characteristic, bounded kernels.
- replaced the RKHS norm with  $L^2(\mathcal{V})$  norm.

In linear-time '**MMD**' and '**HSIC**', respectively:

$$\begin{aligned}\mathbb{P} = \mathbb{Q} &\Leftrightarrow \mu_{\mathbb{P}-\mathbb{Q}} = 0, \\ \mathbb{P} = \mathbb{P}_1 \otimes \mathbb{P}_2 &\Leftrightarrow \mu_{\mathbb{P}-\mathbb{P}_1 \otimes \mathbb{P}_2} = 0.\end{aligned}$$

# Goodness-of-fit

Let  $d = 1$ . Stein operator of model  $p$

$$(T_p f)(x) = \frac{[p(x)f(x)]'}{p(x)} = [\log p(x)]'f(x) + f'(x).$$

# Goodness-of-fit

Let  $d = 1$ . Stein operator of model  $\textcolor{blue}{p}$

$$(T_{\textcolor{blue}{p}} f)(x) = \frac{[\textcolor{blue}{p}(x)f(x)]'}{\textcolor{blue}{p}(x)} = [\log \textcolor{blue}{p}(x)]'f(x) + f'(x).$$

Under  $\lim_{|x| \rightarrow \infty} f(x)p(x) = 0$  (integration by parts):

$$\textcolor{blue}{p} = \textcolor{red}{q} \Rightarrow \mathbb{E}_{x \sim \textcolor{red}{q}}(T_{\textcolor{blue}{p}} f)(x) = 0.$$

# Goodness-of-fit

Let  $d = 1$ . Stein operator of model  $\textcolor{blue}{p}$

$$(T_{\textcolor{blue}{p}} f)(x) = \frac{[\textcolor{blue}{p}(x)f(x)]'}{\textcolor{blue}{p}(x)} = [\log \textcolor{blue}{p}(x)]'f(x) + f'(x).$$

Under  $\lim_{|x| \rightarrow \infty} f(x)p(x) = 0$  (integration by parts):

$$\textcolor{blue}{p} = \textcolor{red}{q} \Rightarrow \mathbb{E}_{x \sim \textcolor{red}{q}}(T_{\textcolor{blue}{p}} f)(x) = 0.$$

Let us take the unit ball of  $\mathcal{H}_k$ :

$$\sup_{\|f\|_{\mathcal{H}_k} \leq 1} \mathbb{E}_{x \sim \textcolor{red}{q}}(T_{\textcolor{blue}{p}} f)(x) = \|g\|_{\mathcal{H}_k}, \quad g(v) = \mathbb{E}_{x \sim \textcolor{red}{q}} \frac{\partial_x [\textcolor{blue}{p}(x)k(x, v)]}{\textcolor{blue}{p}(x)}.$$

## Goodness-of-fit

[Chwialkowski et al., 2016, Liu et al., 2016]

Let  $d = 1$ . Stein operator of model  $p$

$$(T_p f)(x) = \frac{[\log p(x)f(x)]'}{p(x)} = [\log p(x)]'f(x) + f'(x).$$

Under  $\lim_{|x| \rightarrow \infty} f(x)p(x) = 0$  (integration by parts):

$$p = q \Rightarrow \mathbb{E}_{x \sim q}(T_p f)(x) = 0.$$

Let us take the unit ball of  $\mathcal{H}_k$ :

$$\sup_{\|f\|_{\mathcal{H}_k} \leq 1} \mathbb{E}_{x \sim q}(T_p f)(x) = \|g\|_{\mathcal{H}_k}, \quad g(v) = \mathbb{E}_{x \sim q} \frac{\partial_x [\log p(x) k(x, v)]}{p(x)}.$$

For universal  $k$ :

$$p = q \Leftrightarrow g = 0 \text{ (witness)}.$$

Goodness-of-fit [Jitkrittum et al., 2017],

[Chwialkowski et al., 2016, Liu et al., 2016]

Let  $d = 1$ . Stein operator of model  $p$

$$(T_p f)(x) = \frac{[p(x)f(x)]'}{p(x)} = [\log p(x)]'f(x) + f'(x).$$

Under  $\lim_{|x| \rightarrow \infty} f(x)p(x) = 0$  (integration by parts):

$$p = q \Rightarrow \mathbb{E}_{x \sim q}(T_p f)(x) = 0.$$

Let us take the unit ball of  $\mathcal{H}_k$ :

$$\sup_{\|f\|_{\mathcal{H}_k} \leq 1} \mathbb{E}_{x \sim q}(T_p f)(x) = \|g\|_{\mathcal{H}_k}, \quad g(v) = \mathbb{E}_{x \sim q} \frac{\partial_x [p(x)k(x, v)]}{p(x)}.$$

For universal  $k$ :

$$p = q \Leftrightarrow g = 0 \text{ (witness)}.$$

$L^2(\mathcal{V})$  trick goes through.

# Numerical Illustrations

## 2-sample testing: parameter settings

- Gaussian kernel ( $\sigma$ ).  $\alpha = 0.01$ .  $J = 1$ . Repeat 500 trials.
- Report **rejection rate of  $H_0$**
- Compare 4 methods
  - **ME-full**: Optimize  $\mathcal{V}$  and  $\sigma$ .
  - **ME-grid**: Optimize  $\sigma$ . 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 discriminative 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
± vs. ±	201	.010	.012	.018	.008
+ vs. −	201	.998	.656	1.00	.578



- Learned test location (averaged) =

# Independence testing: parameters

- $k_1, k_2$ : 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

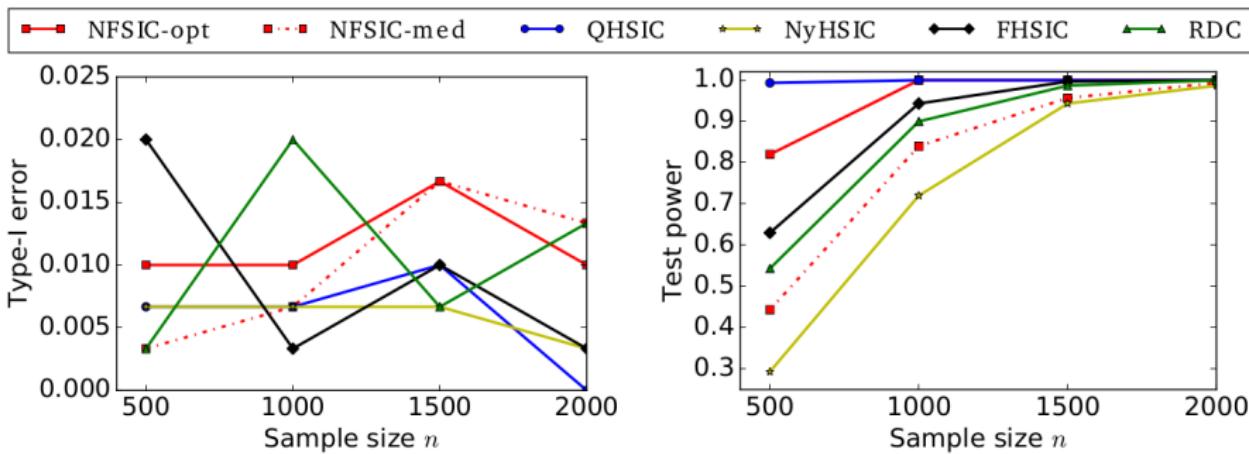
Song ( $x$ ) vs. year of release ( $y$ ).

- Western commercial tracks from 1922 to 2011 [Bertin-Mahieux et al., 2011].
- $x \in \mathbb{R}^{90}$ : audio features.
- **Left**: break  $(x, y)$  pairs, i.e.  $H_0$ ; **right**:  $H_1$  is true.

# Demo-1: million song data

Song ( $x$ ) vs. year of release ( $y$ ).

- Western commercial tracks from 1922 to 2011 [Bertin-Mahieux et al., 2011].
- $x \in \mathbb{R}^{90}$ : audio features.
- **Left:** break ( $x, y$ ) pairs, i.e.  $H_0$ ; **right:**  $H_1$  is true.



## Demo-2: videos and captions

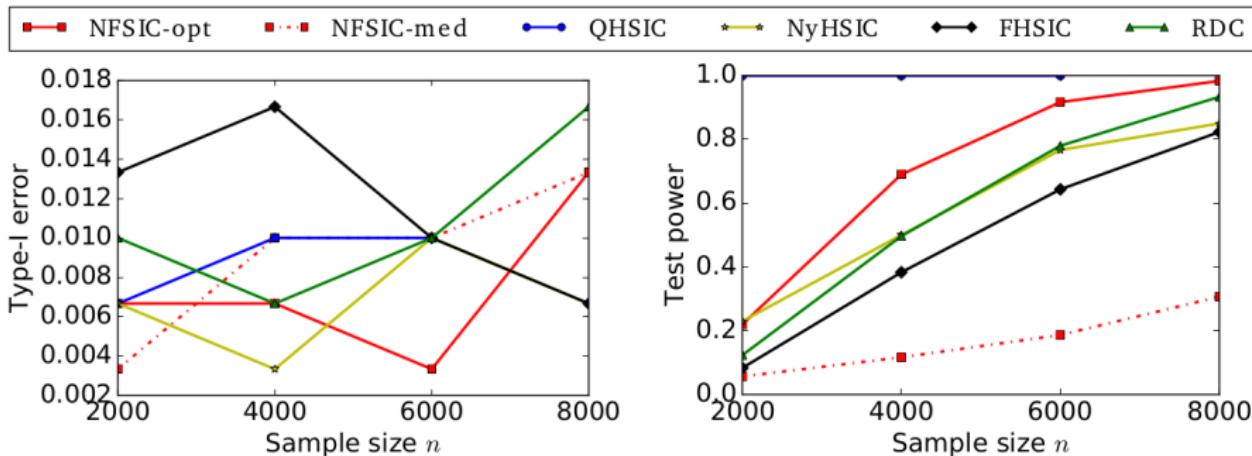
Youtube video ( $x$ ) vs. caption ( $y$ ).

- VideoStory46K [Habibian et al., 2014]
- $x \in \mathbb{R}^{2000}$ : Fisher vector encoding of motion boundary histograms [Wang and Schmid, 2013].
- $y \in \mathbb{R}^{1878}$ : bag of words. TF.
- **Left**: break  $(x, y)$  pairs, i.e.  $H_0$ ; **right**:  $H_1$  is true.

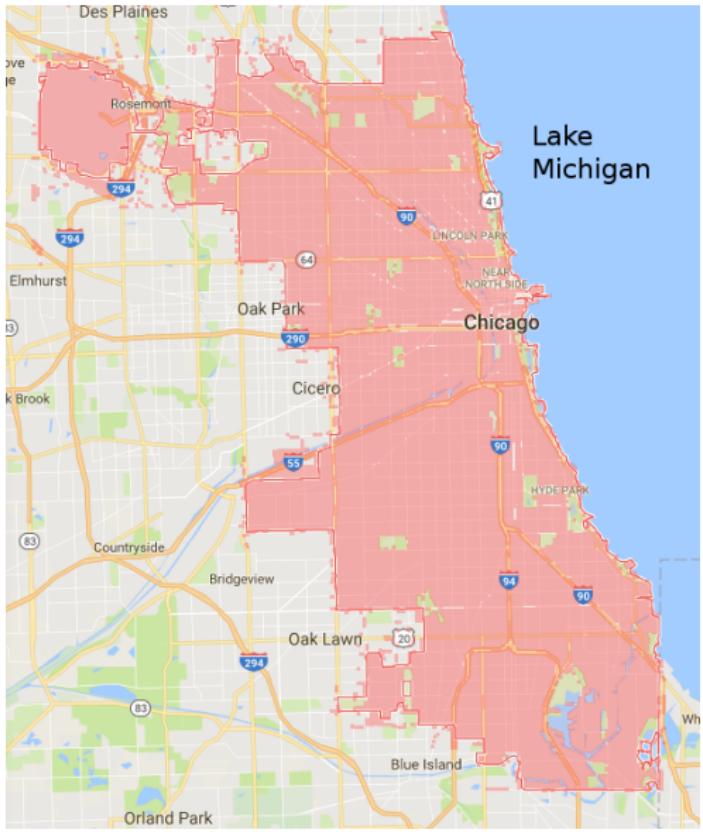
## Demo-2: videos and captions

Youtube video ( $x$ ) vs. caption ( $y$ ).

- VideoStory46K [Habibian et al., 2014]
- $x \in \mathbb{R}^{2000}$ : Fisher vector encoding of motion boundary histograms [Wang and Schmid, 2013].
- $y \in \mathbb{R}^{1878}$ : bag of words. TF.
- **Left:** break  $(x, y)$  pairs, i.e.  $H_0$ ; **right:**  $H_1$  is true.

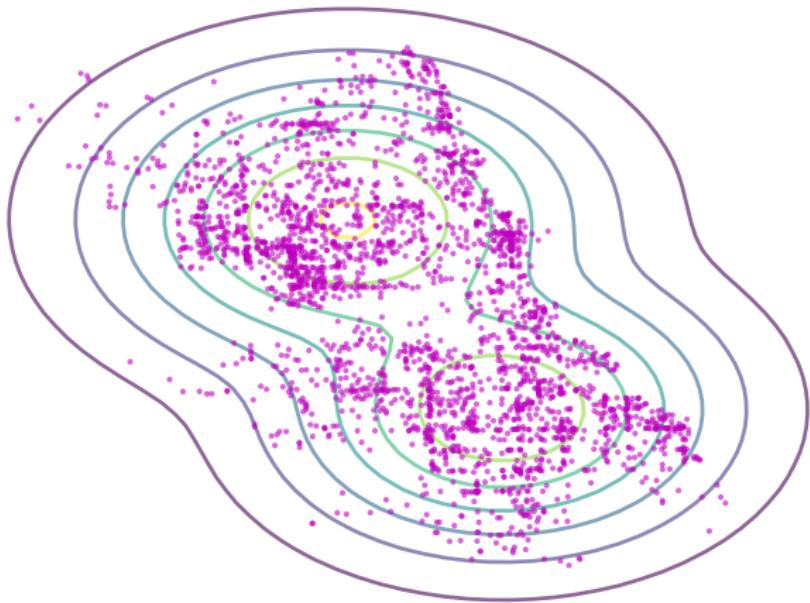


# Goodness-of-Fit Demo

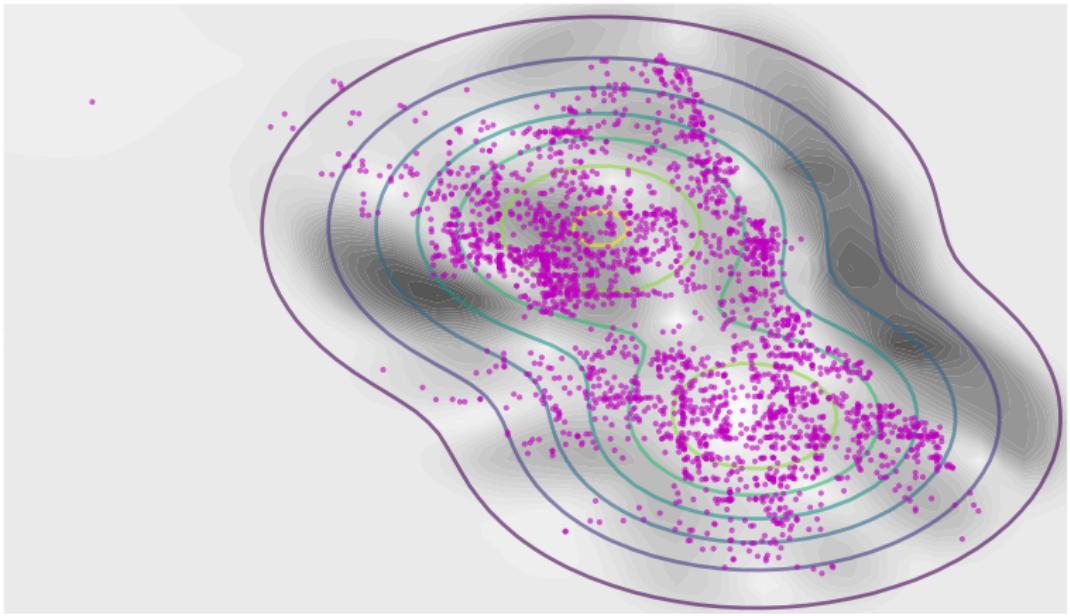




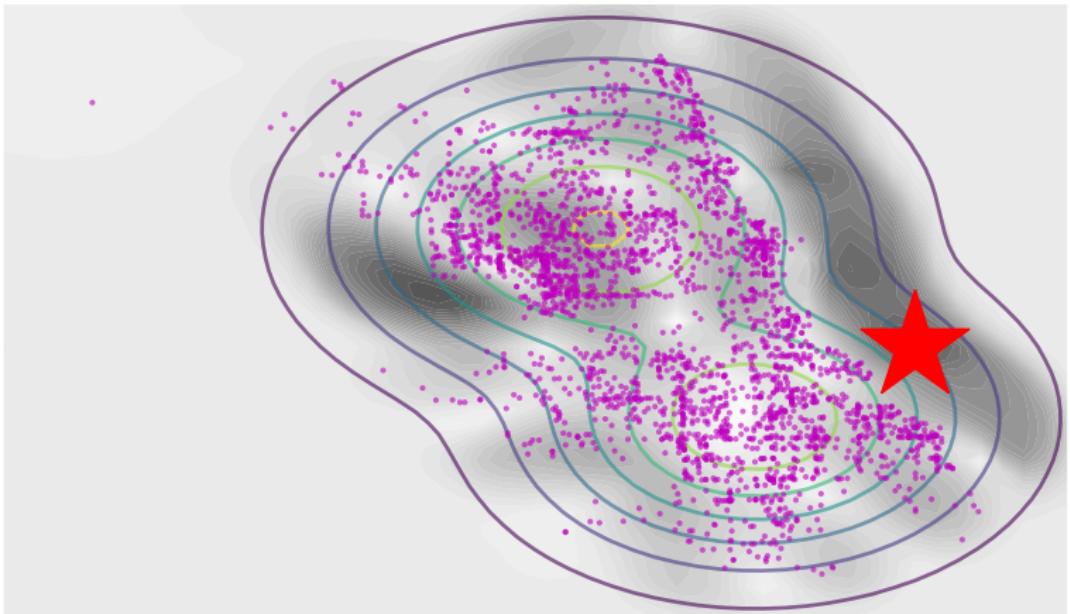
Robbery events (lat/long coordinates)  $\sim q$ .



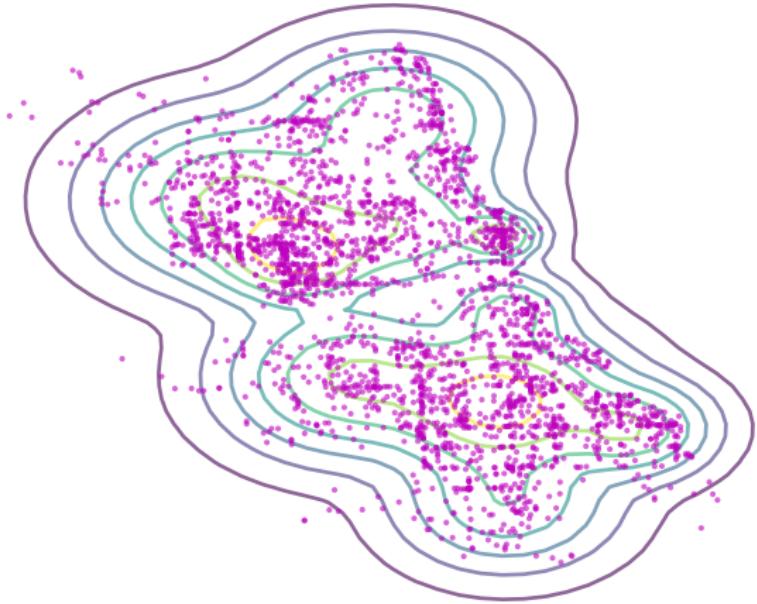
Model  $p$ : 2-component Gaussian mixture.



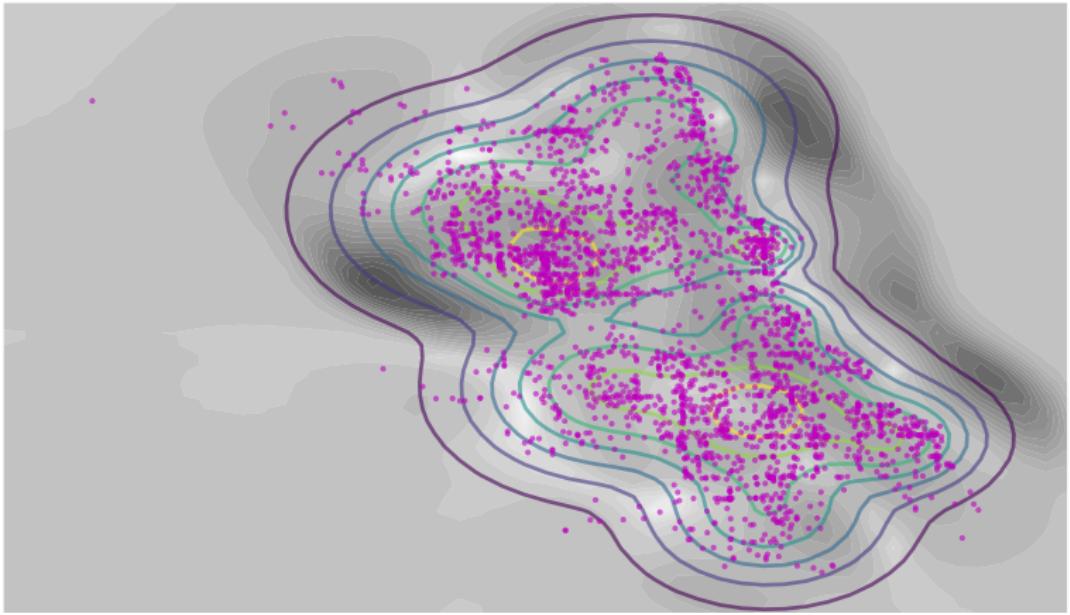
Score surface



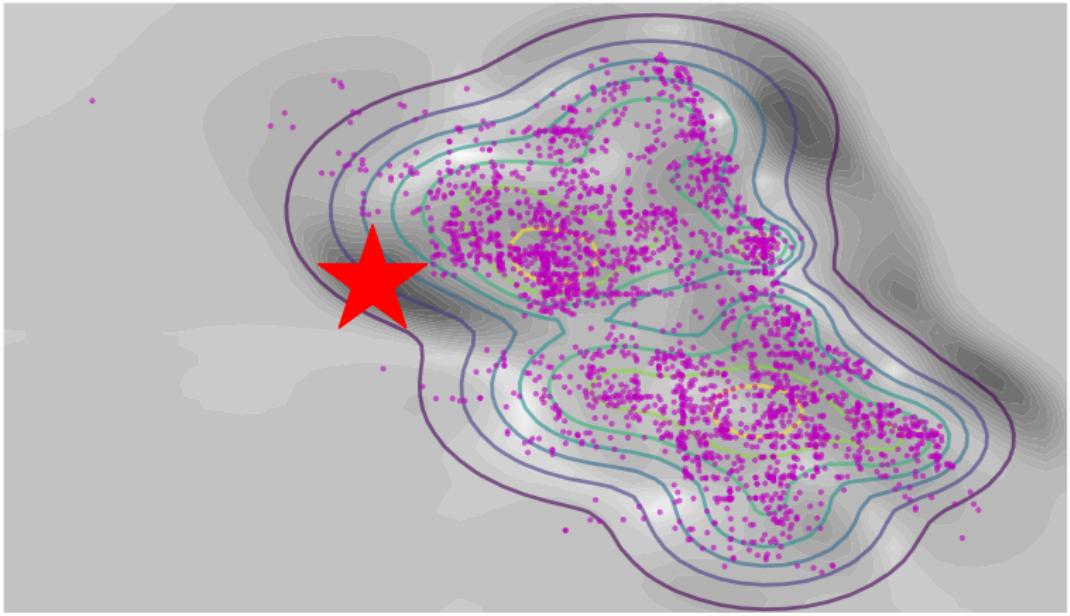
★ = optimized  $v$ .  
No robbery in Lake Michigan.



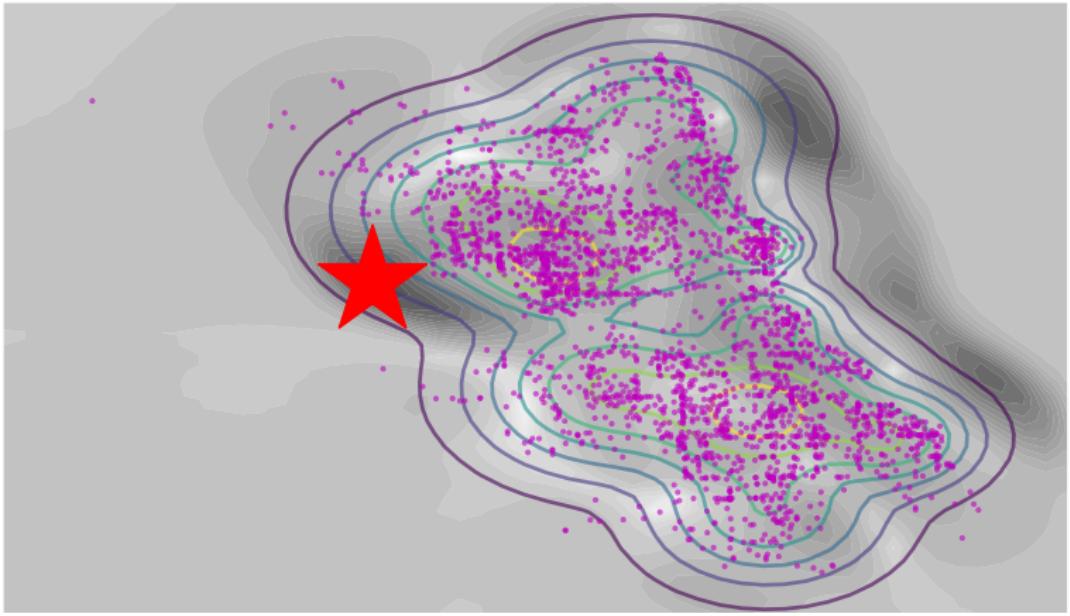
Model  $p$ : 10-component Gaussian mixture.



Capture the right tail better.



Still, does not capture the left tail.



Still, does not capture the left tail.

**Sharp boundary (geography of Chicago)  $\neq$  Gaussian tails.  $\rightarrow$  interpretable features**

- Motivation: infoT objectives, hypothesis testing.
- Kernels, RKHS: definitions, construction.
- Kernel applications: classification, ridge regression, PCA.
- MMD, HSIC, KCCA.
- Characteristic, universal,  $\mathcal{I}$ -characteristic property.
- Hypothesis testing: quadratic & linear-time methods.

Thank you for the attention!



-  Bach, F. R. and Jordan, M. I. (2002).  
Kernel independent component analysis.  
*Journal of Machine Learning Research*, 3:1–48.
-  Balasubramanian, K., Li, T., and Yuan, M. (2017).  
On the optimality of kernel-embedding based goodness-of-fit tests.  
Technical report.  
(<https://arxiv.org/abs/1709.08148>).
-  Berg, C., Christensen, J. P. R., and Ressel, P. (1984).  
*Harmonic Analysis on Semigroups*.  
Springer-Verlag.
-  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., Deshmukh, A. A., Dogan, U., Lee, G., and Scott, C. (2017).  
Domain generalization by marginal transfer learning.  
Technical report.  
(<https://arxiv.org/abs/1711.07910>).
-  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.
-  Borgwardt, K. M., Gretton, A., Rasch, M. J., Kriegel, H.-P., Schölkopf, B., and Smola, A. (2006).  
Integrating structured biological data by kernel maximum mean discrepancy.  
*Bioinformatics*, 22(14):e49–e57.
-  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.
-  Chwialkowski, K. and Gretton, A. (2014).  
A kernel independence test for random processes.  
In *International Conference on Machine Learning (ICML; PMLR)*, volume 32, pages 1422–1430.
-  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.

 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.

-  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. (2009).  
 Characteristic kernels on groups and semigroups.  
 In *Advances in Neural Information Processing Systems (NIPS)*, pages 473–480.
-  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.
-  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>).
-  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. (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 *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>).
-  Hein, M. and Bousquet, O. (2005).  
Hilbertian metrics and positive definite kernels on probability measures.

In *International Conference on AI 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 (PMLR)*,  
volume 37, pages 2982–2990.
-  Jitkrittum, W., Szabó, Z., Chwialkowski, K., and Gretton, A. (2016).  
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. (2017).

An adaptive test of independence with analytic kernel embeddings.

In *International Conference on Machine Learning (ICML; PMLR)*, volume 70, pages 1742–1751.

 Kashima, H. and Koyanagi, T. (2002).

Kernels for semi-structured data.

In *International Conference on Machine Learning (ICML)*, pages 291–298.

 Kim, B., Khanna, R., and Koyejo, O. O. (2016).

Examples are not enough, learn to criticize! criticism for interpretability.

In *Advances in Neural Information Processing Systems (NIPS)*, pages 2280–2288.

 Kondor, R. and Pan, H. (2016).

The multiscale Laplacian graph kernel.

In *Advances in Neural Information Processing Systems (NIPS)*, pages 2982–2990.

-  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.
-  Kybic, J. (2004). High-dimensional mutual information estimation for image registration.  
In *IEEE International Conference on Image Processing (ICIP)*, pages 1779–1782.
-  Law, H. C. L., Sutherland, D. J., Sejdinovic, D., and Flaxman, S. (2018). Bayesian approaches to distribution regression.  
In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 84, pages 1167–1176.
-  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.

-  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.
-  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.
-  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., Muandet, K., Schölkopf, B., and Tolstikhin, I. (2015).  
Towards a learning theory of cause-effect inference.  
*International Conference on Machine Learning (ICML; PMLR)*, 37:1452–1461.
-  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.
-  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.

 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.

 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 398–407.

 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; PMLR)*, volume 31, pages 507–515.

 Rényi, A. (1959).

On measures of dependence.

*Acta Mathematica Academiae Scientiarum Hungaricae*,  
10:441–451.

 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.

-  Schölkopf, B., Herbrich, R., and Smola, A. J. (2001).  
A generalized representer theorem.  
In *Conference on Learning Theory (COLT)*, pages 416–426.
-  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., 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.

- Simon-Gabriel, C.-J. and Schölkopf, B. (2016).  
Kernel distribution embeddings: Universal kernels, characteristic kernels and kernel metrics on distributions.  
Technical report, MPI Tübingen.  
(<https://arxiv.org/abs/1604.05251>).
- 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., Fukumizu, K., Gretton, A., Schölkopf, B., and Lanckriet, G. R. G. (2012).  
On the empirical estimation of integral probability metrics.

-  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; PMLR)*, volume 9, pages 781–788.
-  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.
-  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).

-  Strobl, E. V., Visweswaran, S., and Zhang, K. (2017).  
Approximate kernel-based conditional independence tests for  
fast non-parametric causal discovery.  
*Technical report.*  
(<https://arxiv.org/abs/1702.03877>).
-  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. K. (2018).  
Characteristic and universal tensor product kernels.  
*Journal of Machine Learning Research*, 18:1–29.
-  Szabó, Z., Sriperumbudur, B. K., Póczos, B., and Gretton, A. (2016).  
Learning theory for distribution regression.

-  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.

-  Yamada, M., Umezu, Y., Fukumizu, K., and Takeuchi, I. (2018).  
Post selection inference with kernels.  
In *International Conference on Artificial Intelligence and Statistics (AISTATS; PMLR)*, volume 84, pages 152–160.
-  Yu, Y., Cheng, H., Schuurmans, D., and Szepesvári, C. (2013).  
Characterizing the representer theorem.  
In *International Conference on Machine Learning (ICML; PMLR)*, volume 28, pages 570–578.
-  Zaheer, M., Kottur, S., Ravanbakhsh, S., Póczos, B., Salakhutdinov, R. R., and Smola, A. J. (2017).  
Deep sets.  
In *Advances in Neural Information Processing Systems (NIPS)*, pages 3394–3404.
-  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, K., Schölkopf, B., Muandet, K., and Wang, Z. (2013). Domain adaptation under target and conditional shift. *Journal of Machine Learning Research*, 28(3):819–827.
-  Zhang, Q., Filippi, S., Gretton, A., and Sejdinovic, D. (2017). Large-scale kernel methods for independence testing. *Statistics and Computing*, pages 1–18.
-  Zolotarev, V. M. (1983). Probability metrics. *Theory of Probability and its Applications*, 28:278–302.