

Lecture 2: Mappings of Probabilities to RKHS and Applications

MLSS Tübingen, 2015

Arthur Gretton

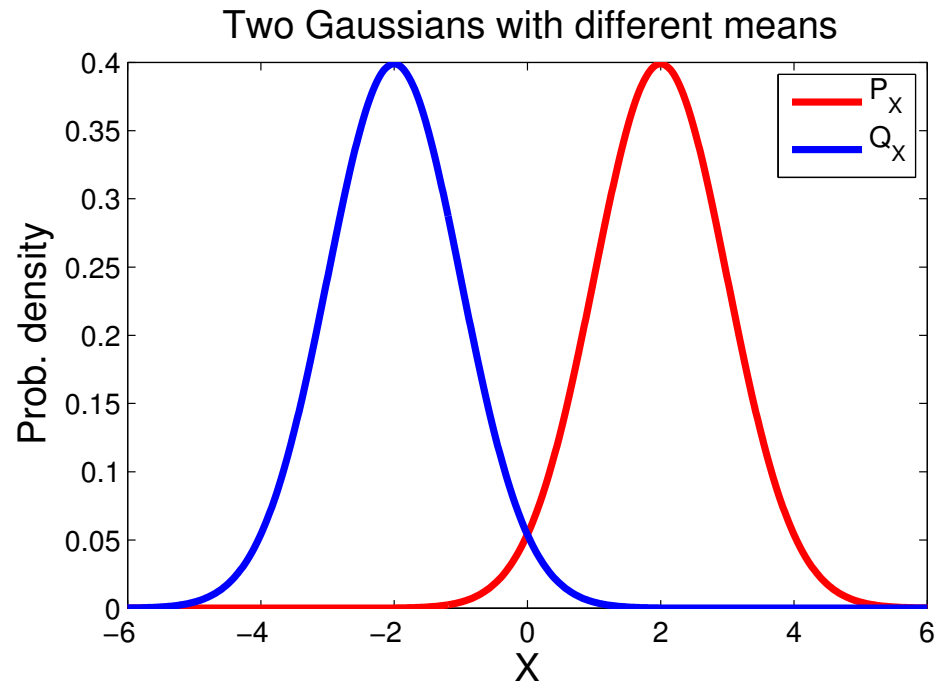
Gatsby Unit, CSML, UCL

Outline

- **Kernel metric** on the space of **probability measures**
 - Function revealing differences in distributions
 - Distance between **means in space of features** (**RKHS**)
 - **Independence measure**: features of joint minus product of marginals
- **Characteristic kernels**: feature space mappings of probabilities **unique**
- **Two-sample, independence tests** for (almost!) any data type
 - distributions on strings, **images**, graphs, groups (rotation matrices), semigroups, . . .
- **Advanced topics**
 - testing on big data, **kernel choice**
 - **Energy distance/distance covariance**: special case of kernel statistic

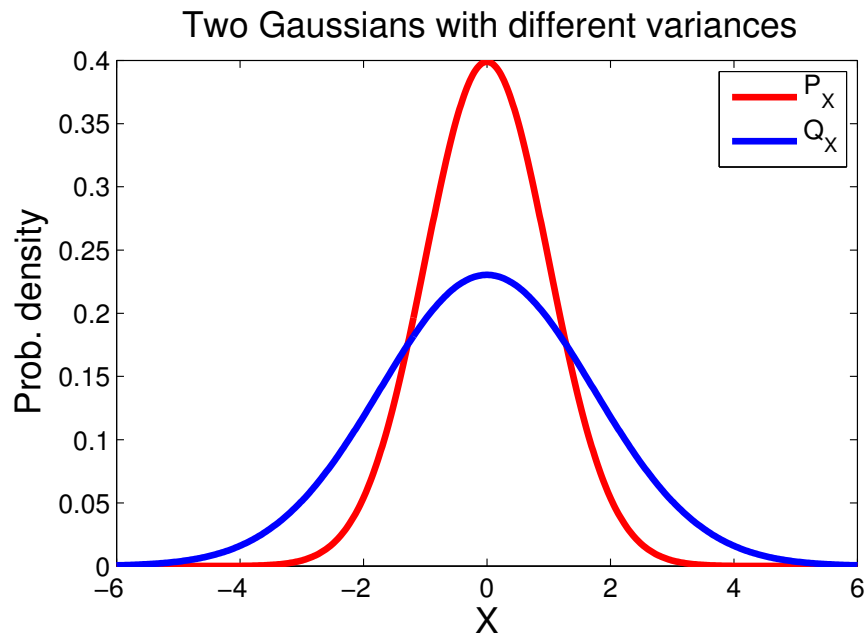
Feature mean difference

- Simple example: 2 Gaussians with different means
- Answer: t -test



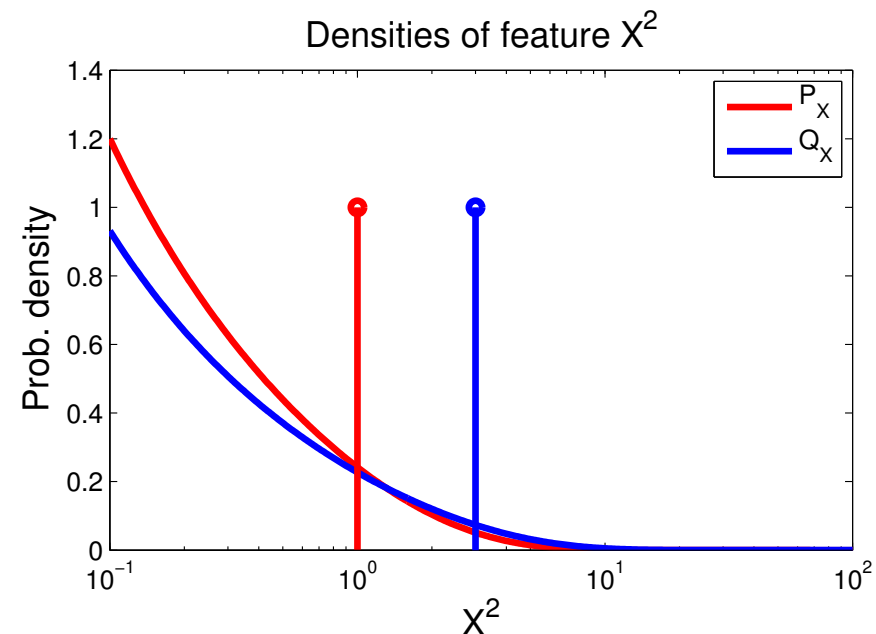
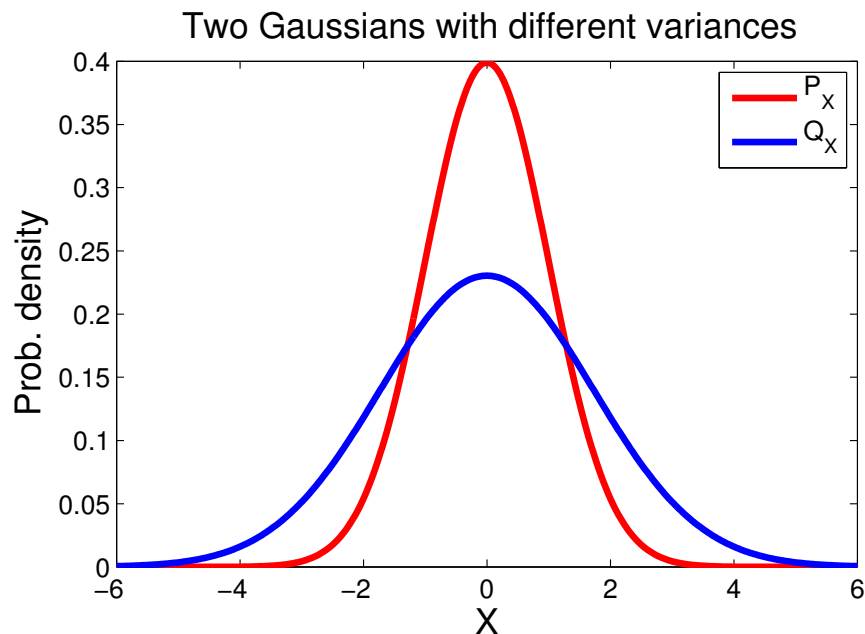
Feature mean difference

- Two Gaussians with same means, different variance
- Idea: look at difference in **means of features** of the RVs
- In Gaussian case: second order features of form $\varphi(x) = x^2$



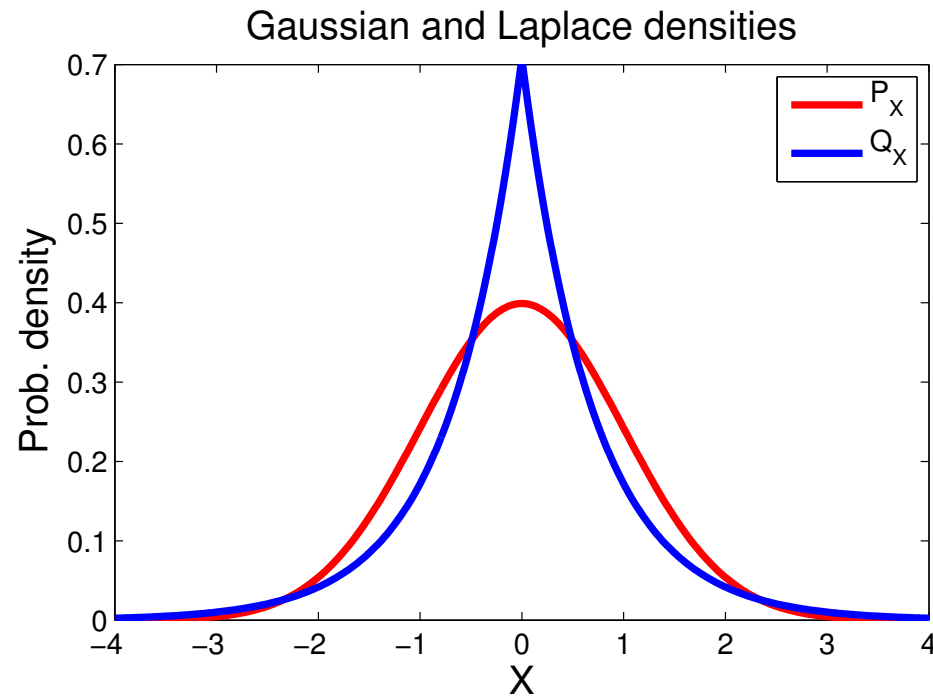
Feature mean difference

- Two Gaussians with same means, different variance
- Idea: look at difference in means of **features** of the RVs
- In Gaussian case: second order features of form $\varphi(x) = x^2$



Feature mean difference

- Gaussian and Laplace distributions
- Same mean *and* same variance
- Difference in means using **higher order features...RKHS**



Probabilities in feature space: the mean trick

The kernel trick

- Given $x \in \mathcal{X}$ for some set \mathcal{X} ,
define **feature map** $\varphi_x \in \mathcal{F}$,

$$\varphi_x = [\dots \varphi_i(x) \dots] \in \ell_2$$

- For **positive definite** $k(x, x')$,

$$k(x, x') = \langle \varphi_x, \varphi_{x'} \rangle_{\mathcal{F}}$$

- The kernel trick: $\forall f \in \mathcal{F}$,

$$f(x) = \langle f, \varphi_x \rangle_{\mathcal{F}}$$

Probabilities in feature space: the mean trick

The kernel trick

- Given $x \in \mathcal{X}$ for some set \mathcal{X} , define feature map $\varphi_x \in \mathcal{F}$,

$$\varphi_x = [\dots \varphi_i(x) \dots] \in \ell_2$$

- For positive definite $k(x, x')$,

$$k(x, x') = \langle \varphi_x, \varphi_{x'} \rangle_{\mathcal{F}}$$

- The kernel trick: $\forall f \in \mathcal{F}$,

$$f(x) = \langle f, \varphi_x \rangle_{\mathcal{F}}$$

The mean trick

- Given \mathbf{P} a Borel probability measure on \mathcal{X} , define feature map $\mu_{\mathbf{P}} \in \mathcal{F}$

$$\mu_{\mathbf{P}} = [\dots \mathbf{E}_{\mathbf{P}} [\varphi_i(\mathbf{x})] \dots]$$

- For positive definite $k(x, x')$,

$$\mathbf{E}_{\mathbf{P}, \mathbf{Q}} k(\mathbf{x}, \mathbf{y}) = \langle \mu_{\mathbf{P}}, \mu_{\mathbf{Q}} \rangle_{\mathcal{F}}$$

for $\mathbf{x} \sim \mathbf{P}$ and $\mathbf{y} \sim \mathbf{Q}$.

- The mean trick: (we call $\mu_{\mathbf{P}}$ a mean/distribution embedding)

$$\mathbf{E}_{\mathbf{P}}(f(\mathbf{x})) =: \langle \mu_{\mathbf{P}}, f \rangle_{\mathcal{F}}$$

What does $\mu_{\mathbf{P}}$ look like?

We plot the function $\mu_{\mathbf{P}}$

- Mean embedding $\mu_{\mathbf{P}} \in \mathcal{F}$

$$\langle \mu_{\mathbf{P}}(\cdot), f(\cdot) \rangle_{\mathcal{F}} = E_{\mathbf{P}} f(x).$$

- What does prob. feature map look like?

$$\begin{aligned} \mu_{\mathbf{P}}(x) &= \langle \mu_{\mathbf{P}}(\cdot), \varphi(x) \rangle_{\mathcal{F}} \\ &= \langle \mu_{\mathbf{P}}(\cdot), k(\cdot, x) \rangle_{\mathcal{F}} = E_{\mathbf{P}} k(x, x). \end{aligned}$$

Expectation of kernel!

- Empirical estimate:

$$\hat{\mu}_{\mathbf{P}}(x) = \frac{1}{m} \sum_{i=1}^m k(x_i, x) \quad x_i \sim \mathbf{P}$$

What does $\mu_{\mathbf{P}}$ look like?

We plot the function $\mu_{\mathbf{P}}$

- Mean embedding $\mu_{\mathbf{P}} \in \mathcal{F}$

$$\langle \mu_{\mathbf{P}}(\cdot), f(\cdot) \rangle_{\mathcal{F}} = E_{\mathbf{P}} f(x).$$

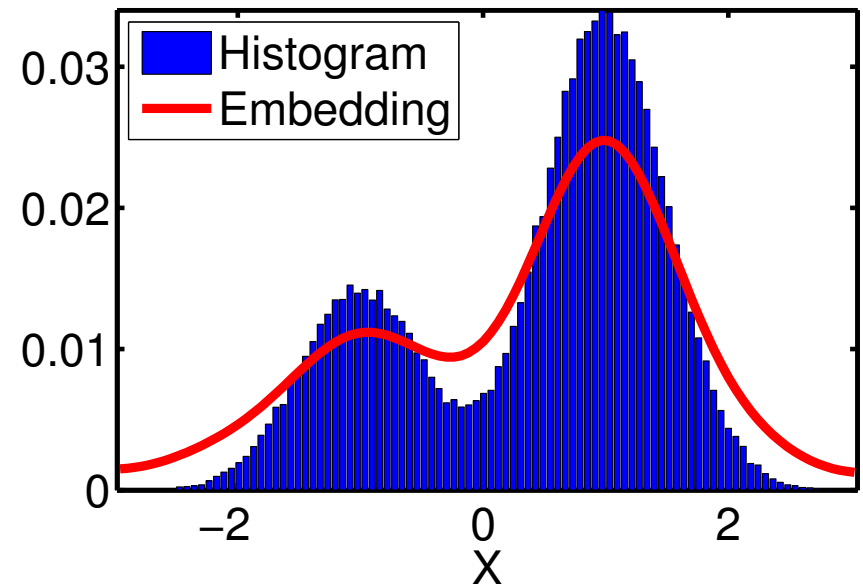
- What does prob. feature map look like?

$$\begin{aligned} \mu_{\mathbf{P}}(x) &= \langle \mu_{\mathbf{P}}(\cdot), \varphi(x) \rangle_{\mathcal{F}} \\ &= \langle \mu_{\mathbf{P}}(\cdot), k(\cdot, x) \rangle_{\mathcal{F}} = E_{\mathbf{P}} k(x, x). \end{aligned}$$

Expectation of kernel!

- Empirical estimate:

$$\hat{\mu}_{\mathbf{P}}(x) = \frac{1}{m} \sum_{i=1}^m k(x_i, x) \quad x_i \sim \mathbf{P}$$



Does the feature space mean exist?

Does there exist an element $\mu_{\mathbf{P}} \in \mathcal{F}$ such that

$$\mathbf{E}_{\mathbf{P}} f(\mathbf{x}) = \mathbf{E}_{\mathbf{P}} \langle f(\cdot), \varphi(\mathbf{x}) \rangle_{\mathcal{F}} = \langle f(\cdot), \mathbf{E}_{\mathbf{P}} \varphi(\mathbf{x}) \rangle_{\mathcal{F}} = \langle f(\cdot), \mu_{\mathbf{P}}(\cdot) \rangle_{\mathcal{F}} \quad \forall f \in \mathcal{F}$$

Does the feature space mean exist?

Does there exist an element $\mu_{\mathbf{P}} \in \mathcal{F}$ such that

$$\mathbf{E}_{\mathbf{P}} f(\mathbf{x}) = \mathbf{E}_{\mathbf{P}} \langle f(\cdot), \varphi(\mathbf{x}) \rangle_{\mathcal{F}} = \langle f(\cdot), \mathbf{E}_{\mathbf{P}} \varphi(\mathbf{x}) \rangle_{\mathcal{F}} = \langle f(\cdot), \mu_{\mathbf{P}}(\cdot) \rangle_{\mathcal{F}} \quad \forall f \in \mathcal{F}$$

Yes: You can exchange expectation and inner product (i.e. $\varphi(\mathbf{x})$ is Bochner integrable [Steinwart and Christmann, 2008]) under the condition

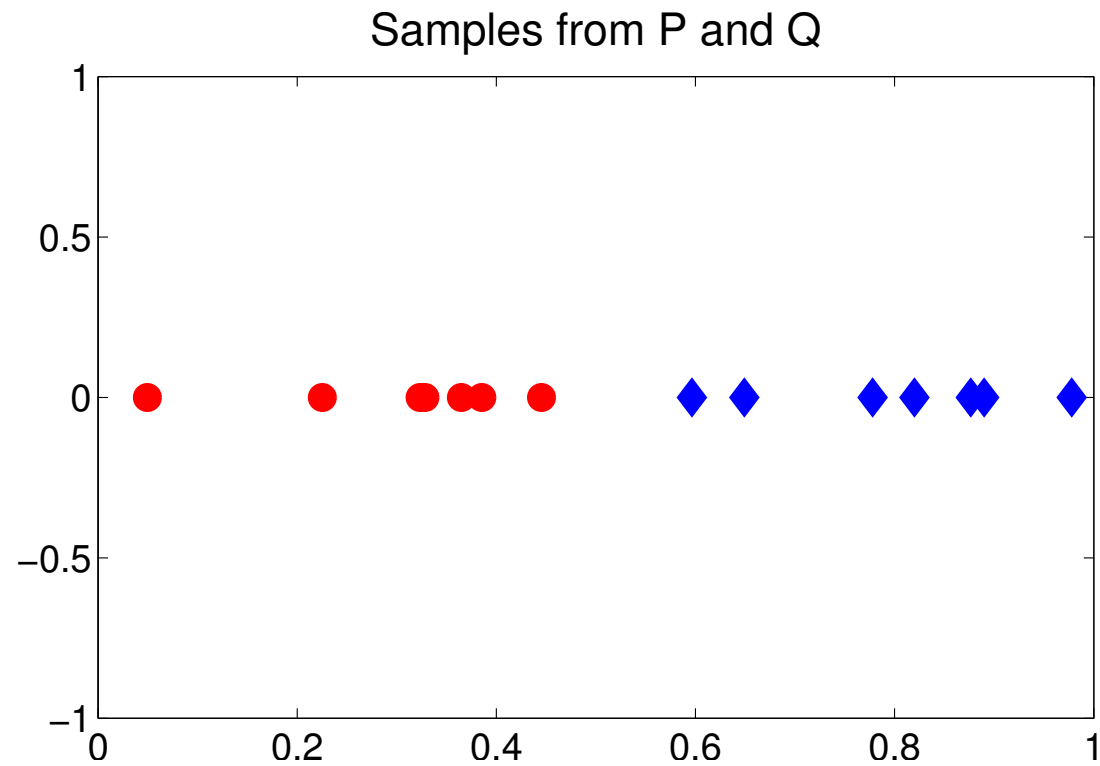
$$\mathbf{E}_{\mathbf{P}} \|\varphi(\mathbf{x})\|_{\mathcal{F}} = \mathbf{E}_{\mathbf{P}} \sqrt{k(\mathbf{x}, \mathbf{x})} < \infty$$

Function Showing Difference in Distributions

- Are **P** and **Q** different?

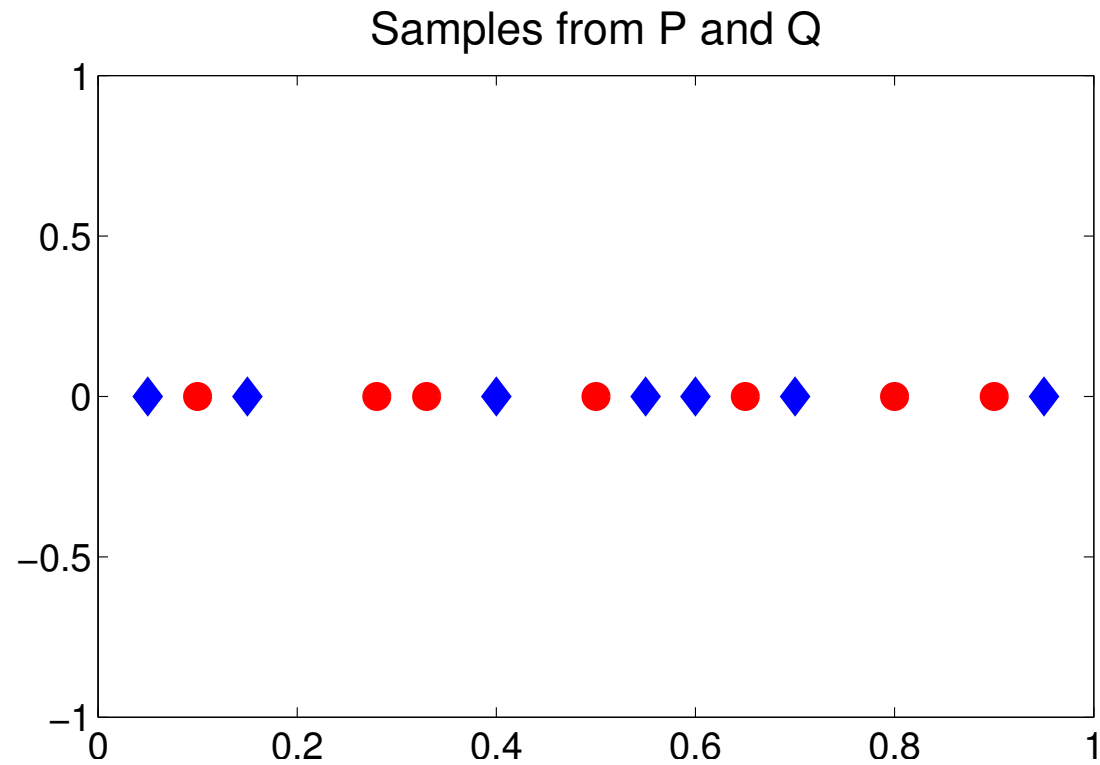
Function Showing Difference in Distributions

- Are **P** and **Q** different?



Function Showing Difference in Distributions

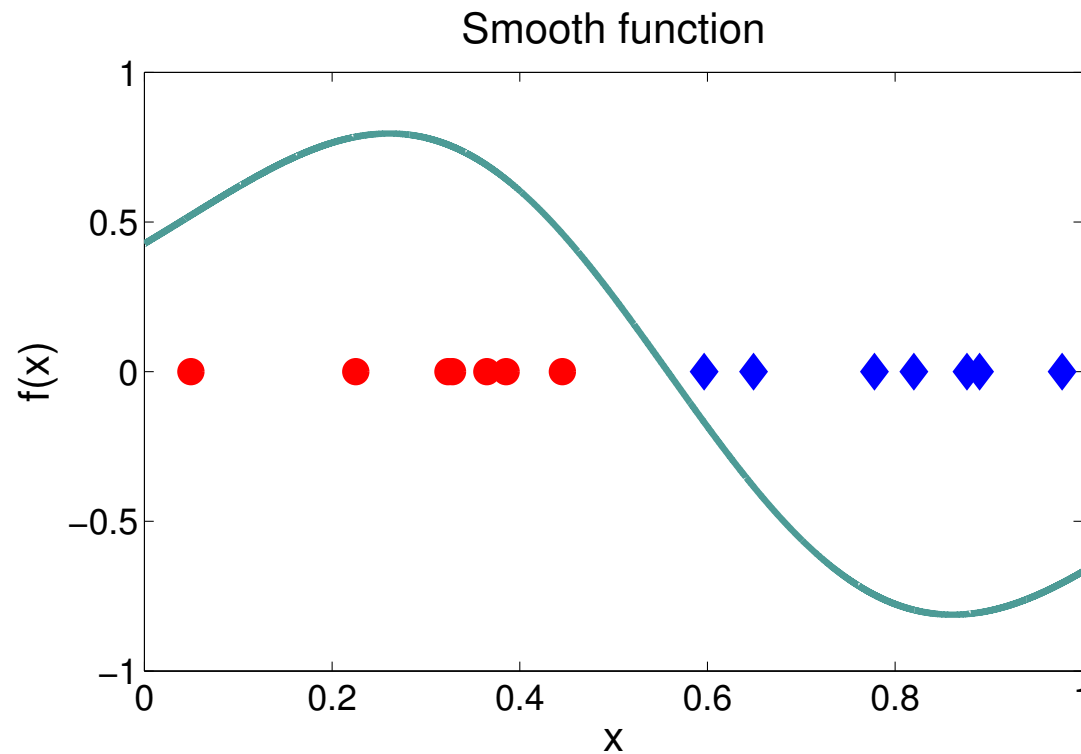
- Are **P** and **Q** different?



Function Showing Difference in Distributions

- **Maximum mean discrepancy**: smooth function for **P** vs **Q**

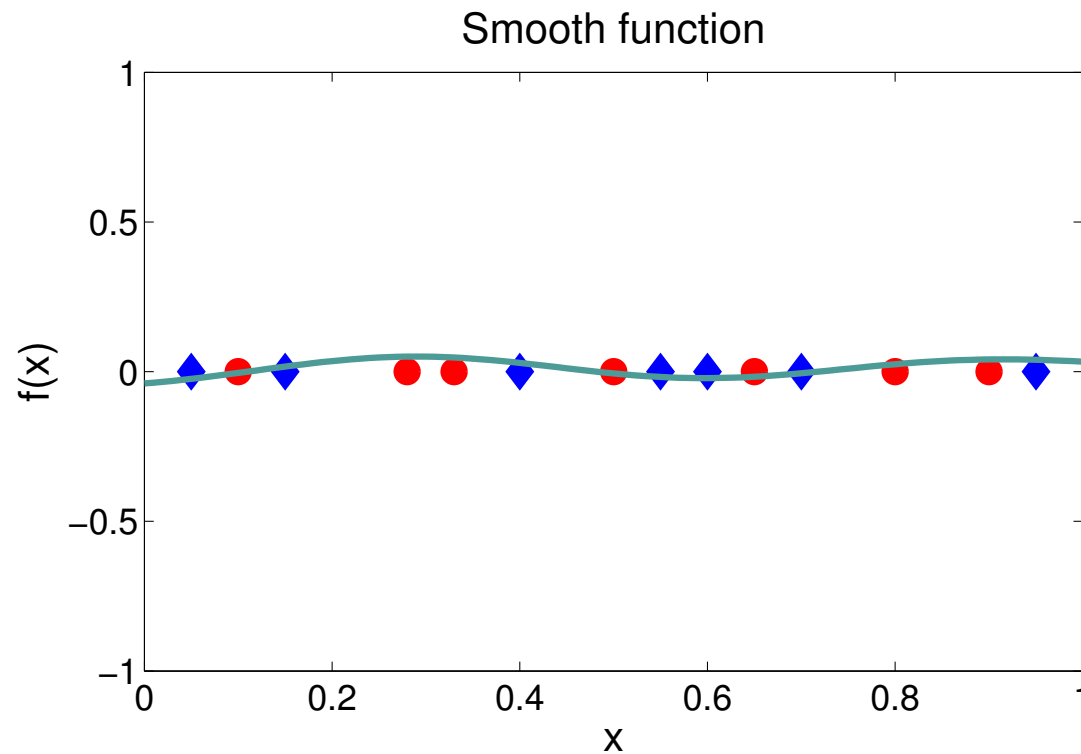
$$\text{MMD}(\mathbf{P}, \mathbf{Q}; F) := \sup_{f \in F} [\mathbf{E}_{\mathbf{P}} f(x) - \mathbf{E}_{\mathbf{Q}} f(y)].$$



Function Showing Difference in Distributions

- Maximum mean discrepancy: smooth function for **P** vs **Q**

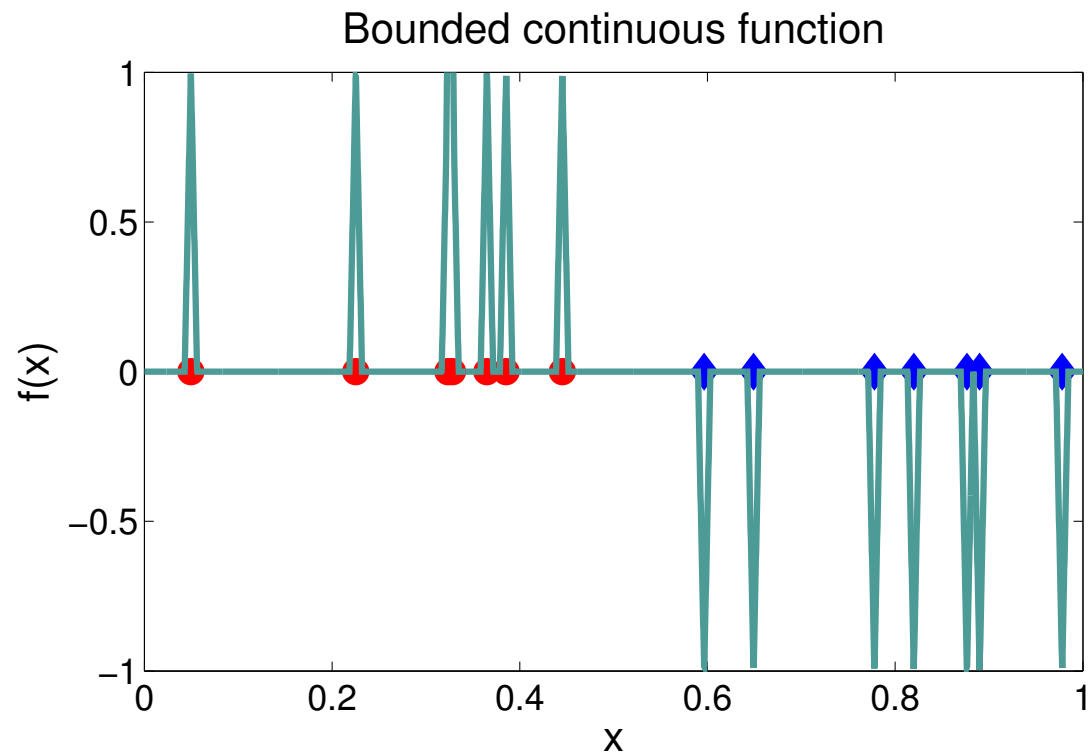
$$\text{MMD}(\mathbf{P}, \mathbf{Q}; F) := \sup_{f \in F} [\mathbf{E}_{\mathbf{P}} f(x) - \mathbf{E}_{\mathbf{Q}} f(y)].$$



Function Showing Difference in Distributions

- What if the function is **not smooth**?

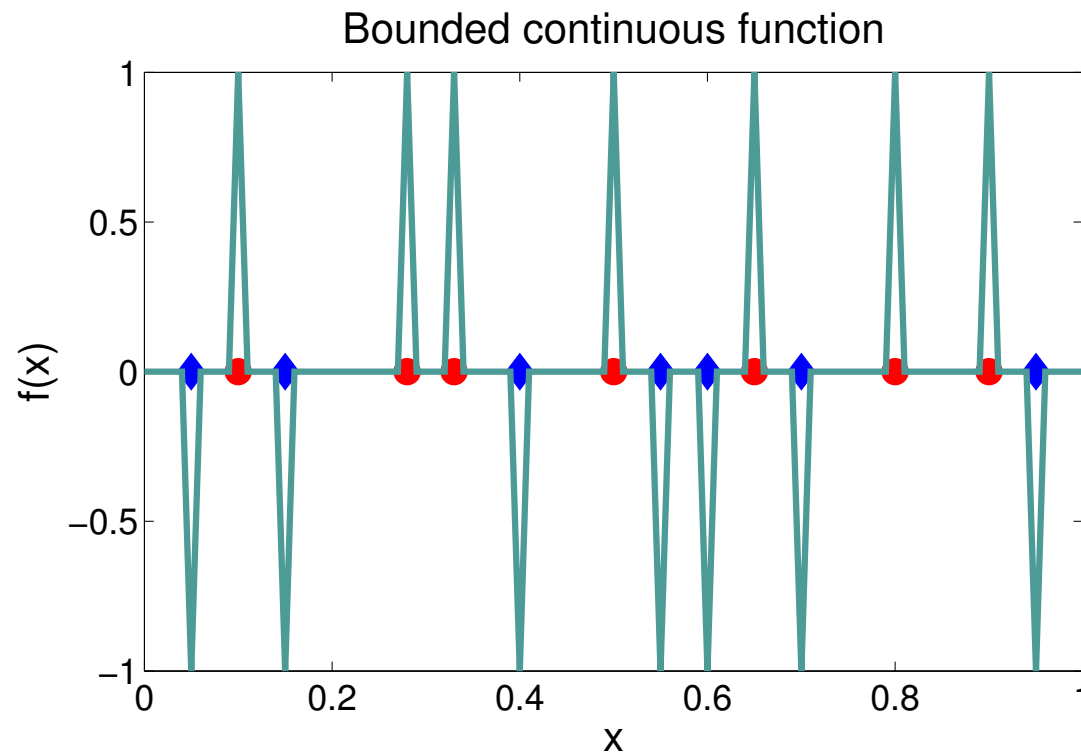
$$\text{MMD}(\mathbf{P}, \mathbf{Q}; F) := \sup_{f \in F} [\mathbf{E}_{\mathbf{P}} f(x) - \mathbf{E}_{\mathbf{Q}} f(y)].$$



Function Showing Difference in Distributions

- What if the function is **not smooth**?

$$\text{MMD}(\mathbf{P}, \mathbf{Q}; F) := \sup_{f \in F} [\mathbf{E}_{\mathbf{P}} f(x) - \mathbf{E}_{\mathbf{Q}} f(y)].$$

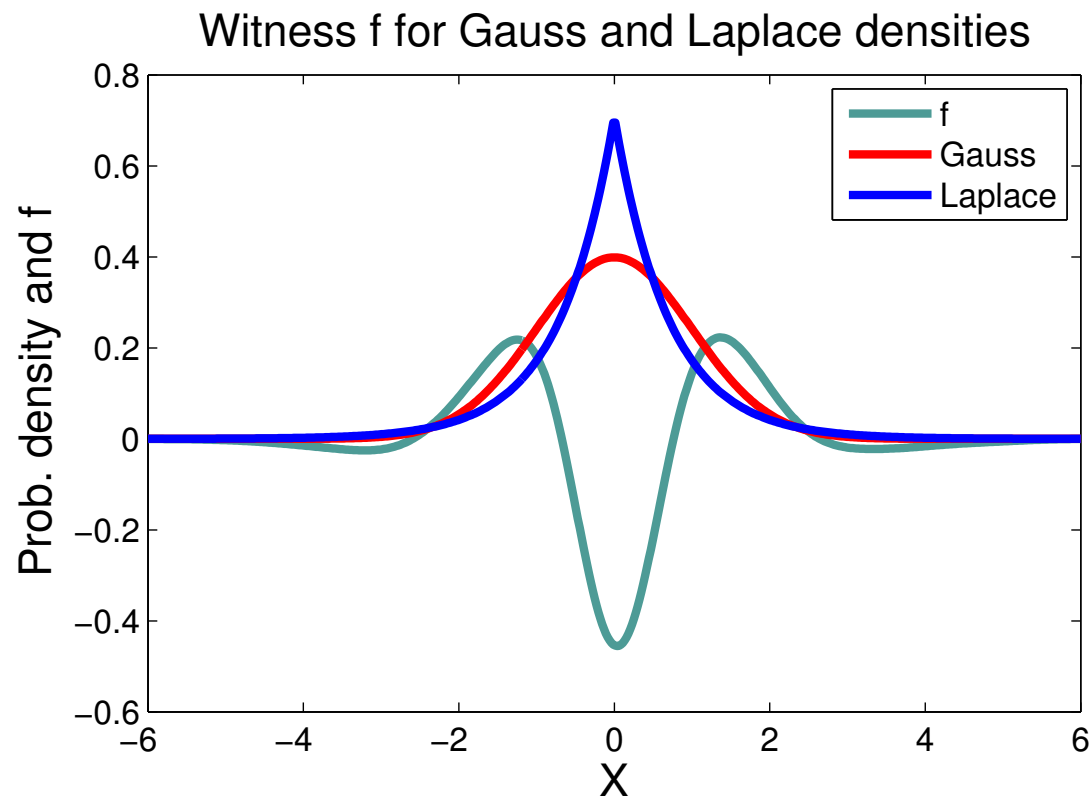


Function Showing Difference in Distributions

- Maximum mean discrepancy: smooth function for **P** vs **Q**

$$\text{MMD}(\mathbf{P}, \mathbf{Q}; F) := \sup_{f \in F} [\mathbf{E}_{\mathbf{P}} f(x) - \mathbf{E}_{\mathbf{Q}} f(y)].$$

- Gauss **P** vs Laplace **Q**



Function Showing Difference in Distributions

- **Maximum mean discrepancy:** smooth function for **P** vs **Q**

$$\text{MMD}(\mathbf{P}, \mathbf{Q}; F) := \sup_{f \in F} [\mathbf{E}_{\mathbf{P}} f(x) - \mathbf{E}_{\mathbf{Q}} f(y)].$$

- **Classical results:** $\text{MMD}(\mathbf{P}, \mathbf{Q}; F) = 0$ iff $\mathbf{P} = \mathbf{Q}$, when
 - $F =$ bounded continuous [Dudley, 2002]
 - $F =$ bounded variation 1 (Kolmogorov metric) [Müller, 1997]
 - $F =$ bounded Lipschitz (Earth mover's distances) [Dudley, 2002]

Function Showing Difference in Distributions

- **Maximum mean discrepancy**: smooth function for **P** vs **Q**

$$\text{MMD}(\mathbf{P}, \mathbf{Q}; F) := \sup_{f \in F} [\mathbf{E}_{\mathbf{P}} f(x) - \mathbf{E}_{\mathbf{Q}} f(y)].$$

- **Classical results**: $\text{MMD}(\mathbf{P}, \mathbf{Q}; F) = 0$ iff $\mathbf{P} = \mathbf{Q}$, when
 - $F =$ bounded continuous [Dudley, 2002]
 - $F =$ bounded variation 1 (Kolmogorov metric) [Müller, 1997]
 - $F =$ bounded Lipschitz (Earth mover's distances) [Dudley, 2002]
- $\text{MMD}(\mathbf{P}, \mathbf{Q}; F) = 0$ iff $\mathbf{P} = \mathbf{Q}$ when $F =$ the **unit ball** in a **characteristic RKHS** \mathcal{F} (coming soon!) [ISMB06, NIPS06a, NIPS07b, NIPS08a, JMLR10]

Function Showing Difference in Distributions

- **Maximum mean discrepancy**: smooth function for **P** vs **Q**

$$\text{MMD}(\mathbf{P}, \mathbf{Q}; F) := \sup_{f \in F} [\mathbf{E}_{\mathbf{P}} f(x) - \mathbf{E}_{\mathbf{Q}} f(y)].$$

- **Classical results**: $\text{MMD}(\mathbf{P}, \mathbf{Q}; F) = 0$ iff $\mathbf{P} = \mathbf{Q}$, when
 - $F =$ bounded continuous [Dudley, 2002]
 - $F =$ bounded variation 1 (Kolmogorov metric) [Müller, 1997]
 - $F =$ bounded Lipschitz (Earth mover's distances) [Dudley, 2002]
- $\text{MMD}(\mathbf{P}, \mathbf{Q}; F) = 0$ iff $\mathbf{P} = \mathbf{Q}$ when $F =$ the **unit ball** in a **characteristic RKHS** \mathcal{F} (coming soon!) [ISMB06, NIPS06a, NIPS07b, NIPS08a, JMLR10]

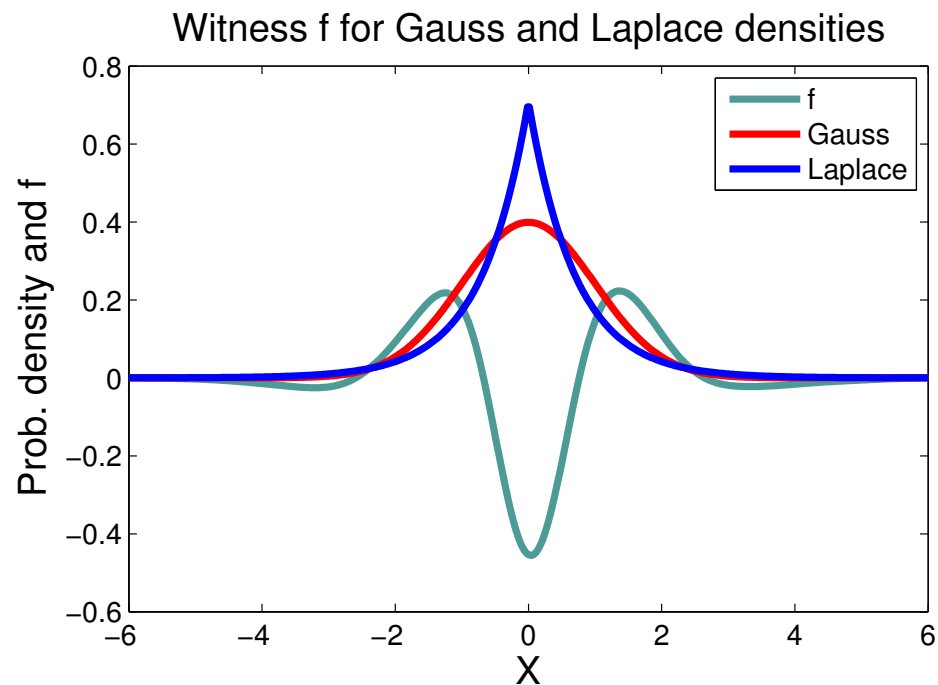
How do **smooth functions** relate to **feature maps**?

Function view vs feature mean view

- The (kernel) MMD: [ISMB06, NIPS06a]

$$\text{MMD}^2(\mathbf{P}, \mathbf{Q}; F)$$

$$= \left(\sup_{f \in F} [\mathbf{E}_{\mathbf{P}} f(x) - \mathbf{E}_{\mathbf{Q}} f(y)] \right)^2$$



Function view vs feature mean view

- The (kernel) MMD: [ISMB06, NIPS06a]

$$\text{MMD}^2(\mathbf{P}, \mathbf{Q}; F)$$

$$= \left(\sup_{f \in F} [\mathbf{E}_{\mathbf{P}} f(x) - \mathbf{E}_{\mathbf{Q}} f(y)] \right)^2$$

use

$$\mathbf{E}_{\mathbf{P}}(f(x)) =: \langle \mu_{\mathbf{P}}, f \rangle_{\mathcal{F}}$$

Function view vs feature mean view

- The (kernel) MMD: [ISMB06, NIPS06a]

$$\text{MMD}^2(\mathbf{P}, \mathbf{Q}; F)$$

$$= \left(\sup_{f \in F} [\mathbf{E}_{\mathbf{P}} f(\mathbf{x}) - \mathbf{E}_{\mathbf{Q}} f(\mathbf{y})] \right)^2$$

use

$$\mathbf{E}_{\mathbf{P}}(f(\mathbf{x})) =: \langle \mu_{\mathbf{P}}, f \rangle_{\mathcal{F}}$$

$$= \left(\sup_{f \in F} \langle f, \mu_{\mathbf{P}} - \mu_{\mathbf{Q}} \rangle_{\mathcal{F}} \right)^2$$

Function view vs feature mean view

- The (kernel) MMD: [ISMB06, NIPS06a]

$$\text{MMD}^2(\mathbf{P}, \mathbf{Q}; F)$$

$$= \left(\sup_{f \in F} [\mathbf{E}_{\mathbf{P}} f(\mathbf{x}) - \mathbf{E}_{\mathbf{Q}} f(\mathbf{y})] \right)^2$$

use

$$= \left(\sup_{f \in F} \langle f, \mu_{\mathbf{P}} - \mu_{\mathbf{Q}} \rangle_{\mathcal{F}} \right)^2$$

$$\|\theta\|_{\mathcal{F}} = \sup_{f \in F} \langle f, \theta \rangle_{\mathcal{F}}$$

$$= \|\mu_{\mathbf{P}} - \mu_{\mathbf{Q}}\|_{\mathcal{F}}^2$$

Function view and feature view **equivalent**

Empirical estimate of MMD

- An unbiased empirical estimate: for $\{x_i\}_{i=1}^m \sim \mathbf{P}$ and $\{y_i\}_{i=1}^m \sim \mathbf{Q}$,

$$\widehat{MMD}^2 = \frac{1}{m(m-1)} \sum_{i=1}^m \sum_{j \neq i}^m [k(x_i, x_j) + k(y_i, y_j)] \\ - \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m [k(y_i, x_j) + k(x_i, y_j)]$$

Empirical estimate of MMD

- An unbiased **empirical estimate**: for $\{x_i\}_{i=1}^m \sim \mathbf{P}$ and $\{y_i\}_{i=1}^m \sim \mathbf{Q}$,

$$\widehat{MMD}^2 = \frac{1}{m(m-1)} \sum_{i=1}^m \sum_{j \neq i}^m [k(x_i, x_j) + k(y_i, y_j)] \\ - \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m [k(y_i, x_j) + k(x_i, y_j)]$$

- **Proof:**

$$\|\mu_{\mathbf{P}} - \mu_{\mathbf{Q}}\|_{\mathcal{F}}^2 = \langle \mu_{\mathbf{P}} - \mu_{\mathbf{Q}}, \mu_{\mathbf{P}} - \mu_{\mathbf{Q}} \rangle_{\mathcal{F}} \\ = \langle \mu_{\mathbf{P}}, \mu_{\mathbf{P}} \rangle + \langle \mu_{\mathbf{Q}}, \mu_{\mathbf{Q}} \rangle - 2 \langle \mu_{\mathbf{P}}, \mu_{\mathbf{Q}} \rangle$$

Empirical estimate of MMD

- An unbiased empirical estimate: for $\{x_i\}_{i=1}^m \sim \mathbf{P}$ and $\{y_i\}_{i=1}^m \sim \mathbf{Q}$,

$$\widehat{MMD}^2 = \frac{1}{m(m-1)} \sum_{i=1}^m \sum_{j \neq i}^m [k(x_i, x_j) + k(y_i, y_j)] \\ - \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m [k(y_i, x_j) + k(x_i, y_j)]$$

- Proof:

$$\|\mu_{\mathbf{P}} - \mu_{\mathbf{Q}}\|_{\mathcal{F}}^2 = \langle \mu_{\mathbf{P}} - \mu_{\mathbf{Q}}, \mu_{\mathbf{P}} - \mu_{\mathbf{Q}} \rangle_{\mathcal{F}} \\ = \langle \mu_{\mathbf{P}}, \mu_{\mathbf{P}} \rangle + \langle \mu_{\mathbf{Q}}, \mu_{\mathbf{Q}} \rangle - 2 \langle \mu_{\mathbf{P}}, \mu_{\mathbf{Q}} \rangle$$

Empirical estimate of MMD

- An unbiased empirical estimate: for $\{x_i\}_{i=1}^m \sim \mathbf{P}$ and $\{y_i\}_{i=1}^m \sim \mathbf{Q}$,

$$\begin{aligned}\widehat{MMD}^2 &= \frac{1}{m(m-1)} \sum_{i=1}^m \sum_{j \neq i}^m [k(x_i, x_j) + k(y_i, y_j)] \\ &\quad - \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m [k(y_i, x_j) + k(x_i, y_j)]\end{aligned}$$

- Proof:

$$\begin{aligned}\|\mu_{\mathbf{P}} - \mu_{\mathbf{Q}}\|_{\mathcal{F}}^2 &= \langle \mu_{\mathbf{P}} - \mu_{\mathbf{Q}}, \mu_{\mathbf{P}} - \mu_{\mathbf{Q}} \rangle_{\mathcal{F}} \\ &= \langle \mu_{\mathbf{P}}, \mu_{\mathbf{P}} \rangle + \langle \mu_{\mathbf{Q}}, \mu_{\mathbf{Q}} \rangle - 2 \langle \mu_{\mathbf{P}}, \mu_{\mathbf{Q}} \rangle \\ &= \mathbf{E}_{\mathbf{P}}[\mu_{\mathbf{P}}(x)] + \dots\end{aligned}$$

Empirical estimate of MMD

- An unbiased **empirical estimate**: for $\{x_i\}_{i=1}^m \sim \mathbf{P}$ and $\{y_i\}_{i=1}^m \sim \mathbf{Q}$,

$$\begin{aligned}\widehat{MMD}^2 &= \frac{1}{m(m-1)} \sum_{i=1}^m \sum_{j \neq i}^m [k(x_i, x_j) + k(y_i, y_j)] \\ &\quad - \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m [k(y_i, x_j) + k(x_i, y_j)]\end{aligned}$$

- **Proof:**

$$\begin{aligned}\|\mu_{\mathbf{P}} - \mu_{\mathbf{Q}}\|_{\mathcal{F}}^2 &= \langle \mu_{\mathbf{P}} - \mu_{\mathbf{Q}}, \mu_{\mathbf{P}} - \mu_{\mathbf{Q}} \rangle_{\mathcal{F}} \\ &= \langle \mu_{\mathbf{P}}, \mu_{\mathbf{P}} \rangle + \langle \mu_{\mathbf{Q}}, \mu_{\mathbf{Q}} \rangle - 2 \langle \mu_{\mathbf{P}}, \mu_{\mathbf{Q}} \rangle \\ &= \mathbf{E}_{\mathbf{P}}[\mu_{\mathbf{P}}(\mathbf{x})] + \dots \\ &= \mathbf{E}_{\mathbf{P}} \langle \mu_{\mathbf{P}}(\cdot), \varphi(\mathbf{x}) \rangle + \dots\end{aligned}$$

Empirical estimate of MMD

- An unbiased **empirical estimate**: for $\{x_i\}_{i=1}^m \sim \mathbf{P}$ and $\{y_i\}_{i=1}^m \sim \mathbf{Q}$,

$$\begin{aligned}\widehat{MMD}^2 &= \frac{1}{m(m-1)} \sum_{i=1}^m \sum_{j \neq i}^m [k(x_i, x_j) + k(y_i, y_j)] \\ &\quad - \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m [k(y_i, x_j) + k(x_i, y_j)]\end{aligned}$$

- **Proof:**

$$\begin{aligned}\|\mu_{\mathbf{P}} - \mu_{\mathbf{Q}}\|_{\mathcal{F}}^2 &= \langle \mu_{\mathbf{P}} - \mu_{\mathbf{Q}}, \mu_{\mathbf{P}} - \mu_{\mathbf{Q}} \rangle_{\mathcal{F}} \\ &= \langle \mu_{\mathbf{P}}, \mu_{\mathbf{P}} \rangle + \langle \mu_{\mathbf{Q}}, \mu_{\mathbf{Q}} \rangle - 2 \langle \mu_{\mathbf{P}}, \mu_{\mathbf{Q}} \rangle \\ &= \mathbf{E}_{\mathbf{P}}[\mu_{\mathbf{P}}(x)] + \dots \\ &= \mathbf{E}_{\mathbf{P}} \langle \mu_{\mathbf{P}}(\cdot), k(x, \cdot) \rangle + \dots\end{aligned}$$

Empirical estimate of MMD

- An unbiased **empirical estimate**: for $\{x_i\}_{i=1}^m \sim \mathbf{P}$ and $\{y_i\}_{i=1}^m \sim \mathbf{Q}$,

$$\widehat{MMD}^2 = \frac{1}{m(m-1)} \sum_{i=1}^m \sum_{j \neq i}^m [k(x_i, x_j) + k(y_i, y_j)] \\ - \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m [k(y_i, x_j) + k(x_i, y_j)]$$

- **Proof:**

$$\begin{aligned} \|\mu_{\mathbf{P}} - \mu_{\mathbf{Q}}\|_{\mathcal{F}}^2 &= \langle \mu_{\mathbf{P}} - \mu_{\mathbf{Q}}, \mu_{\mathbf{P}} - \mu_{\mathbf{Q}} \rangle_{\mathcal{F}} \\ &= \langle \mu_{\mathbf{P}}, \mu_{\mathbf{P}} \rangle + \langle \mu_{\mathbf{Q}}, \mu_{\mathbf{Q}} \rangle - 2 \langle \mu_{\mathbf{P}}, \mu_{\mathbf{Q}} \rangle \\ &= \mathbf{E}_{\mathbf{P}}[\mu_{\mathbf{P}}(x)] + \dots \\ &= \mathbf{E}_{\mathbf{P}} \langle \mu_{\mathbf{P}}(\cdot), k(x, \cdot) \rangle + \dots \\ &= \mathbf{E}_{\mathbf{P}} k(x, x') + \mathbf{E}_{\mathbf{Q}} k(y, y') - 2\mathbf{E}_{\mathbf{P}, \mathbf{Q}} k(x, y) \end{aligned}$$

Empirical estimate of MMD

- An unbiased **empirical estimate**: for $\{x_i\}_{i=1}^m \sim \mathbf{P}$ and $\{y_i\}_{i=1}^m \sim \mathbf{Q}$,

$$\begin{aligned}\widehat{MMD}^2 &= \frac{1}{m(m-1)} \sum_{i=1}^m \sum_{j \neq i}^m [k(x_i, x_j) + k(y_i, y_j)] \\ &\quad - \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m [k(y_i, x_j) + k(x_i, y_j)]\end{aligned}$$

- **Proof:**

$$\begin{aligned}\|\mu_{\mathbf{P}} - \mu_{\mathbf{Q}}\|_{\mathcal{F}}^2 &= \langle \mu_{\mathbf{P}} - \mu_{\mathbf{Q}}, \mu_{\mathbf{P}} - \mu_{\mathbf{Q}} \rangle_{\mathcal{F}} \\ &= \langle \mu_{\mathbf{P}}, \mu_{\mathbf{P}} \rangle + \langle \mu_{\mathbf{Q}}, \mu_{\mathbf{Q}} \rangle - 2 \langle \mu_{\mathbf{P}}, \mu_{\mathbf{Q}} \rangle \\ &= \mathbf{E}_{\mathbf{P}}[\mu_{\mathbf{P}}(x)] + \dots \\ &= \mathbf{E}_{\mathbf{P}} \langle \mu_{\mathbf{P}}(\cdot), k(x, \cdot) \rangle + \dots \\ &= \mathbf{E}_{\mathbf{P}} k(x, x') + \mathbf{E}_{\mathbf{Q}} k(y, y') - 2\mathbf{E}_{\mathbf{P}, \mathbf{Q}} k(x, y)\end{aligned}$$

Then $\widehat{\mathbf{E}}k(x, x') = \frac{1}{m(m-1)} \sum_{i=1}^m \sum_{j \neq i}^m k(x_i, x_j)$

MMD for independence: HSIC

- Dependence measure: the **Hilbert Schmidt Independence Criterion** [ALT05, NIPS07a, ALT07, ALT08, JMLR10]

Related to [Feuerverger, 1993] and [Székely and Rizzo, 2009, Székely et al., 2007]

$$HSIC(\mathbf{P}_{XY}, \mathbf{P}_X \mathbf{P}_Y) := \|\mu_{\mathbf{P}_{XY}} - \mu_{\mathbf{P}_X \mathbf{P}_Y}\|^2$$

MMD for independence: HSIC

- Dependence measure: the **Hilbert Schmidt Independence Criterion** [ALT05, NIPS07a, ALT07, ALT08, JMLR10]

Related to [Feuerverger, 1993] and [Székely and Rizzo, 2009, Székely et al., 2007]

$$HSIC(\mathbf{P}_{XY}, \mathbf{P}_X \mathbf{P}_Y) := \|\mu_{\mathbf{P}_{XY}} - \mu_{\mathbf{P}_X \mathbf{P}_Y}\|^2$$

$$\begin{aligned} & k\left(\begin{array}{|c|} \hline \text{1} \\ \hline \end{array}, \begin{array}{|c|} \hline \text{2} \\ \hline \end{array}\right) \quad l\left(\begin{array}{|c|} \hline \text{1} \\ \hline \end{array}, \begin{array}{|c|} \hline \text{2} \\ \hline \end{array}\right) \\ & \quad \quad \quad \downarrow \\ & \mathcal{K}\left(\begin{array}{|c|} \hline \text{1} \\ \hline \end{array} \begin{array}{|c|} \hline \text{1} \\ \hline \end{array}, \begin{array}{|c|} \hline \text{2} \\ \hline \end{array} \begin{array}{|c|} \hline \text{2} \\ \hline \end{array}\right) = \\ & k\left(\begin{array}{|c|} \hline \text{1} \\ \hline \end{array}, \begin{array}{|c|} \hline \text{2} \\ \hline \end{array}\right) \times l\left(\begin{array}{|c|} \hline \text{1} \\ \hline \end{array}, \begin{array}{|c|} \hline \text{2} \\ \hline \end{array}\right) \end{aligned}$$

MMD for independence: HSIC

- Dependence measure: the **Hilbert Schmidt Independence Criterion** [ALT05, NIPS07a, ALT07, ALT08, JMLR10]

Related to [Feuerverger, 1993] and [Székely and Rizzo, 2009, Székely et al., 2007]

$$HSIC(\mathbf{P}_{XY}, \mathbf{P}_X \mathbf{P}_Y) := \|\mu_{\mathbf{P}_{XY}} - \mu_{\mathbf{P}_X \mathbf{P}_Y}\|^2$$

HSIC using expectations of kernels:

Define RKHS \mathcal{F} on \mathcal{X} with kernel k , RKHS \mathcal{G} on \mathcal{Y} with kernel l . Then

$$\begin{aligned} HSIC(\mathbf{P}_{XY}, \mathbf{P}_X \mathbf{P}_Y) &= \mathbf{E}_{XY} \mathbf{E}_{X'Y'} k(x, x') l(y, y') + \mathbf{E}_X \mathbf{E}_{X'} k(x, x') \mathbf{E}_Y \mathbf{E}_{Y'} l(y, y') \\ &\quad - 2 \mathbf{E}_{X'Y'} [\mathbf{E}_X k(x, x') \mathbf{E}_Y l(y, y')] . \end{aligned}$$

HSIC: empirical estimate and intuition



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.

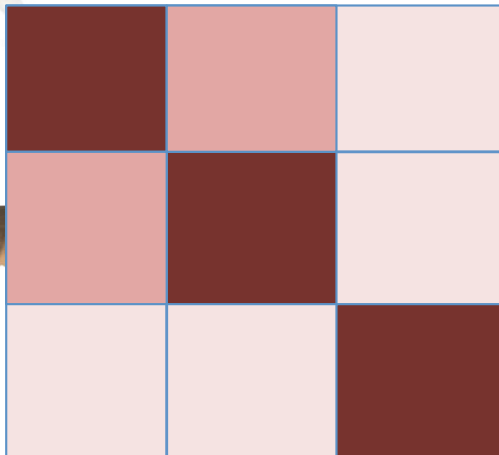


Known for their curiosity, intelligence, and excellent communication skills, the Javanese breed is perfect if you want a responsive, interactive pet, one that will blow in your ear and follow you everywhere.

HSIC: empirical estimate and intuition

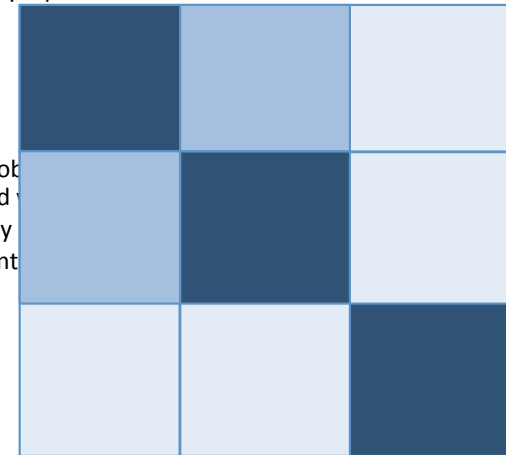


K



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.

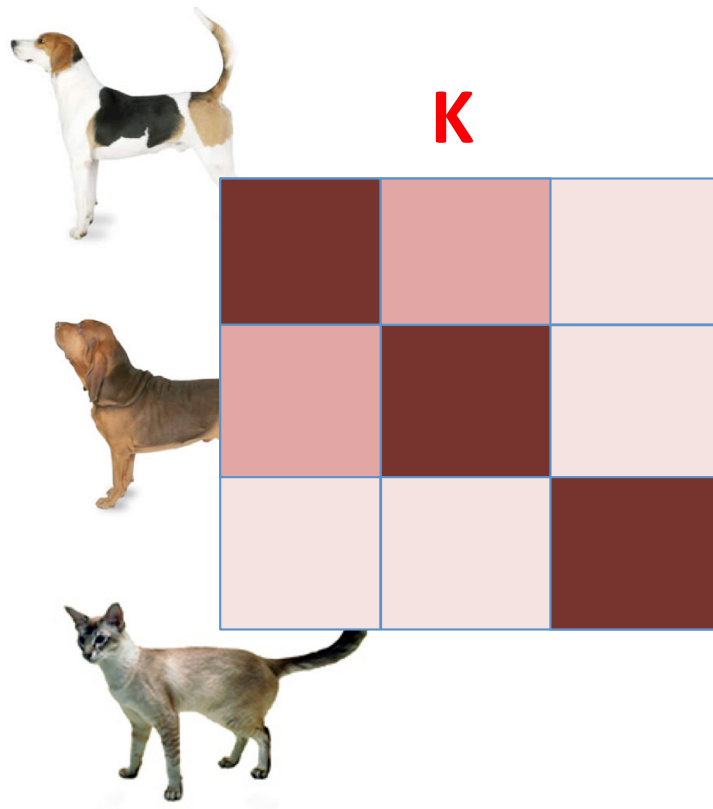
L



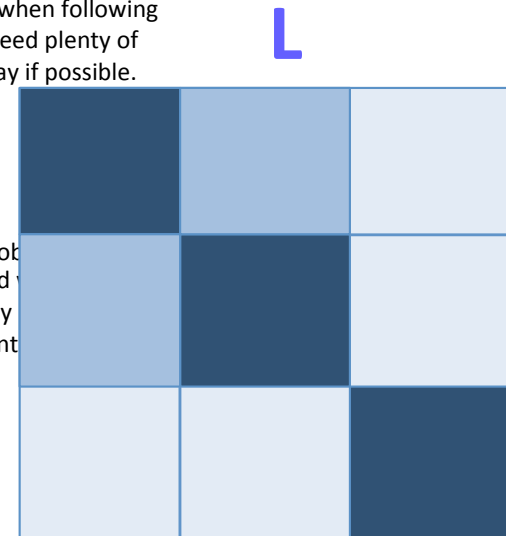
A large animal who slings slobbery, distinctive houndy odor, and more than to follow his nose. They need a lot of exercise and mental stimulation.

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.

HSIC: empirical estimate and intuition



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 slobbery, distinctive houndy odor, and is more than willing to follow his nose. They need a lot of exercise and mental stimulation.

Known for their curiosity, intelligence, and excellent communication skills, the Javanese breed is perfect if you want a responsive, interactive pet, one that will blow in your ear and follow you everywhere.

Text from dogtime.com and petfinder.com

Empirical $HSIC(\mathbf{P}_{XY}, \mathbf{P}_X \mathbf{P}_Y)$:

$$\frac{1}{n^2} (H K H \circ H L H)_{++}$$

Characteristic kernels (Via Fourier, on the torus \mathbb{T})

Characteristic Kernels (via Fourier)

Reminder:

Characteristic: MMD a **metric** (MMD = 0 iff **P** = **Q**) [NIPS07b, JMLR10]

In the next slides:

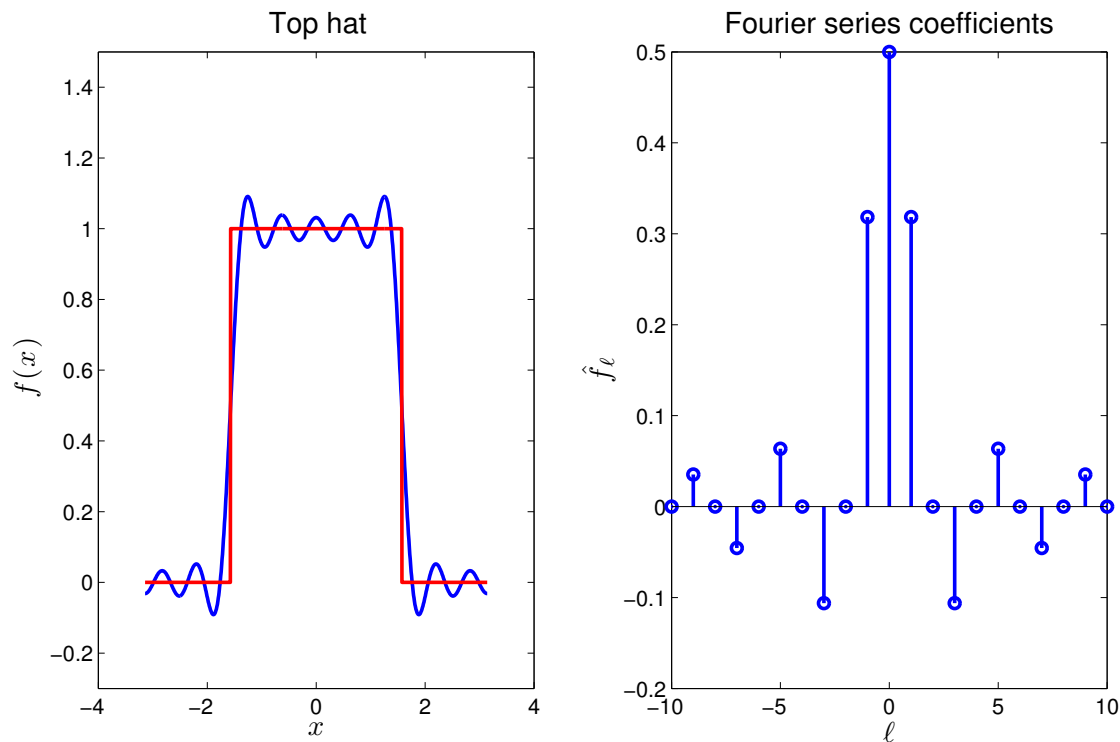
1. Characteristic property on $[-\pi, \pi]$ with periodic boundary
2. Characteristic property on \mathbb{R}^d

Characteristic Kernels (via Fourier)

Reminder: **Fourier series**

- Function $[-\pi, \pi]$ with periodic boundary.

$$f(x) = \sum_{l=-\infty}^{\infty} \hat{f}_l \exp(\imath l x) = \sum_{l=-\infty}^{\infty} \hat{f}_l (\cos(lx) + \imath \sin(lx)).$$



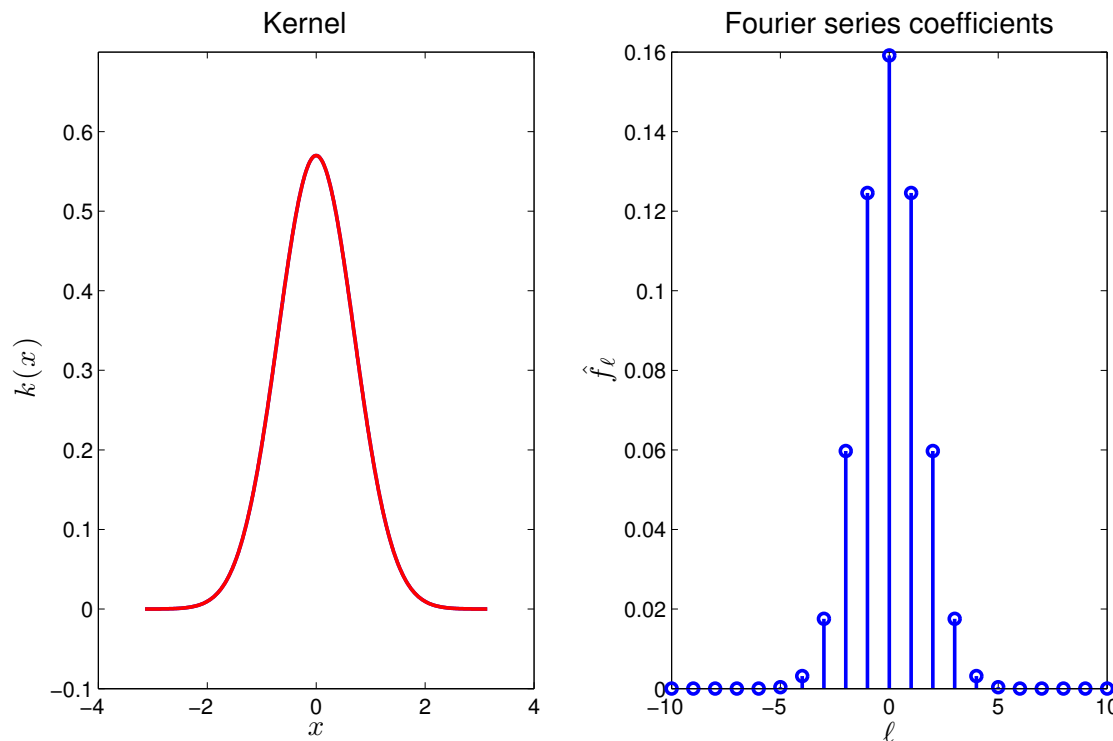
Characteristic Kernels (via Fourier)

Reminder: **Fourier series of kernel**

$$k(x, y) = k(x - y) = k(z), \quad k(z) = \sum_{\ell=-\infty}^{\infty} \hat{k}_{\ell} \exp(i\ell z),$$

E.g., $k(x) = \frac{1}{2\pi} \vartheta\left(\frac{x}{2\pi}, \frac{i\sigma^2}{2\pi}\right), \quad \hat{k}_{\ell} = \frac{1}{2\pi} \exp\left(\frac{-\sigma^2 \ell^2}{2}\right).$

ϑ is the Jacobi theta function, close to Gaussian when σ^2 sufficiently narrower than $[-\pi, \pi]$.



Characteristic Kernels (via Fourier)

Maximum mean embedding via Fourier series:

- Fourier series for \mathbf{P} is **characteristic function** $\bar{\phi}_{\mathbf{P}}$
- Fourier series for mean embedding is **product** of Fourier series!
(convolution theorem)

$$\mu_{\mathbf{P}}(x) = E_{\mathbf{P}}k(x - x) = \int_{-\pi}^{\pi} k(x - t)d\mathbf{P}(t) \quad \hat{\mu}_{\mathbf{P},\ell} = \hat{k}_{\ell} \times \bar{\phi}_{\mathbf{P},\ell}$$

Characteristic Kernels (via Fourier)

Maximum mean embedding via Fourier series:

- Fourier series for \mathbf{P} is **characteristic function** $\bar{\phi}_{\mathbf{P}}$
- Fourier series for mean embedding is **product** of Fourier series!
(convolution theorem)

$$\mu_{\mathbf{P}}(x) = E_{\mathbf{P}}k(x - x) = \int_{-\pi}^{\pi} k(x - t)d\mathbf{P}(t) \quad \hat{\mu}_{\mathbf{P},\ell} = \hat{k}_{\ell} \times \bar{\phi}_{\mathbf{P},\ell}$$

- MMD can be written in terms of Fourier series:

$$\text{MMD}(\mathbf{P}, \mathbf{Q}; F) := \left\| \sum_{\ell=-\infty}^{\infty} [(\bar{\phi}_{\mathbf{P},\ell} - \bar{\phi}_{\mathbf{Q},\ell}) \hat{k}_{\ell}] \exp(i\ell x) \right\|_{\mathcal{F}}$$

A simpler Fourier expression for MMD

- From previous slide,

$$\text{MMD}(\mathbf{P}, \mathbf{Q}; F) := \left\| \sum_{l=-\infty}^{\infty} [(\bar{\phi}_{\mathbf{P},l} - \bar{\phi}_{\mathbf{Q},l}) \hat{k}_l] \exp(\imath lx) \right\|_{\mathcal{F}}$$

- The squared norm of a function f in \mathcal{F} is:

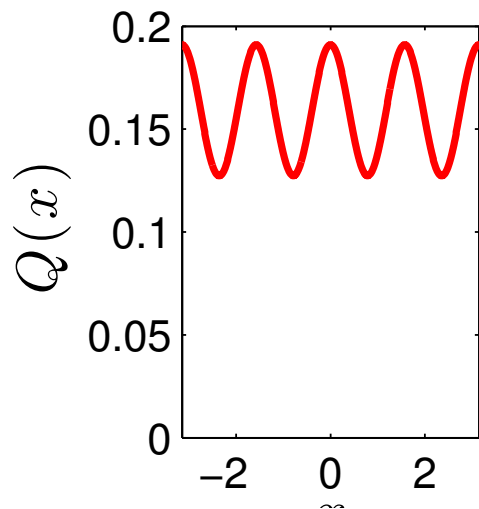
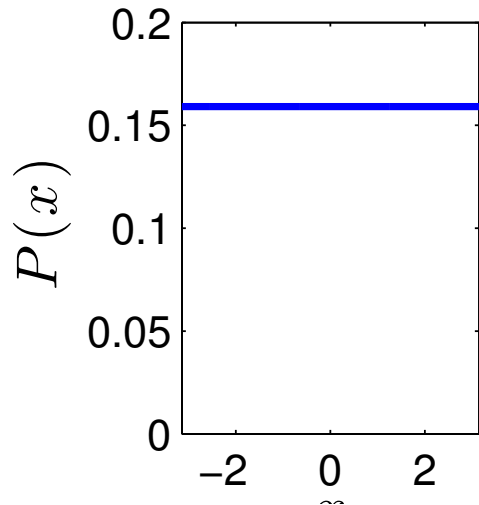
$$\|f\|_{\mathcal{F}}^2 = \langle f, f \rangle_{\mathcal{F}} = \sum_{l=-\infty}^{\infty} \frac{|\hat{f}_l|^2}{\hat{k}_l}.$$

- Simple, interpretable expression for squared MMD:

$$\text{MMD}^2(\mathbf{P}, \mathbf{Q}; F) = \sum_{l=-\infty}^{\infty} \frac{[|\phi_{\mathbf{P},l} - \phi_{\mathbf{Q},l}|^2 \hat{k}_l]^2}{\hat{k}_l} = \sum_{l=-\infty}^{\infty} |\phi_{\mathbf{P},l} - \phi_{\mathbf{Q},l}|^2 \hat{k}_l$$

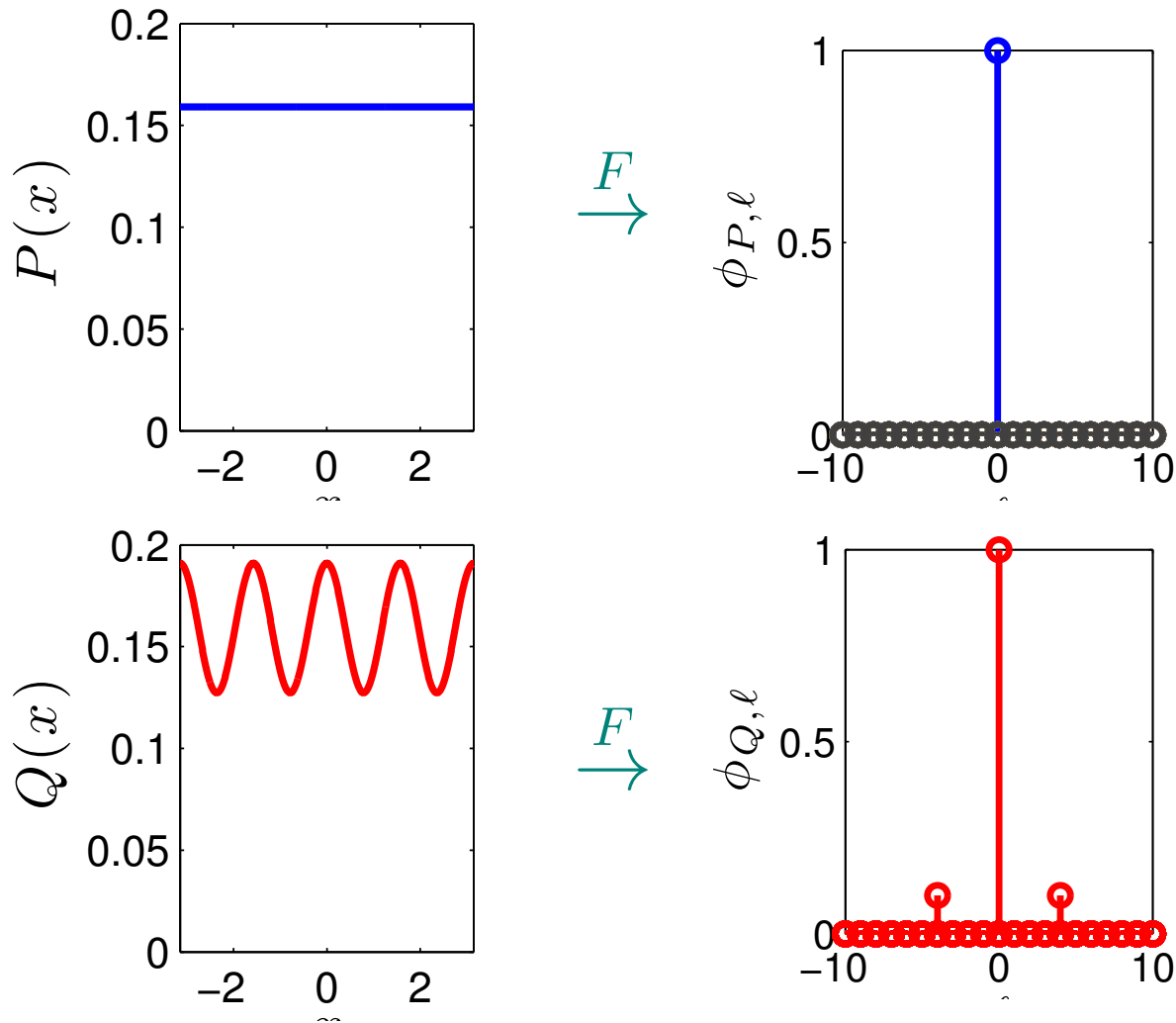
Example

- Example: **P** differs from **Q** at one frequency



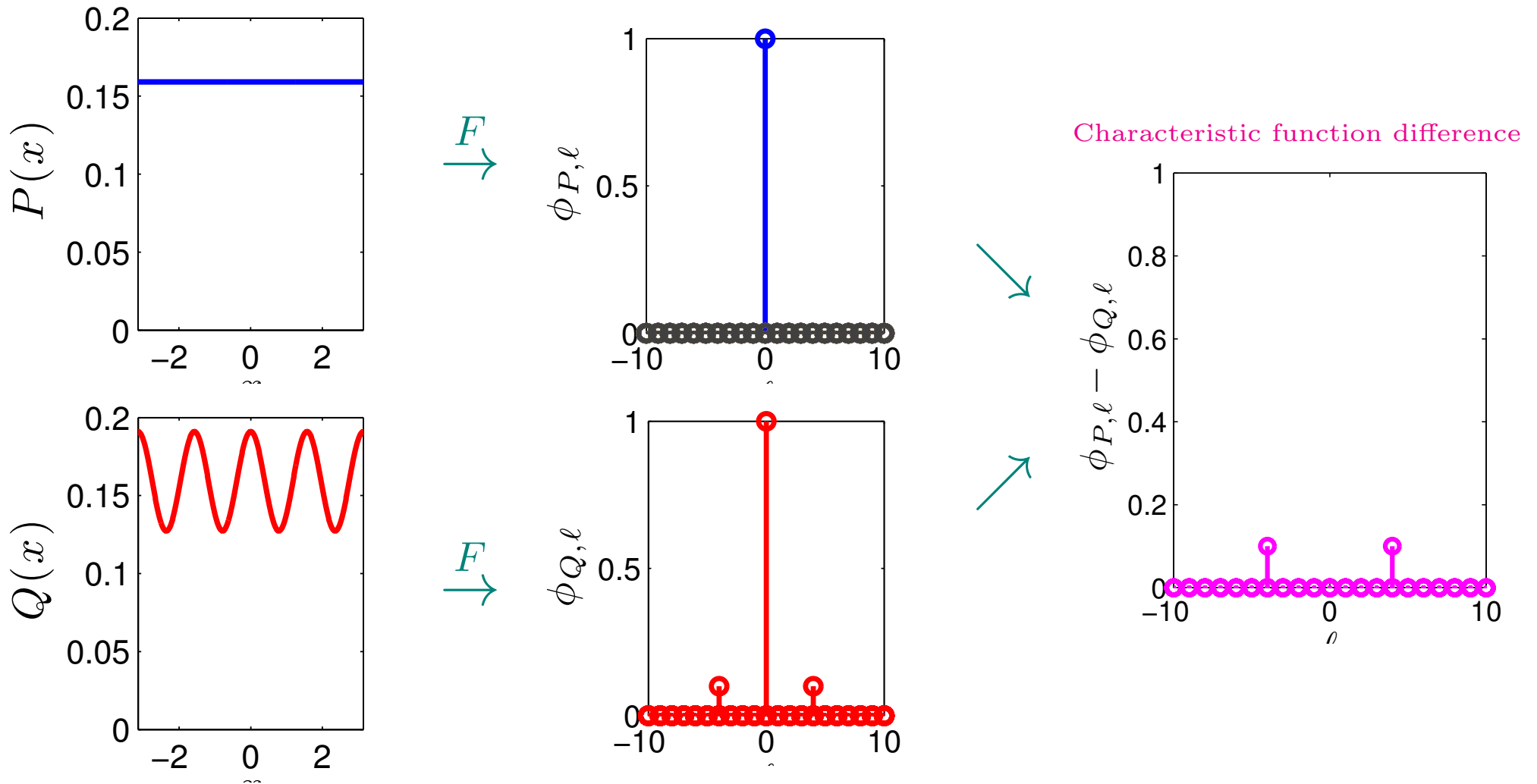
Characteristic Kernels (2)

- Example: **P** differs from **Q** at (roughly) one frequency



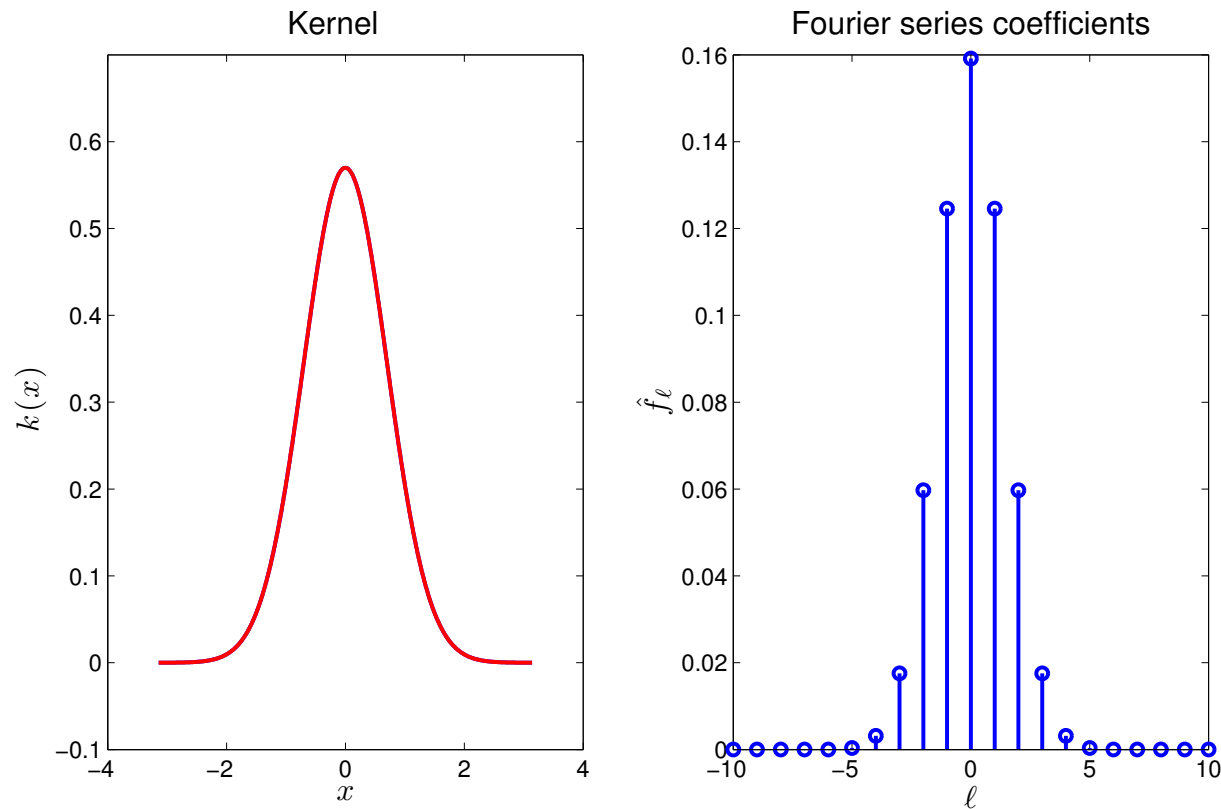
Characteristic Kernels (2)

- Example: **P** differs from **Q** at (roughly) one frequency



Example

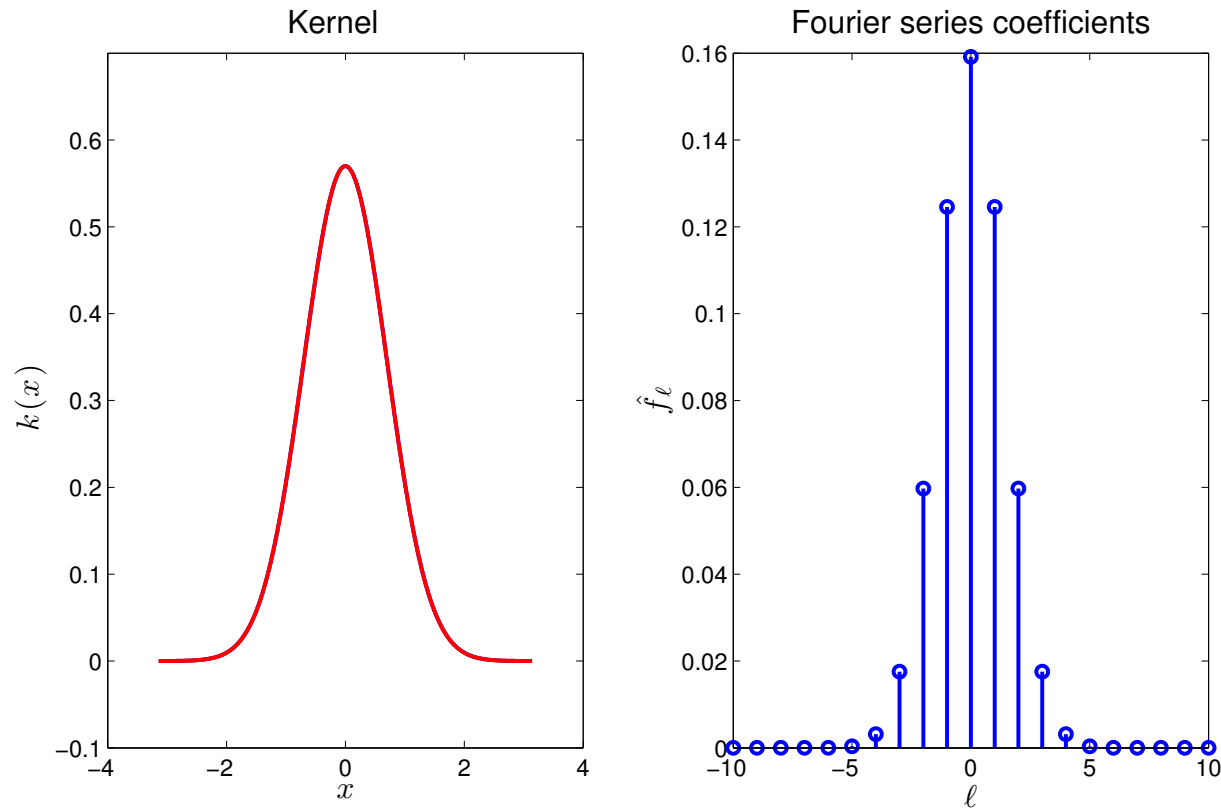
Is the **Gaussian-spectrum** kernel characteristic?



$$\text{MMD}^2(\mathbf{P}, \mathbf{Q}; F) := \sum_{l=-\infty}^{\infty} |\phi_{\mathbf{P},l} - \phi_{\mathbf{Q},l}|^2 \hat{k}_l$$

Example

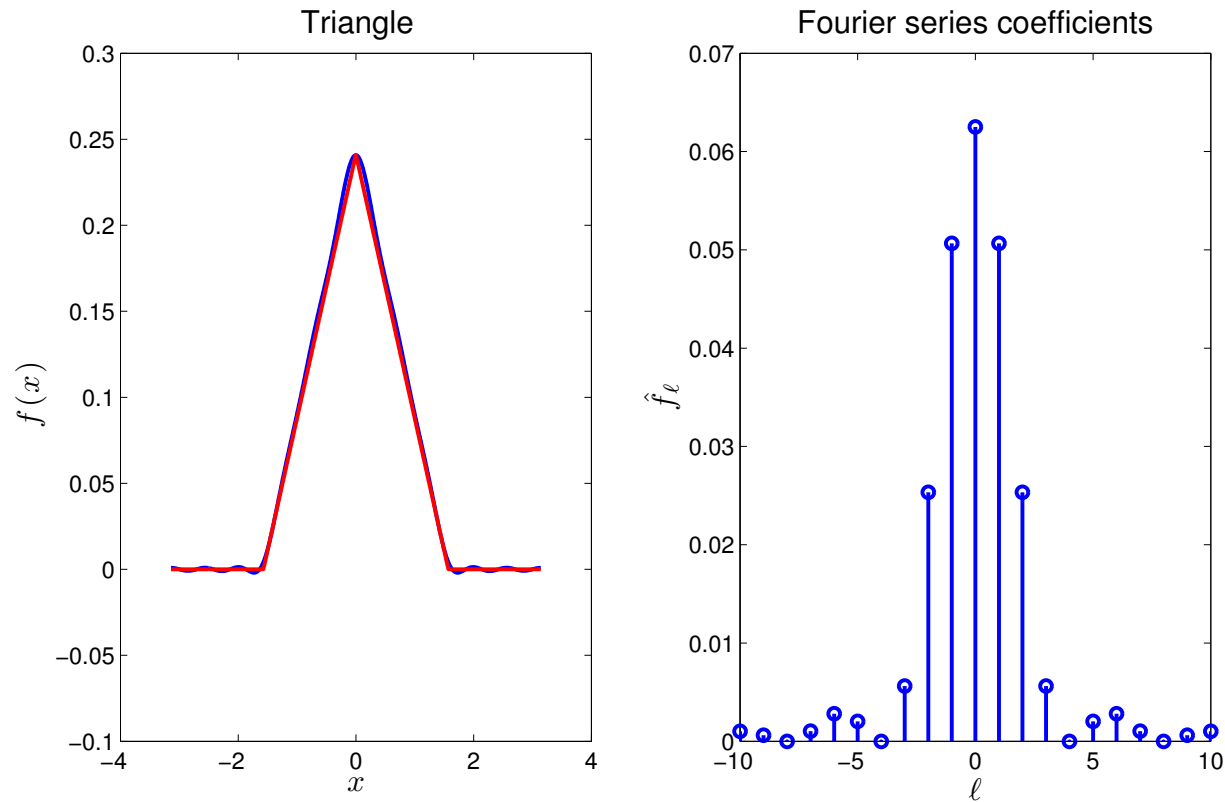
Is the Gaussian-spectrum kernel characteristic? **YES**



$$\text{MMD}^2(\mathbf{P}, \mathbf{Q}; F) := \sum_{l=-\infty}^{\infty} |\phi_{\mathbf{P},l} - \phi_{\mathbf{Q},l}|^2 \hat{k}_\ell$$

Example

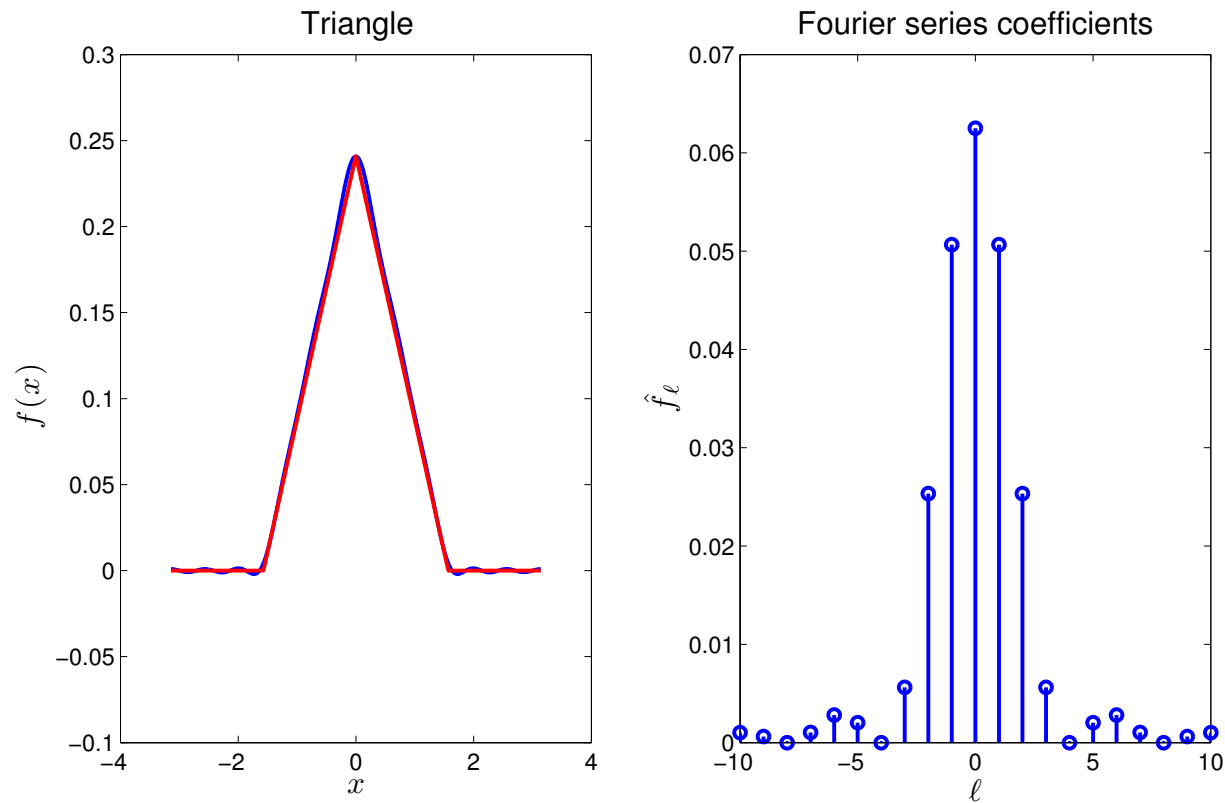
Is the **triangle** kernel characteristic?



$$\text{MMD}^2(\mathbf{P}, \mathbf{Q}; F) := \sum_{l=-\infty}^{\infty} |\phi_{\mathbf{P}, l} - \phi_{\mathbf{Q}, l}|^2 \hat{k}_\ell$$

Example

Is the **triangle** kernel characteristic? **NO**



$$\text{MMD}^2(\mathbf{P}, \mathbf{Q}; F) := \sum_{l=-\infty}^{\infty} |\phi_{\mathbf{P}, l} - \phi_{\mathbf{Q}, l}|^2 \hat{k}_l$$

Characteristic kernels (Via Fourier, on \mathbb{R}^d)

Characteristic Kernels (via Fourier)

- Can we prove **characteristic** on \mathbb{R}^d ?

Characteristic Kernels (via Fourier)

- Can we prove **characteristic on \mathbb{R}^d** ?
- **Characteristic function of \mathbf{P} via Fourier transform**

$$\phi_{\mathbf{P}}(\omega) = \int_{\mathbb{R}^d} e^{ix^\top \omega} d\mathbf{P}(x)$$

Characteristic Kernels (via Fourier)

- Can we prove **characteristic on \mathbb{R}^d** ?
- **Characteristic function of \mathbf{P} via Fourier transform**

$$\phi_{\mathbf{P}}(\omega) = \int_{\mathbb{R}^d} e^{ix^\top \omega} d\mathbf{P}(x)$$

- **Translation invariant kernels: $k(x, y) = k(x - y) = k(z)$**
- **Bochner's theorem:**

$$k(z) = \int_{\mathbb{R}^d} e^{-iz^\top \omega} d\Lambda(\omega)$$

- Λ finite non-negative Borel measure

Characteristic Kernels (via Fourier)

- Can we prove **characteristic on \mathbb{R}^d** ?
- **Characteristic function of \mathbf{P} via Fourier transform**

$$\phi_{\mathbf{P}}(\omega) = \int_{\mathbb{R}^d} e^{ix^\top \omega} d\mathbf{P}(x)$$

- **Translation invariant kernels: $k(x, y) = k(x - y) = k(z)$**
- **Bochner's theorem:**

$$k(z) = \int_{\mathbb{R}^d} e^{-iz^\top \omega} d\Lambda(\omega)$$

- Λ finite non-negative Borel measure

Characteristic Kernels (via Fourier)

- Fourier representation of MMD:

$$\text{MMD}(\mathbf{P}, \mathbf{Q}; F) := \int \int |\phi_{\mathbf{P}}(\omega) - \phi_{\mathbf{Q}}(\omega)|^2 d\Lambda(\omega)$$

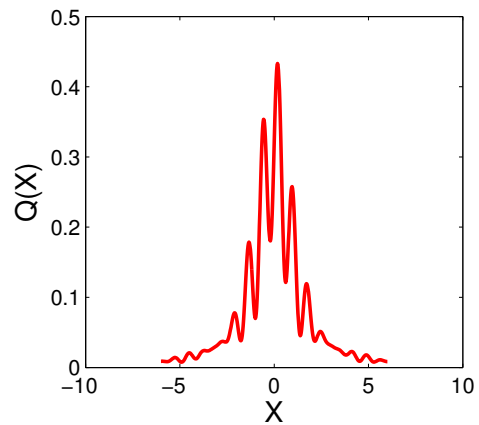
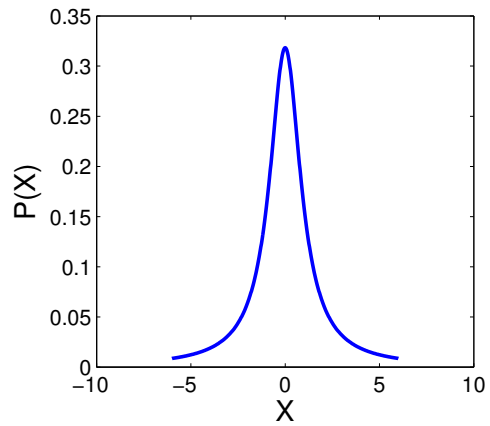
- $\phi_{\mathbf{P}}$ characteristic function of \mathbf{P}

Proof: Using Bochner's theorem (a) and Fubini's theorem (b),

$$\begin{aligned} \text{MMD}(\mathbf{P}, \mathbf{Q}) &= \int \int_{\mathbb{R}^d} k(x - y) d(\mathbf{P} - \mathbf{Q})(x) d(\mathbf{P} - \mathbf{Q})(y) \\ &\stackrel{(a)}{=} \int \int \int_{\mathbb{R}^d} e^{-i(x-y)^T \omega} d\Lambda(\omega) d(\mathbf{P} - \mathbf{Q})(x) d(\mathbf{P} - \mathbf{Q})(y) \\ &\stackrel{(b)}{=} \int \int_{\mathbb{R}^d} e^{-ix^T \omega} d(\mathbf{P} - \mathbf{Q})(x) \int_{\mathbb{R}^d} e^{iy^T \omega} d(\mathbf{P} - \mathbf{Q})(y) d\Lambda(\omega) \\ &= \int |\phi_{\mathbf{P}}(\omega) - \phi_{\mathbf{Q}}(\omega)|^2 d\Lambda(\omega) \end{aligned}$$

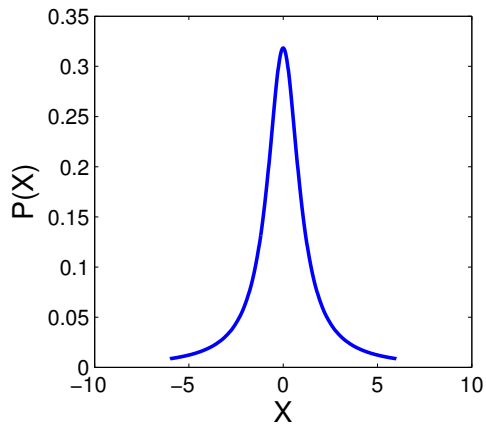
Example

- Example: **P** differs from **Q** at (roughly) one frequency

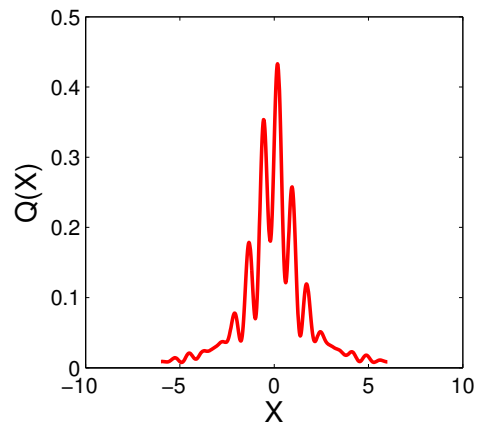
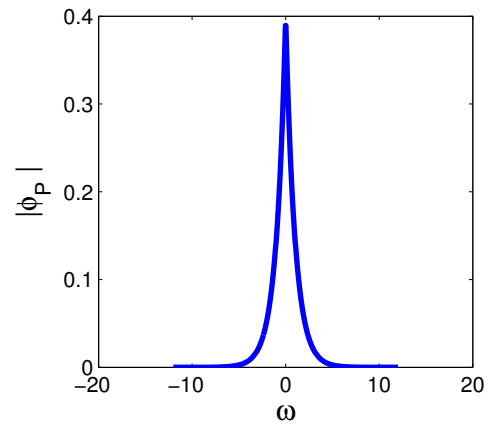


Example

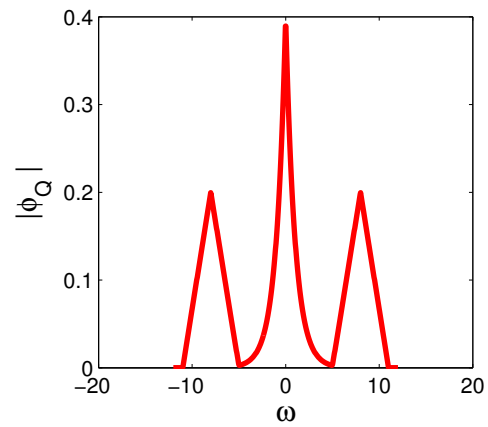
- Example: **P** differs from **Q** at (roughly) one frequency



\xrightarrow{F}

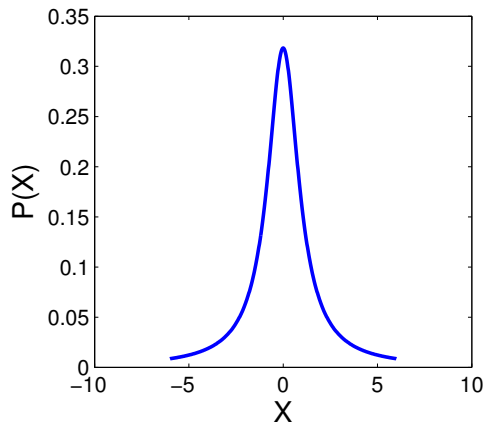


\xrightarrow{F}

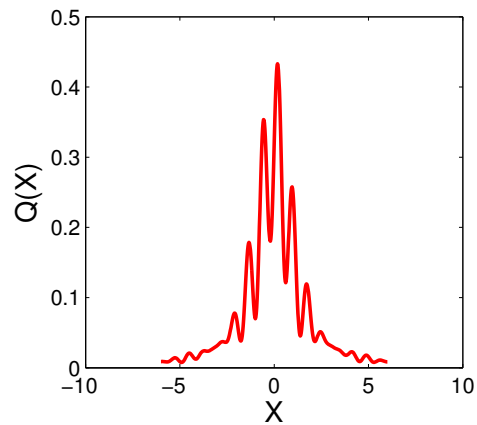
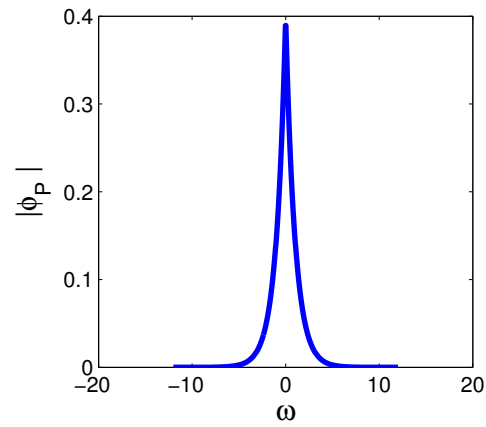


Example

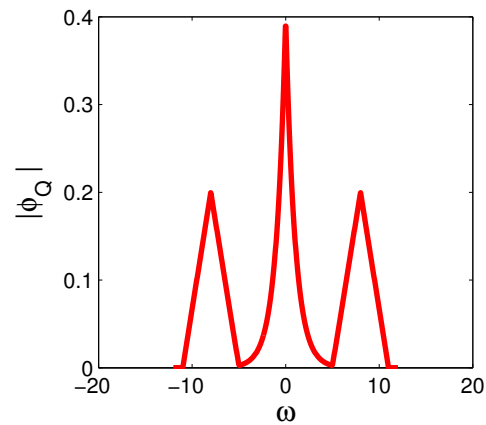
- Example: **P** differs from **Q** at (roughly) one frequency



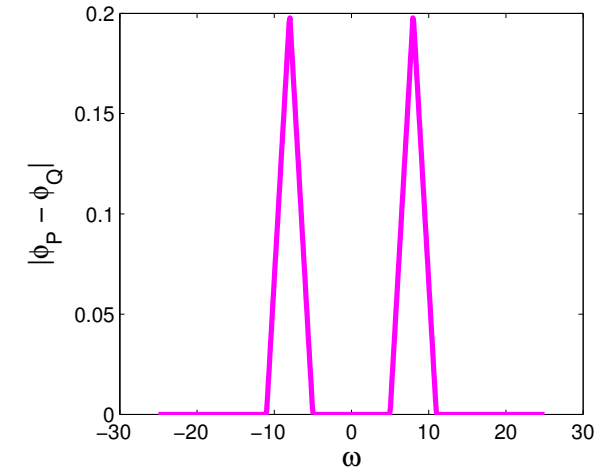
\xrightarrow{F}



\xrightarrow{F}



Characteristic function difference

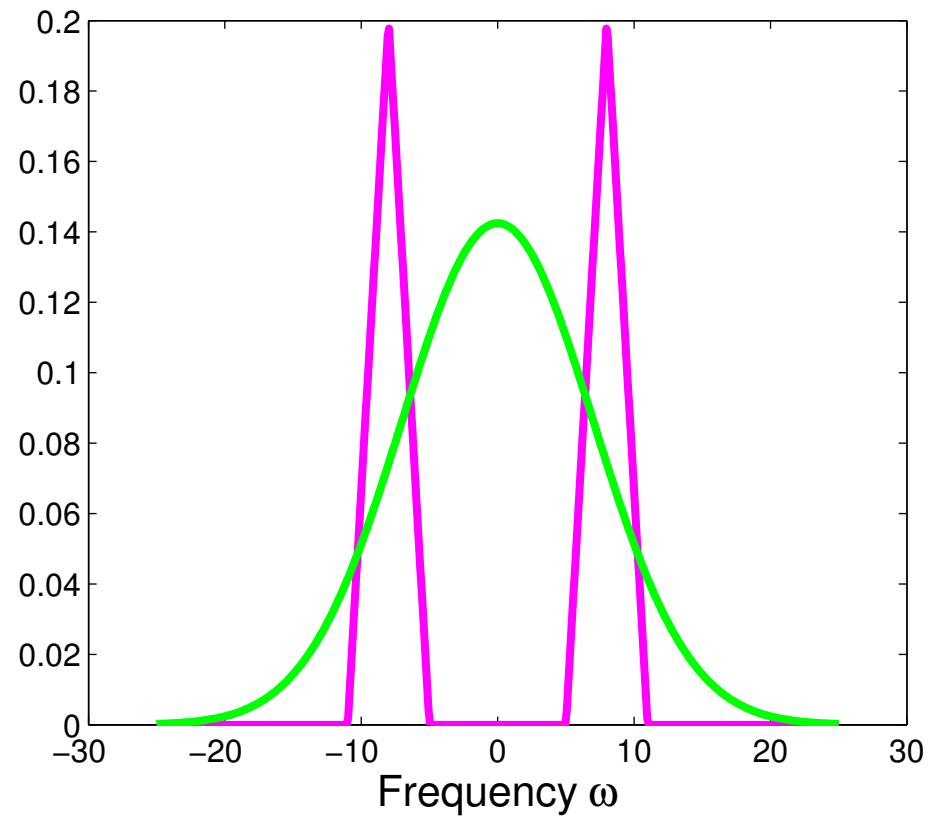


Example

- Example: **P** differs from **Q** at (roughly) one frequency

Gaussian kernel

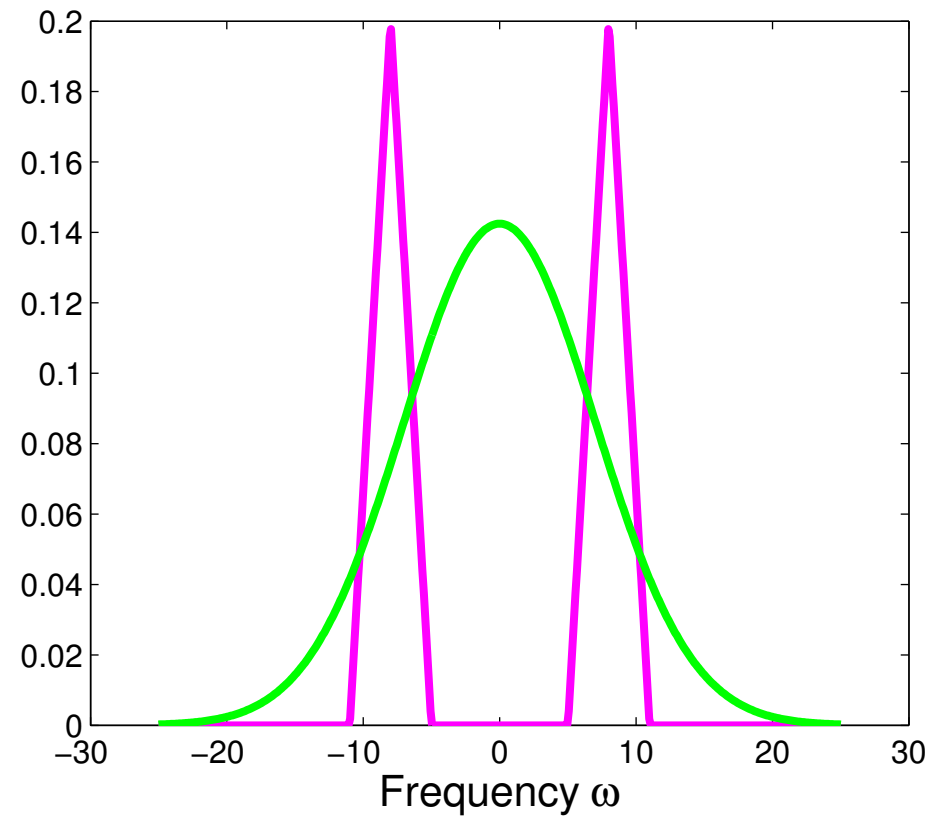
Difference $|\phi_P - \phi_Q|$



Example

- Example: **P** differs from **Q** at (roughly) one frequency

Characteristic

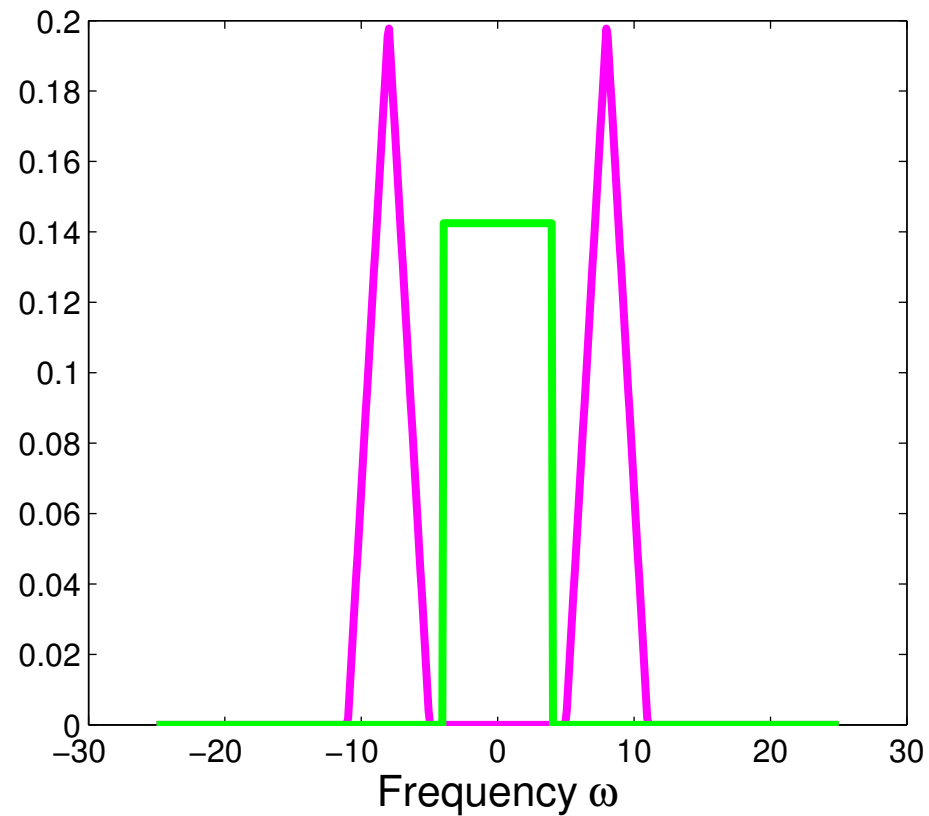


Example

- Example: **P** differs from **Q** at (roughly) one frequency

Sinc kernel

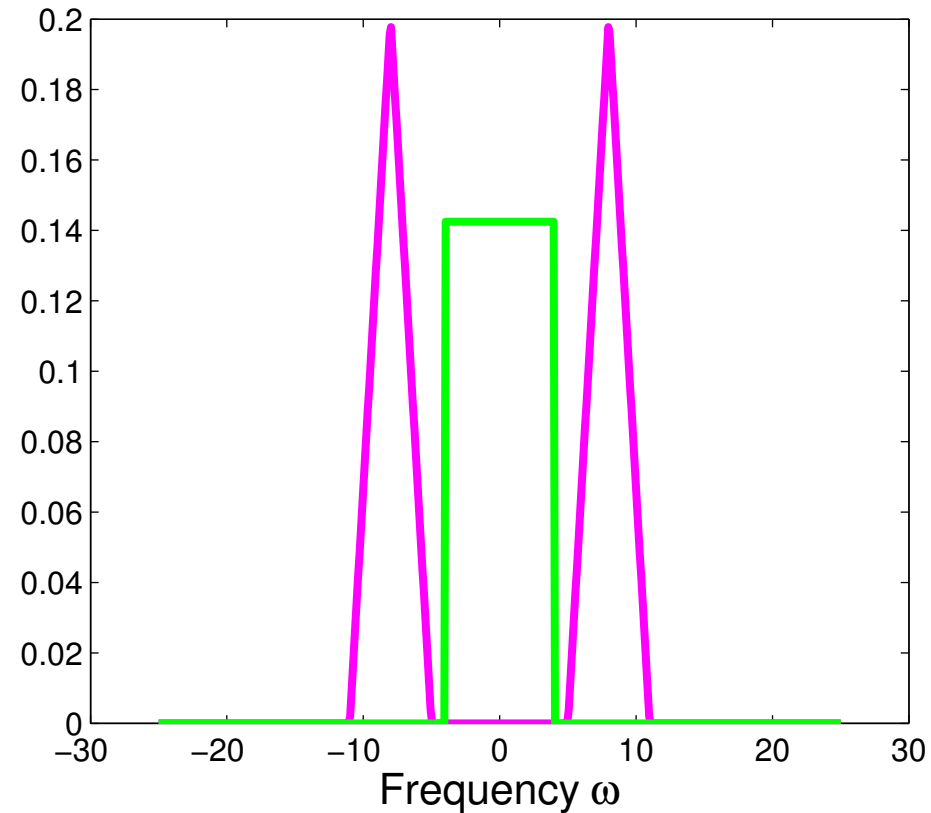
Difference $|\phi_P - \phi_Q|$



Example

- Example: **P** differs from **Q** at (roughly) one frequency

NOT characteristic

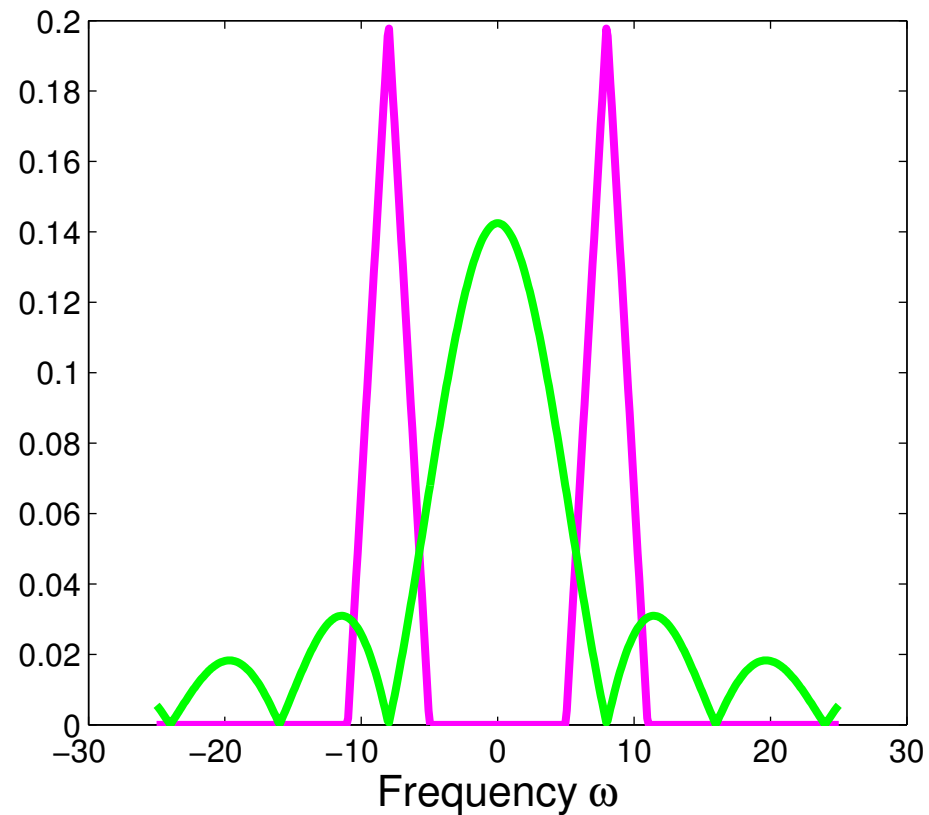


Example

- Example: **P** differs from **Q** at (roughly) one frequency

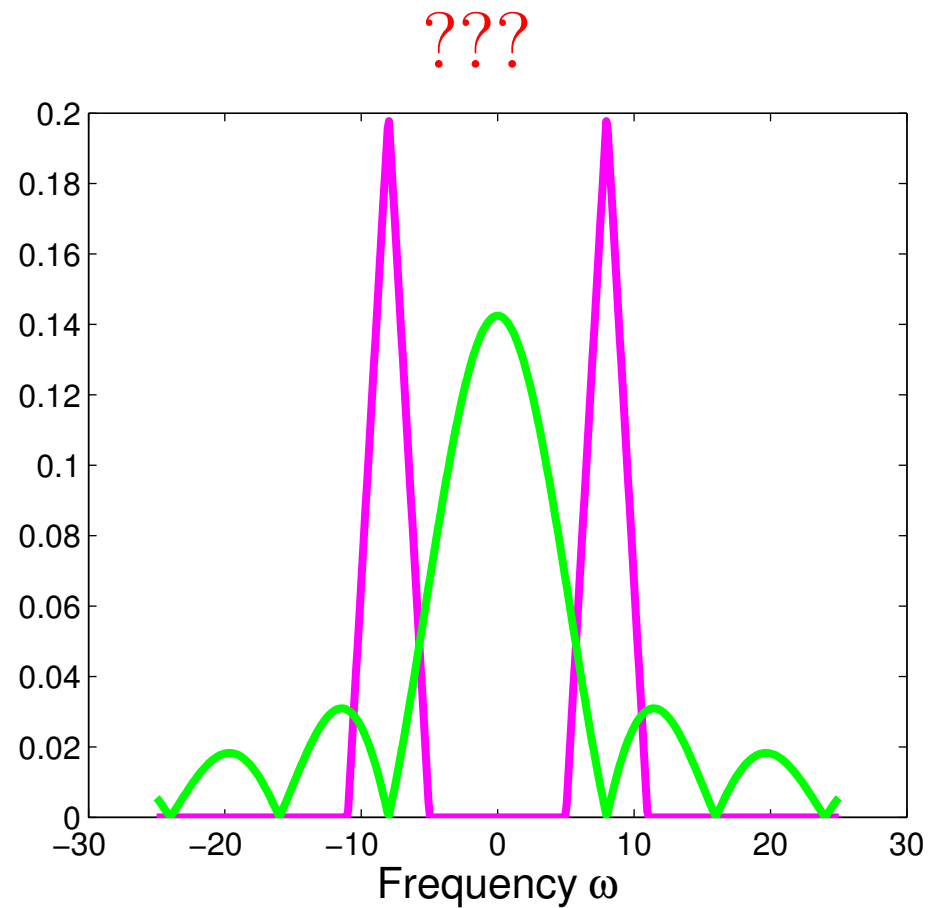
Triangle (B-spline) kernel

Difference $|\phi_P - \phi_Q|$



Example

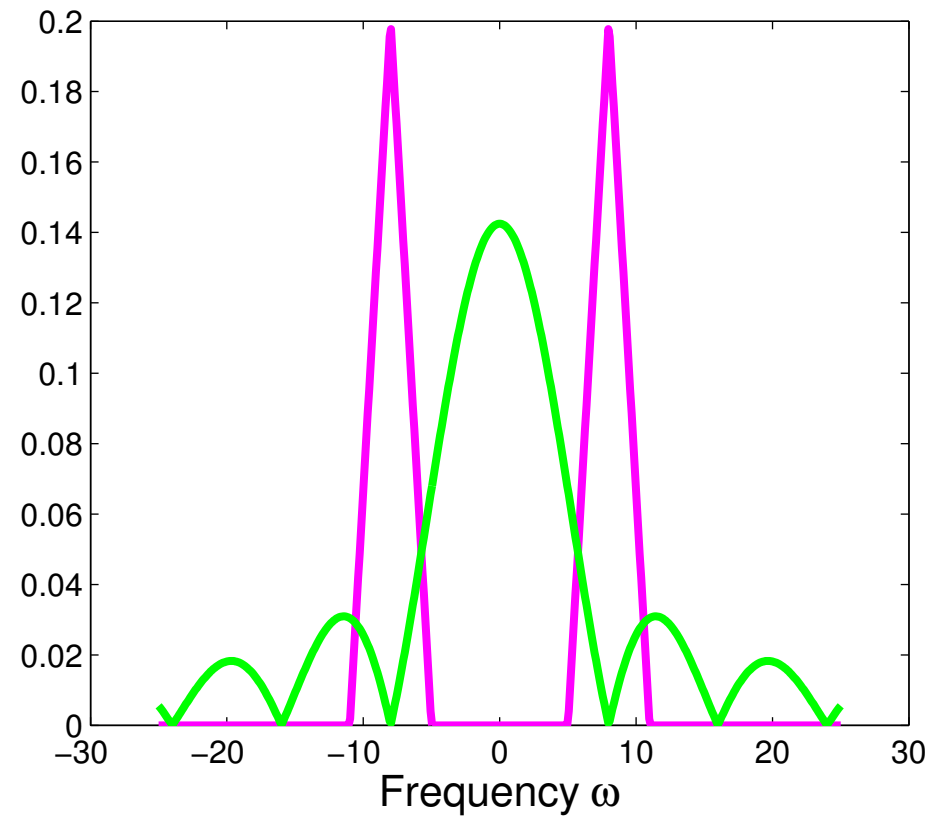
- Example: **P** differs from **Q** at (roughly) one frequency



Example

- Example: **P** differs from **Q** at (roughly) one frequency

Characteristic



Summary: Characteristic Kernels

Characteristic kernel: (MMD = 0 iff $\mathbf{P} = \mathbf{Q}$) [NIPS07b, COLT08]

Main theorem: A translation invariant k characteristic for prob. measures on \mathbb{R}^d if and only if $\text{supp}(\Lambda) = \mathbb{R}^d$ (i.e. support zero on at most a countable set)

[COLT08, JMLR10]

Corollary: continuous, compactly supported k characteristic (since Fourier spectrum $\Lambda(\omega)$ cannot be zero on an interval). 1-D proof sketch from [Mallat, 1999,

Theorem 2.6] proof on \mathbb{R}^d via distribution theory in [Sriperumbudur et al., 2010, Corollary 10 p. 1535]

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

Proof: $\text{supp} \{ \Lambda \} = \mathbb{R}^d \implies k$ characteristic:

Recall Fourier definition of MMD:

$$\text{MMD}^2(\mathbf{P}, \mathbf{Q}) = \int_{\mathbb{R}^d} |\phi_{\mathbf{P}}(\omega) - \phi_{\mathbf{Q}}(\omega)|^2 d\Lambda(\omega).$$

Characteristic functions $\phi_{\mathbf{P}}(\omega)$ and $\phi_{\mathbf{Q}}(\omega)$ **uniformly continuous**, hence their difference cannot be non-zero only on a countable set.

Map $\phi_{\mathbf{P}}$ uniformly continuous: $\forall \epsilon > 0, \exists \delta > 0$ such that $\forall (\omega_1, \omega_2) \in \Omega$ for which $d(\omega_1, \omega_2) < \delta$, we have

$d(\phi_{\mathbf{P}}(\omega_1), \phi_{\mathbf{P}}(\omega_2)) < \epsilon$. **Uniform:** δ depends only on ϵ , not on ω_1, ω_2 .

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

Proof: k characteristic $\implies \text{supp} \{ \Lambda \} = \mathbb{R}^d$:

Proof by contrapositive.

Given $\text{supp} \{ \Lambda \} \subsetneq \mathbb{R}^d$, hence \exists open interval U such that $\Lambda(\omega)$ zero on U .

Construct densities $p(x), q(x)$ such that $\phi_{\mathbf{P}}, \phi_{\mathbf{Q}}$ differ only inside U

Further extensions

- Similar reasoning wherever extensions of **Bochner's theorem** exist:

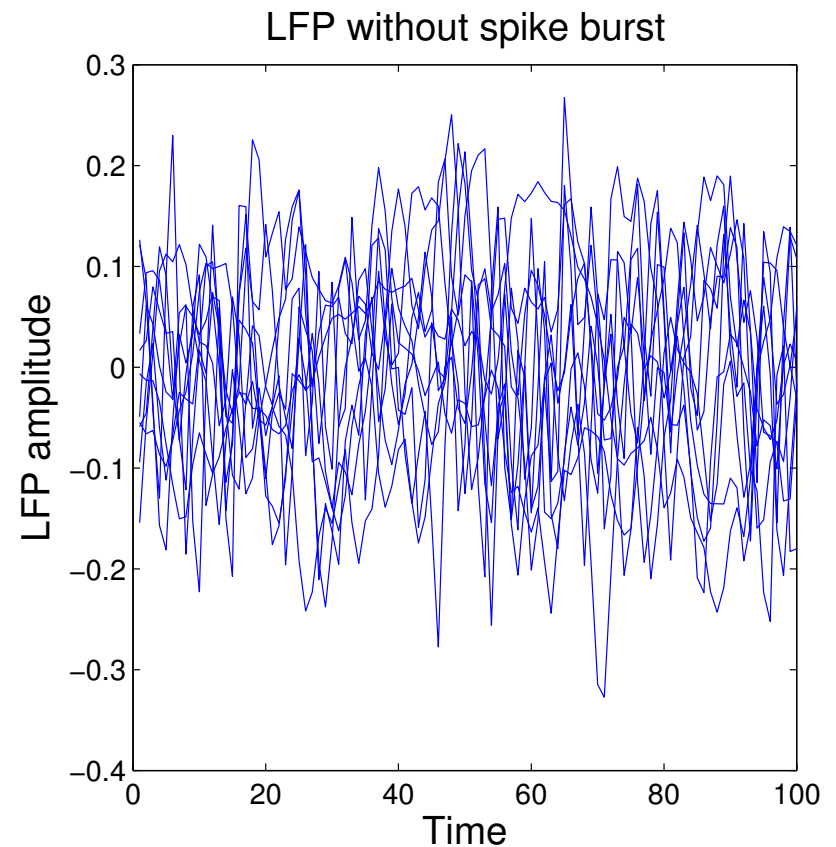
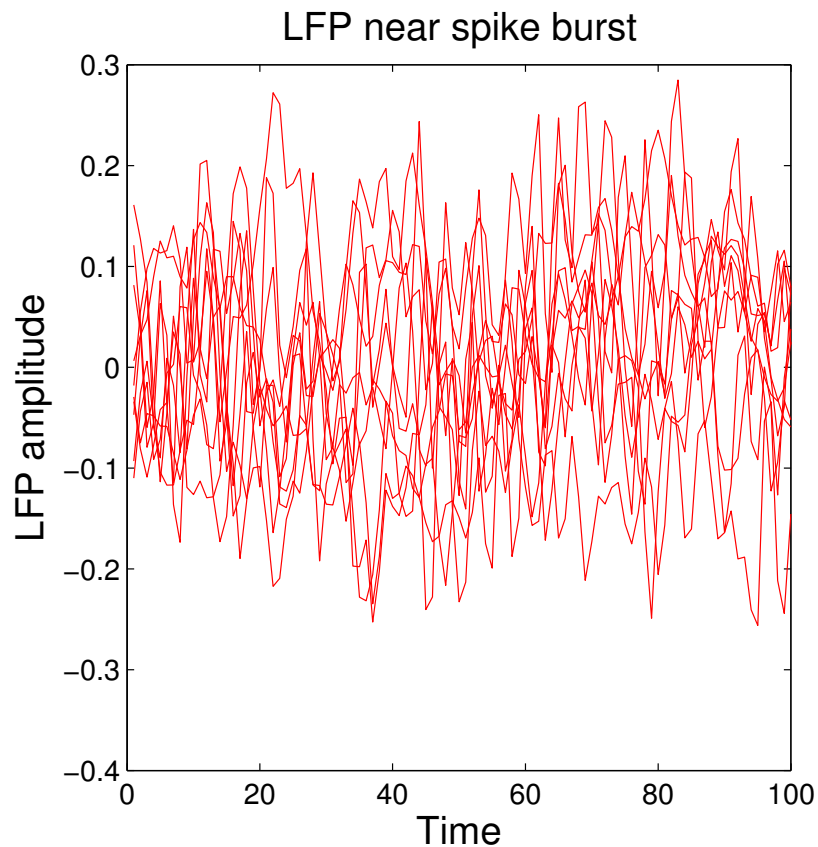
[Fukumizu et al., 2009]

- Locally compact Abelian groups (periodic domains, as we saw)
 - Compact, non-Abelian groups (orthogonal matrices)
 - The semigroup \mathbb{R}_n^+ (histograms)
- **Related kernel statistics:** Fisher statistic [Harchaoui et al., 2008] (zero iff $\mathbf{P} = \mathbf{Q}$ for characteristic kernels), other distances [Zhou and Chellappa, 2006] (not yet shown to establish whether $\mathbf{P} = \mathbf{Q}$), energy distances

Statistical hypothesis testing

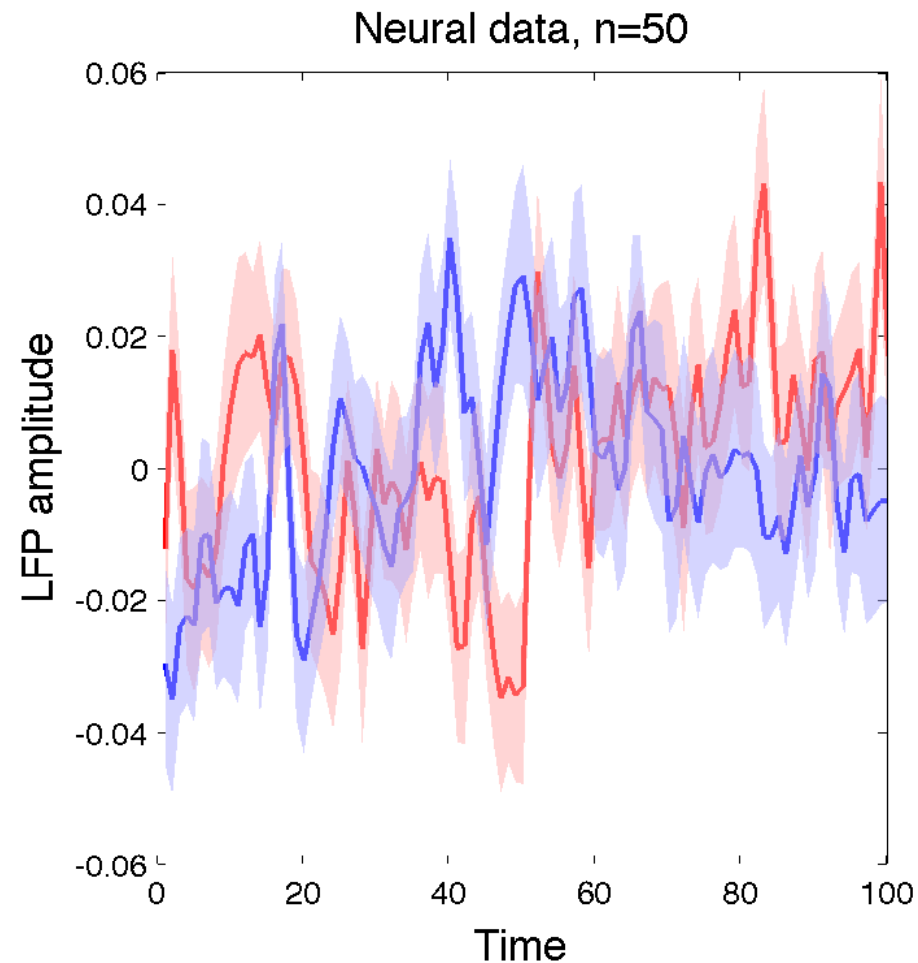
Motivating question: differences in brain signals

The problem: Do local field potential (LFP) signals change when measured near a spike burst?



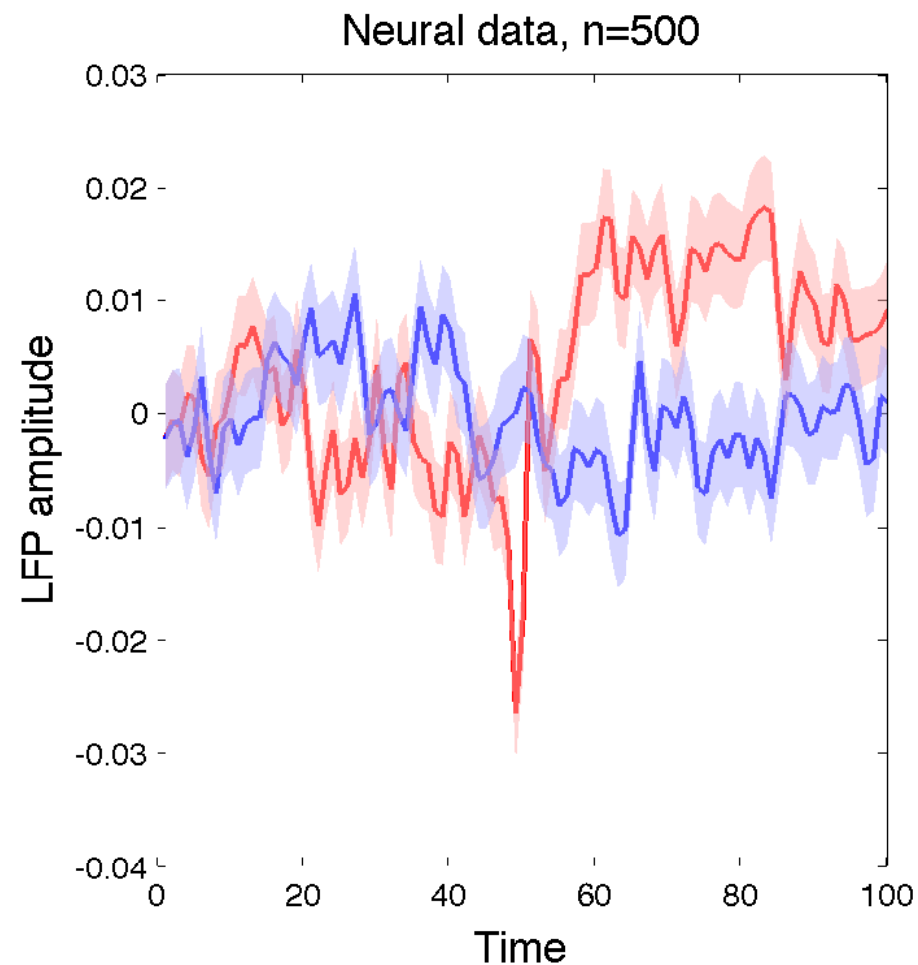
Motivating question: differences in brain signals

The problem: Do local field potential (LFP) signals change when measured near a spike burst?



Motivating question: differences in brain signals

The problem: Do local field potential (LFP) signals change when measured near a spike burst?



Statistical test using MMD (1)

- Two hypotheses:
 - H_0 : null hypothesis ($\mathbf{P} = \mathbf{Q}$)
 - H_1 : alternative hypothesis ($\mathbf{P} \neq \mathbf{Q}$)

Statistical test using MMD (1)

- Two hypotheses:
 - H_0 : null hypothesis ($\mathbf{P} = \mathbf{Q}$)
 - H_1 : alternative hypothesis ($\mathbf{P} \neq \mathbf{Q}$)
- Observe samples $\mathbf{x} := \{x_1, \dots, x_n\}$ from \mathbf{P} and \mathbf{y} from \mathbf{Q}
- If empirical MMD($\mathbf{x}, \mathbf{y}; F$) is
 - “far from zero”: reject H_0
 - “close to zero”: accept H_0

Statistical test using MMD (2)

- “far from zero” vs “close to zero” - threshold?
- **One answer:** asymptotic distribution of $\widehat{\text{MMD}}^2$

Statistical test using MMD (2)

- “far from zero” vs “close to zero” - threshold?
- **One answer:** asymptotic distribution of $\widehat{\text{MMD}}^2$
- An unbiased **empirical estimate** (quadratic cost):

$$\widehat{\text{MMD}}^2 = \frac{1}{n(n-1)} \sum_{i \neq j} \underbrace{k(x_i, x_j) - k(x_i, y_j) - k(y_i, x_j) + k(y_i, y_j)}_{h((x_i, y_i), (x_j, y_j))}$$

Statistical test using MMD (2)

- “far from zero” vs “close to zero” - threshold?
- **One answer:** asymptotic distribution of $\widehat{\text{MMD}}^2$
- An unbiased **empirical estimate** (quadratic cost):

$$\widehat{\text{MMD}}^2 = \frac{1}{n(n-1)} \sum_{i \neq j} \underbrace{k(x_i, x_j) - k(x_i, y_j) - k(y_i, x_j) + k(y_i, y_j)}_{h((x_i, y_i), (x_j, y_j))}$$

- When $\mathbf{P} \neq \mathbf{Q}$, **asymptotically normal**
 $(\sqrt{n}) \left(\widehat{\text{MMD}}^2 - \text{MMD}^2 \right) \sim \mathcal{N}(0, \sigma_u^2)$

[Hoeffding, 1948, Serfling, 1980]

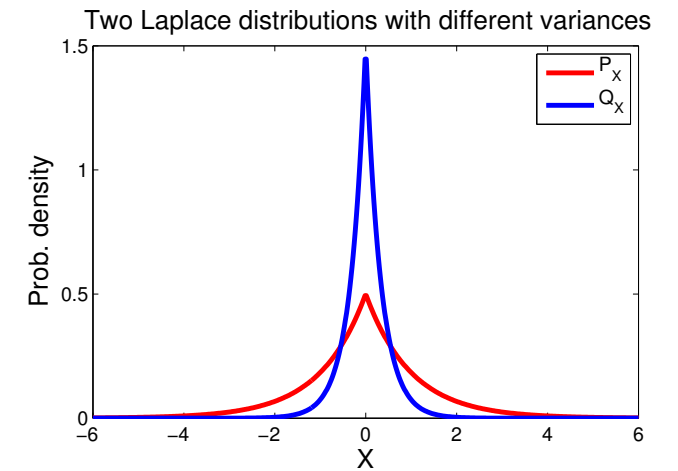
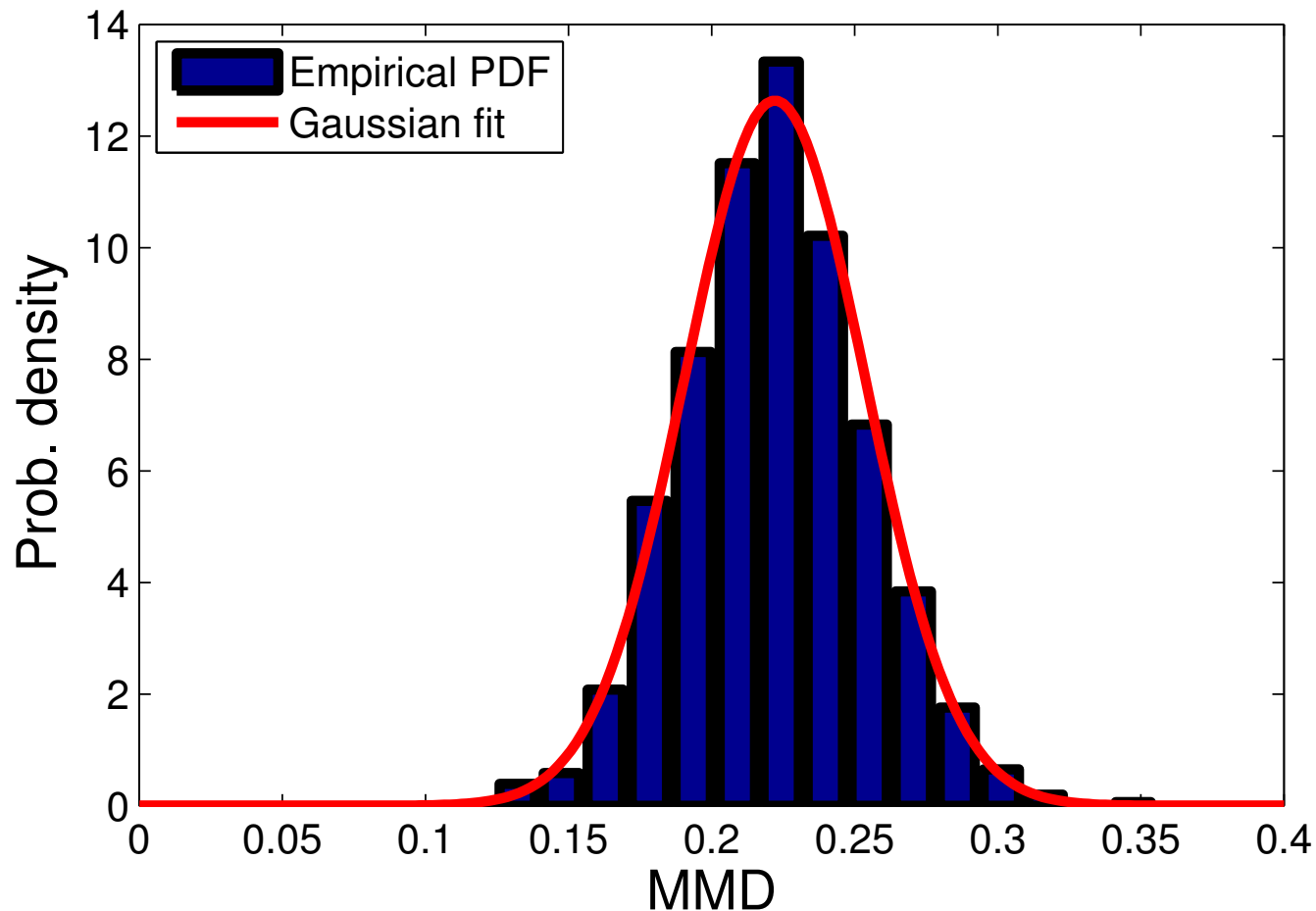
- Expression for the **variance**: $z_i := (x_i, y_i)$

$$\sigma_u^2 = 4 \left(\mathbb{E}_{\mathbf{z}} \left[\left(\mathbb{E}_{\mathbf{z}'} h(\mathbf{z}, \mathbf{z}') \right)^2 \right] - \left[\mathbb{E}_{\mathbf{z}, \mathbf{z}'} (h(\mathbf{z}, \mathbf{z}')) \right]^2 \right)$$

Statistical test using MMD (3)

- Example: laplace distributions with different variance

MMD distribution and Gaussian fit under H1



Statistical test using MMD (4)

- When $\mathbf{P} = \mathbf{Q}$, U-statistic degenerate: $\mathbb{E}_{z'}[h(z, z')] = 0$ [Anderson et al., 1994]

- Distribution is

$$n\text{MMD}(\mathbf{x}, \mathbf{y}; F) \sim \sum_{l=1}^{\infty} \lambda_l [z_l^2 - 2]$$

- where

- $z_l \sim \mathcal{N}(0, 2)$ i.i.d

- $\int_{\mathcal{X}} \underbrace{\tilde{k}(x, x')}_{\text{centred}} \psi_i(x) d\mathbf{P}(x) = \lambda_i \psi_i(x')$

Statistical test using MMD (4)

- When $\mathbf{P} = \mathbf{Q}$, U-statistic degenerate: $\mathbb{E}_{z'}[h(z, z')] = 0$ [Anderson et al., 1994]

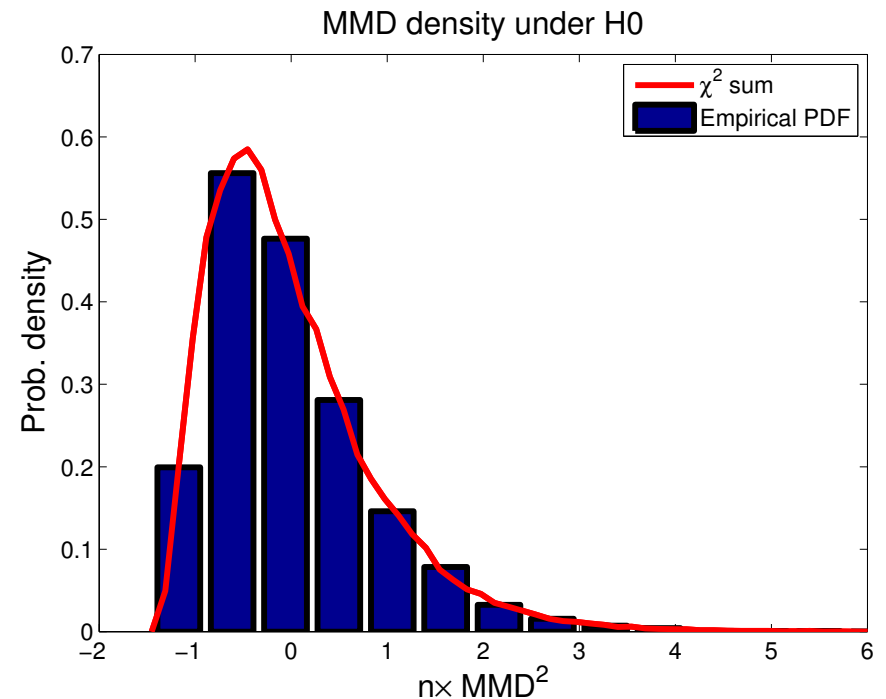
- Distribution is

$$n\text{MMD}(\mathbf{x}, \mathbf{y}; F) \sim \sum_{l=1}^{\infty} \lambda_l [z_l^2 - 2]$$

- where

- $z_l \sim \mathcal{N}(0, 2)$ i.i.d

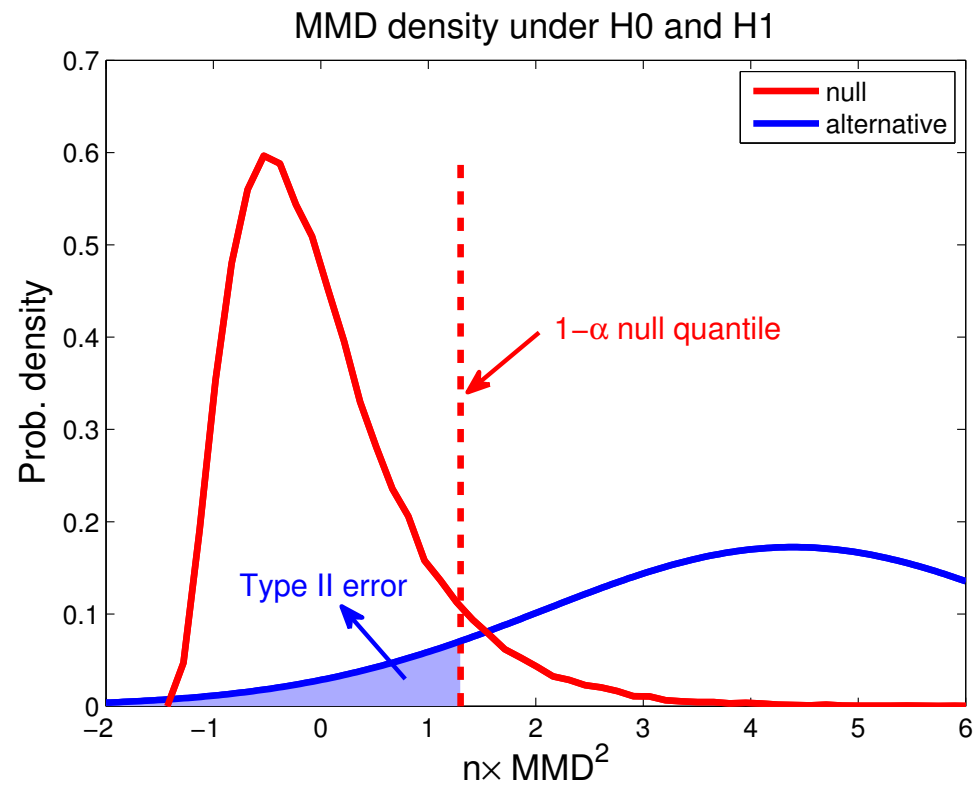
- $\int_{\mathcal{X}} \underbrace{\tilde{k}(x, x')}_{\text{centred}} \psi_i(x) d\mathbf{P}(x) = \lambda_i \psi_i(x')$



Statistical test using MMD (5)

- Given $\mathbf{P} = \mathbf{Q}$, want threshold T such that $\mathbf{P}(\text{MMD} > T) \leq 0.05$

$$\widehat{\text{MMD}}^2 = \overline{K_{P,P}} + \overline{K_{Q,Q}} - 2\overline{K_{P,Q}}$$



Statistical test using MMD (5)

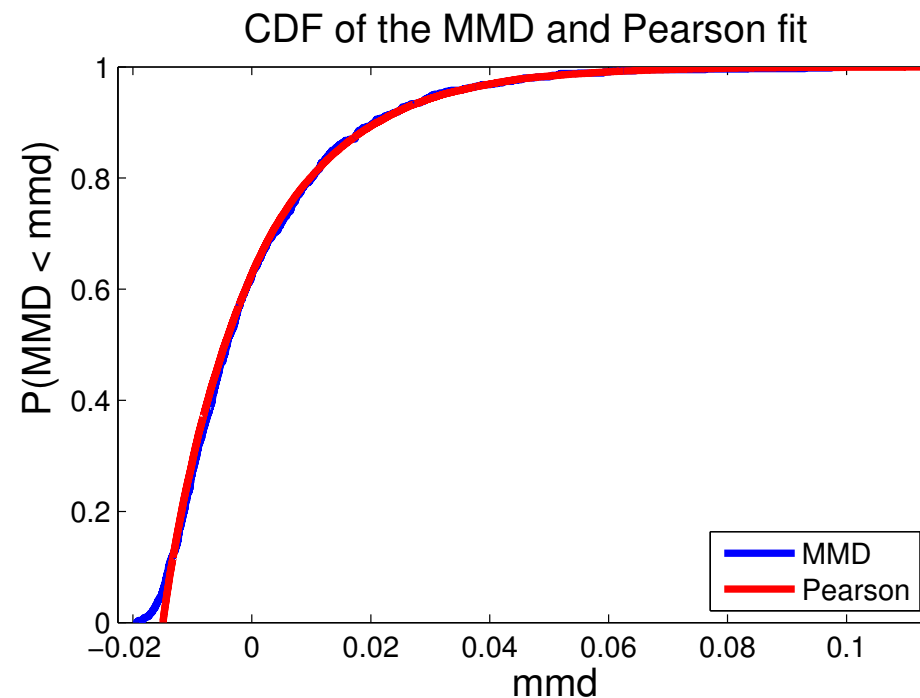
- Given $\mathbf{P} = \mathbf{Q}$, want threshold T such that $\mathbf{P}(\text{MMD} > T) \leq 0.05$

Statistical test using MMD (5)

- Given $\mathbf{P} = \mathbf{Q}$, want threshold T such that $\mathbf{P}(\text{MMD} > T) \leq 0.05$
- [Permutation](#) for empirical CDF [Arcones and Giné, 1992, Alba Fernández et al., 2008]
- [Pearson curves](#) by matching first four moments [Johnson et al., 1994]
- [Large deviation bounds](#) [Hoeffding, 1963, McDiarmid, 1989]
- [Consistent test](#) using kernel eigenspectrum [NIPS09b]

Statistical test using MMD (5)

- Given $\mathbf{P} = \mathbf{Q}$, want threshold T such that $\mathbf{P}(\text{MMD} > T) \leq 0.05$
- [Permutation](#) for empirical CDF [Arcones and Giné, 1992, Alba Fernández et al., 2008]
- [Pearson curves](#) by matching first four moments [Johnson et al., 1994]
- [Large deviation bounds](#) [Hoeffding, 1963, McDiarmid, 1989]
- [Consistent test](#) using kernel eigenspectrum [NIPS09b]



Approximate null distribution of \widehat{MMD} via permutation

Empirical MMD:

$$w = \left(\underbrace{1, 1, 1, \dots, 1}_n, \underbrace{-1, \dots, -1, -1, -1}_n \right)^\top$$

$$\frac{1}{n^2} \sum \left(\begin{bmatrix} K_{P,P} & K_{P,Q} \\ K_{Q,P} & K_{Q,Q} \end{bmatrix} \odot [ww^\top] \right) \approx \widehat{MMD}^2$$

Approximate null distribution of \widehat{MMD} via permutation

Permuted case: [Alba Fernández et al., 2008]

$$w = \left(\underbrace{1, -1, 1, \dots, 1}_n, \underbrace{-1, \dots, 1, -1, -1}_n \right)^\top$$

(equal number of +1 and -1)

$$\frac{1}{n^2} \sum \left(\begin{bmatrix} K_{P,P} & K_{P,Q} \\ K_{Q,P} & K_{Q,Q} \end{bmatrix} \odot [ww^\top] \right) = [?]$$

Approximate null distribution of \widehat{MMD} via permutation

Permuted case: [Alba Fernández et al., 2008]

$$w = \underbrace{(1, -1, 1, \dots, 1)}_n, \underbrace{(-1, \dots, 1, -1, -1)}_n)^\top$$

(equal number of +1 and -1)

$$\frac{1}{n^2} \sum \left(\begin{bmatrix} K_{P,P} & K_{P,Q} \\ K_{Q,P} & K_{Q,Q} \end{bmatrix} \odot [ww^\top] \right) = [?]$$

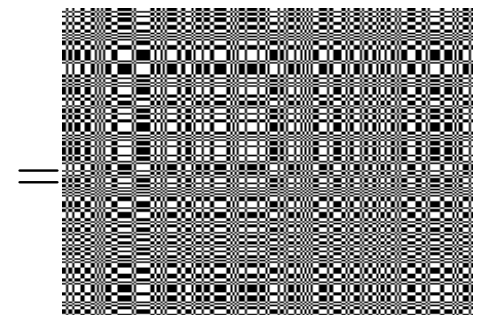
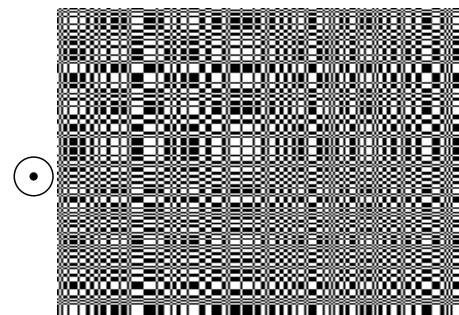
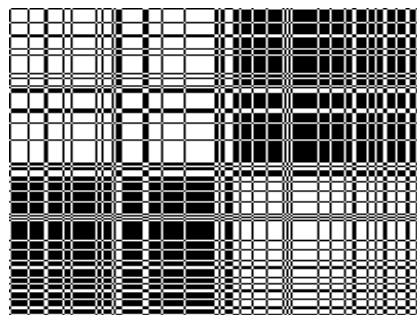
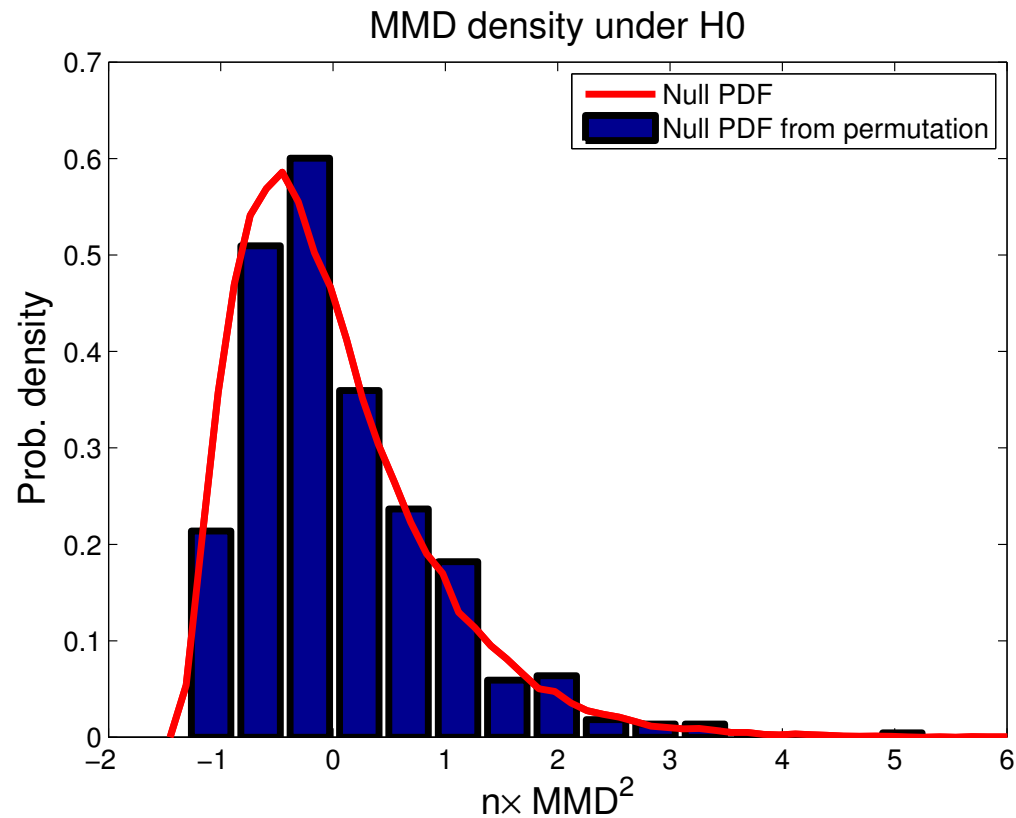


Figure thanks to Kacper Chwialkowski.

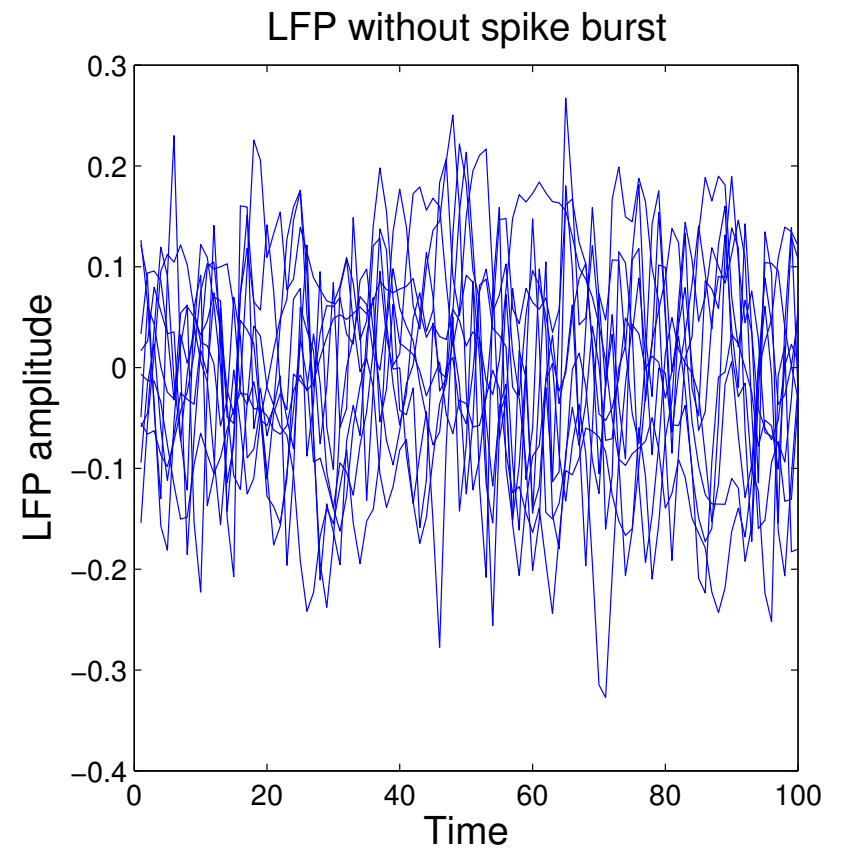
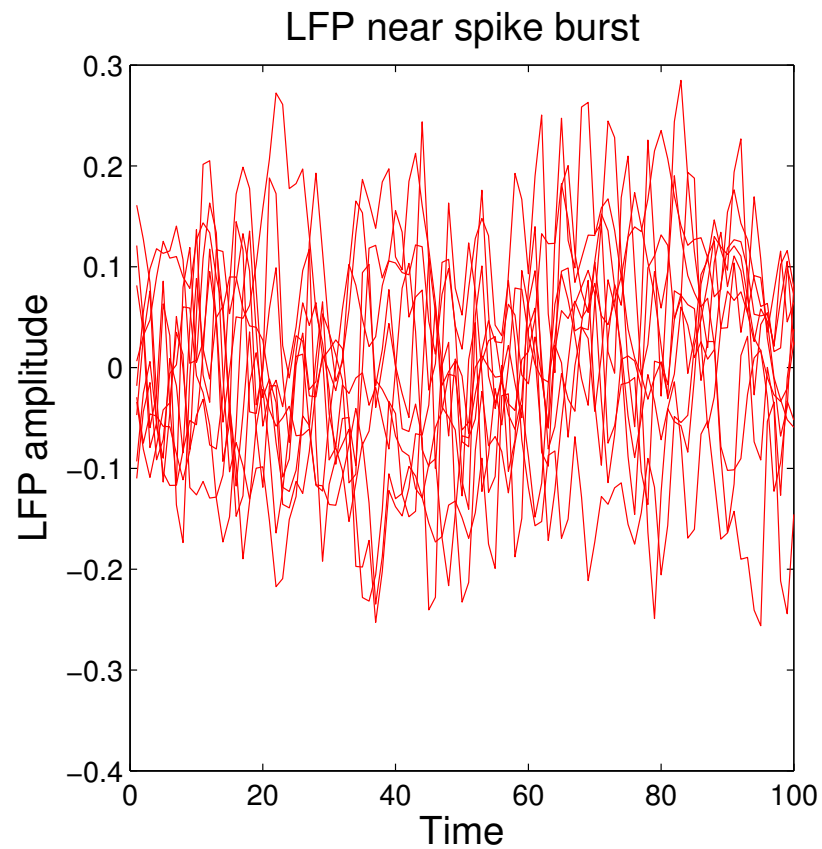
Approximate null distribution of \widehat{MMD}^2 via permutation

$$\widehat{MMD}_p^2 \approx \frac{1}{n^2} \sum \left(\begin{bmatrix} K_{P,P} & K_{P,Q} \\ K_{Q,P} & K_{Q,Q} \end{bmatrix} \odot [ww^\top] \right)$$



Detecting differences in brain signals

Do local field potential (LFP) signals change when measured near a spike burst?



Nero data: consistent test w/o permutation

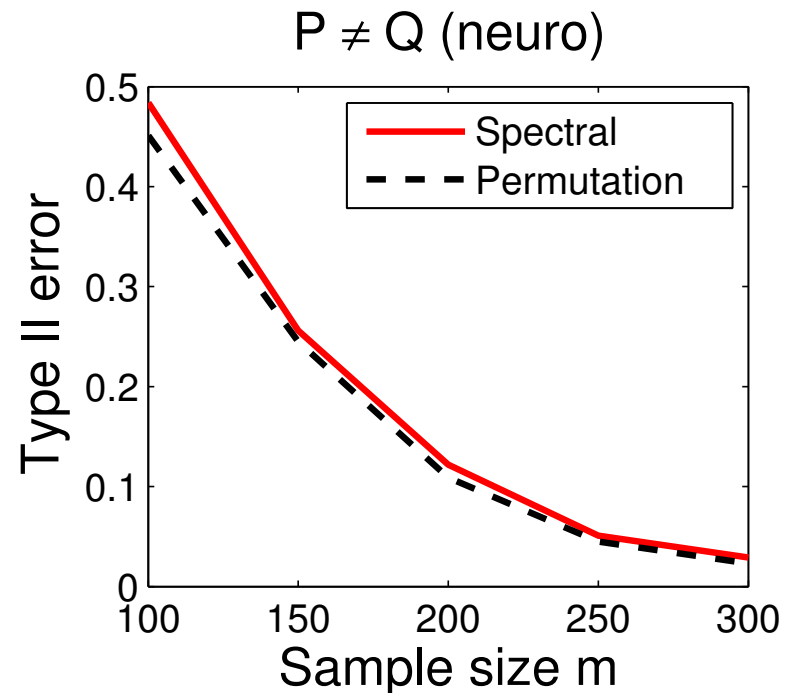
- Maximum mean discrepancy (MMD): distance between **P** and **Q**

$$\text{MMD}(\mathbf{P}, \mathbf{Q}; F) := \|\mu_{\mathbf{P}} - \mu_{\mathbf{Q}}\|_{\mathcal{F}}^2$$

- Is $\widehat{\text{MMD}}$ significantly > 0 ?
- **P** = **Q**, null distrib. of $\widehat{\text{MMD}}$:

$$n\widehat{\text{MMD}} \xrightarrow{D} \sum_{l=1}^{\infty} \lambda_l (z_l^2 - 2),$$

- λ_l is l th eigenvalue of kernel $\tilde{k}(x_i, x_j)$



Use Gram matrix spectrum for $\hat{\lambda}_l$: consistent test without permutation

Hypothesis testing with HSIC

Distribution of HSIC at independence

- (Biased) empirical HSIC a v-statistic

$$HSIC_b = \frac{1}{n^2} \text{trace}(KHLH)$$

- Statistical testing: How do we find when this is larger enough that the null hypothesis $\mathbf{P} = \mathbf{P}_x \mathbf{P}_y$ is unlikely?
- Formally: given $\mathbf{P} = \mathbf{P}_x \mathbf{P}_y$, what is the threshold T such that $\mathbf{P}(\text{HSIC} > T) < \alpha$ for small α ?

Distribution of HSIC at independence

- (Biased) empirical HSIC a v-statistic

$$HSIC_b = \frac{1}{n^2} \text{trace}(KHLH)$$

- Associated U-statistic degenerate when $\mathbf{P} = \mathbf{P}_x \mathbf{P}_y$ [Serfling, 1980]:

$$nHSIC_b \xrightarrow{D} \sum_{l=1}^{\infty} \lambda_l z_l^2, \quad z_l \sim \mathcal{N}(0, 1) \text{i.i.d.}$$

$$\lambda_l \psi_l(z_j) = \int h_{ijqr} \psi_l(z_i) dF_{i,q,r}, \quad h_{ijqr} = \frac{1}{4!} \sum_{(t,u,v,w)}^{(i,j,q,r)} k_{tu} l_{tu} + k_{tu} l_{vw} - 2k_{tu} l_{tv}$$

Distribution of HSIC at independence

- (Biased) empirical HSIC a v-statistic

$$HSIC_b = \frac{1}{n^2} \text{trace}(KHLH)$$

- Associated U-statistic degenerate when $\mathbf{P} = \mathbf{P}_x \mathbf{P}_y$ [Serfling, 1980]:

$$nHSIC_b \xrightarrow{D} \sum_{l=1}^{\infty} \lambda_l z_l^2, \quad z_l \sim \mathcal{N}(0, 1) \text{i.i.d.}$$

$$\lambda_l \psi_l(z_j) = \int h_{ijqr} \psi_l(z_i) dF_{i,q,r}, \quad h_{ijqr} = \frac{1}{4!} \sum_{(t,u,v,w)}^{(i,j,q,r)} k_{tu} l_{tu} + k_{tu} l_{vw} - 2k_{tu} l_{tv}$$

- First two moments [NIPS07b]

$$\mathbf{E}(HSIC_b) = \frac{1}{n} \text{Tr} C_{xx} \text{Tr} C_{yy}$$
$$\text{var}(HSIC_b) = \frac{2(n-4)(n-5)}{(n)_4} \|C_{xx}\|_{\text{HS}}^2 \|C_{yy}\|_{\text{HS}}^2 + O(n^{-3}).$$

Statistical testing with HSIC

- Given $\mathbf{P} = \mathbf{P}_x \mathbf{P}_y$, what is the threshold T such that $\mathbf{P}(\text{HSIC} > T) < \alpha$ for small α ?
- Null distribution via **permutation** [Feuerverger, 1993]
 - Compute HSIC for $\{x_i, y_{\pi(i)}\}_{i=1}^n$ for random permutation π of indices $\{1, \dots, n\}$. This gives HSIC for independent variables.
 - Repeat for many different permutations, get empirical CDF
 - Threshold T is $1 - \alpha$ quantile of empirical CDF

Statistical testing with HSIC

- Given $\mathbf{P} = \mathbf{P}_x \mathbf{P}_y$, what is the threshold T such that $\mathbf{P}(\text{HSIC} > T) < \alpha$ for small α ?
- Null distribution via **permutation** [Feuerverger, 1993]
 - Compute HSIC for $\{x_i, y_{\pi(i)}\}_{i=1}^n$ for random permutation π of indices $\{1, \dots, n\}$. This gives HSIC for independent variables.
 - Repeat for many different permutations, get empirical CDF
 - Threshold T is $1 - \alpha$ quantile of empirical CDF
- Approximate null distribution via **moment matching** [Kankainen, 1995]:

$$n\text{HSIC}_b(Z) \sim \frac{x^{\alpha-1} e^{-x/\beta}}{\beta^\alpha \Gamma(\alpha)}$$

where

$$\alpha = \frac{(\mathbf{E}(\text{HSIC}_b))^2}{\text{var}(\text{HSIC}_b)}, \quad \beta = \frac{\text{var}(\text{HSIC}_b)}{n\mathbf{E}(\text{HSIC}_b)}.$$

Experiment: dependence testing for translation

Are the French text extracts translations of English?

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.

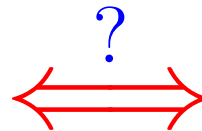
X_2 : No doubt there is great pressure on provincial and municipal governments in relation to the issue of child care, but the reality is that there have been no cuts to child care funding from the federal government to the provinces. In fact, we have increased federal investments for early childhood 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.

...



Experiment: dependence testing for translation

- (Biased) empirical HSIC:

$$HSIC_b = \frac{1}{n^2} \text{trace}(KHLH)$$

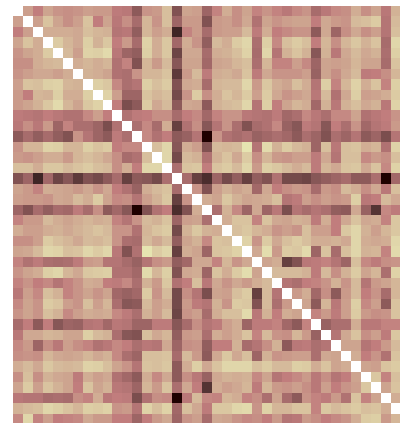
- Translation example: [NIPS07b]

Canadian Hansard
(agriculture)

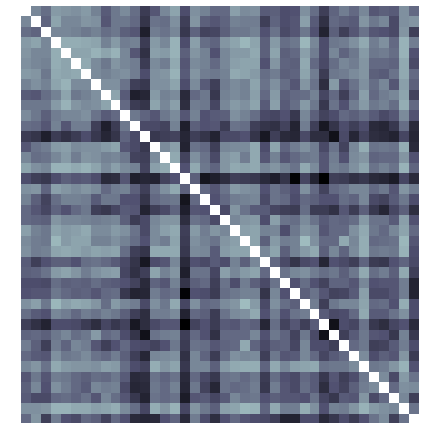
- 5-line extracts,
k-spectrum kernel, $k = 10$,
repetitions=300,
sample size 10

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

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



K



L

⇒ HSIC ⇐

- *k*-spectrum kernel: average Type II error 0 ($\alpha = 0.05$)

Experiment: dependence testing for translation

- (Biased) empirical HSIC:

$$HSIC_b = \frac{1}{n^2} \text{trace}(KHLH)$$

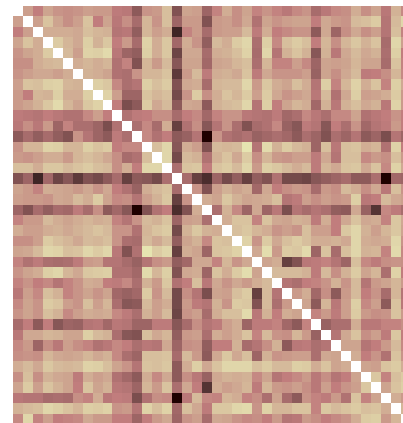
- Translation example: [NIPS07b]

Canadian Hansard
(agriculture)

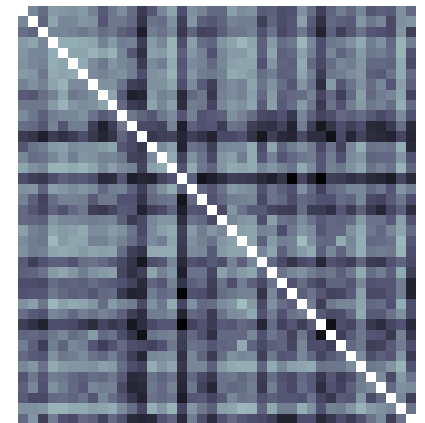
- 5-line extracts,
 k -spectrum kernel, $k = 10$,
repetitions=300,
sample size 10

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

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



K



L

\Rightarrow HSIC \Leftarrow

- k -spectrum kernel: average Type II error 0 ($\alpha = 0.05$)
- Bag of words kernel: average Type II error 0.18

Kernel two-sample tests for big data, optimal kernel choice

Quadratic time estimate of MMD

$$\text{MMD}^2 = \|\mu_{\mathbf{P}} - \mu_{\mathbf{Q}}\|_{\mathcal{F}}^2 = \mathbf{E}_{\mathbf{P}}k(x, x') + \mathbf{E}_{\mathbf{Q}}k(y, y') - 2\mathbf{E}_{\mathbf{P}, \mathbf{Q}}k(x, y)$$

Quadratic time estimate of MMD

$$\text{MMD}^2 = \|\mu_{\mathbf{P}} - \mu_{\mathbf{Q}}\|_{\mathcal{F}}^2 = \mathbf{E}_{\mathbf{P}}k(x, x') + \mathbf{E}_{\mathbf{Q}}k(y, y') - 2\mathbf{E}_{\mathbf{P}, \mathbf{Q}}k(x, y)$$

Given i.i.d. $X := \{x_1, \dots, x_m\}$ and $Y := \{y_1, \dots, y_m\}$ from \mathbf{P}, \mathbf{Q} , respectively:

The earlier estimate: (quadratic time)

$$\hat{\mathbf{E}}_{\mathbf{P}}k(x, x') = \frac{1}{m(m-1)} \sum_{i=1}^m \sum_{j \neq i}^m k(x_i, x_j)$$

Quadratic time estimate of MMD

$$\text{MMD}^2 = \|\mu_{\mathbf{P}} - \mu_{\mathbf{Q}}\|_{\mathcal{F}}^2 = \mathbf{E}_{\mathbf{P}}k(x, x') + \mathbf{E}_{\mathbf{Q}}k(y, y') - 2\mathbf{E}_{\mathbf{P}, \mathbf{Q}}k(x, y)$$

Given i.i.d. $X := \{x_1, \dots, x_m\}$ and $Y := \{y_1, \dots, y_m\}$ from \mathbf{P}, \mathbf{Q} , respectively:

The earlier estimate: (quadratic time)

$$\widehat{\mathbf{E}}_{\mathbf{P}}k(x, x') = \frac{1}{m(m-1)} \sum_{i=1}^m \sum_{j \neq i}^m k(x_i, x_j)$$

New, linear time estimate:

$$\begin{aligned} \widehat{\mathbf{E}}_{\mathbf{P}}k(x, x') &= \frac{2}{m} [k(x_1, x_2) + k(x_3, x_4) + \dots] \\ &= \frac{2}{m} \sum_{i=1}^{m/2} k(x_{2i-1}, x_{2i}) \end{aligned}$$

Linear time MMD

Shorter expression with explicit k dependence:

$$\text{MMD}^2 =: \eta_k(p, q) = \mathbb{E}_{xx'yy'} h_k(x, x', y, y') =: \mathbb{E}_v h_k(v),$$

where

$$h_k(x, x', y, y') = k(x, x') + k(y, y') - k(x, y') - k(x', y),$$

and $v := [x, x', y, y']$.

Linear time MMD

Shorter expression with explicit k dependence:

$$\text{MMD}^2 =: \eta_k(p, q) = \mathbb{E}_{xx'yy'} h_k(x, x', y, y') =: \mathbb{E}_v h_k(v),$$

where

$$h_k(x, x', y, y') = k(x, x') + k(y, y') - k(x, y') - k(x', y),$$

and $v := [x, x', y, y']$.

The linear time estimate again:

$$\check{\eta}_k = \frac{2}{m} \sum_{i=1}^{m/2} h_k(v_i),$$

where $v_i := [x_{2i-1}, x_{2i}, y_{2i-1}, y_{2i}]$ and

$$h_k(v_i) := k(x_{2i-1}, x_{2i}) + k(y_{2i-1}, y_{2i}) - k(x_{2i-1}, y_{2i}) - k(x_{2i}, y_{2i-1})$$

Linear time vs quadratic time MMD

Disadvantages of linear time MMD vs quadratic time MMD

- Much higher variance for a given m , hence...
- ...a much less powerful test for a given m

Linear time vs quadratic time MMD

Disadvantages of linear time MMD vs quadratic time MMD

- Much higher variance for a given m , hence...
- ...a much less powerful test for a given m

Advantages of the linear time MMD vs quadratic time MMD

- Very simple asymptotic null distribution (a Gaussian, vs an infinite weighted sum of χ^2)
- Both test statistic and threshold computable in $O(m)$, with storage $O(1)$.
- Given unlimited data, a given Type II error can be attained with less computation

Asymptotics of linear time MMD

By central limit theorem,

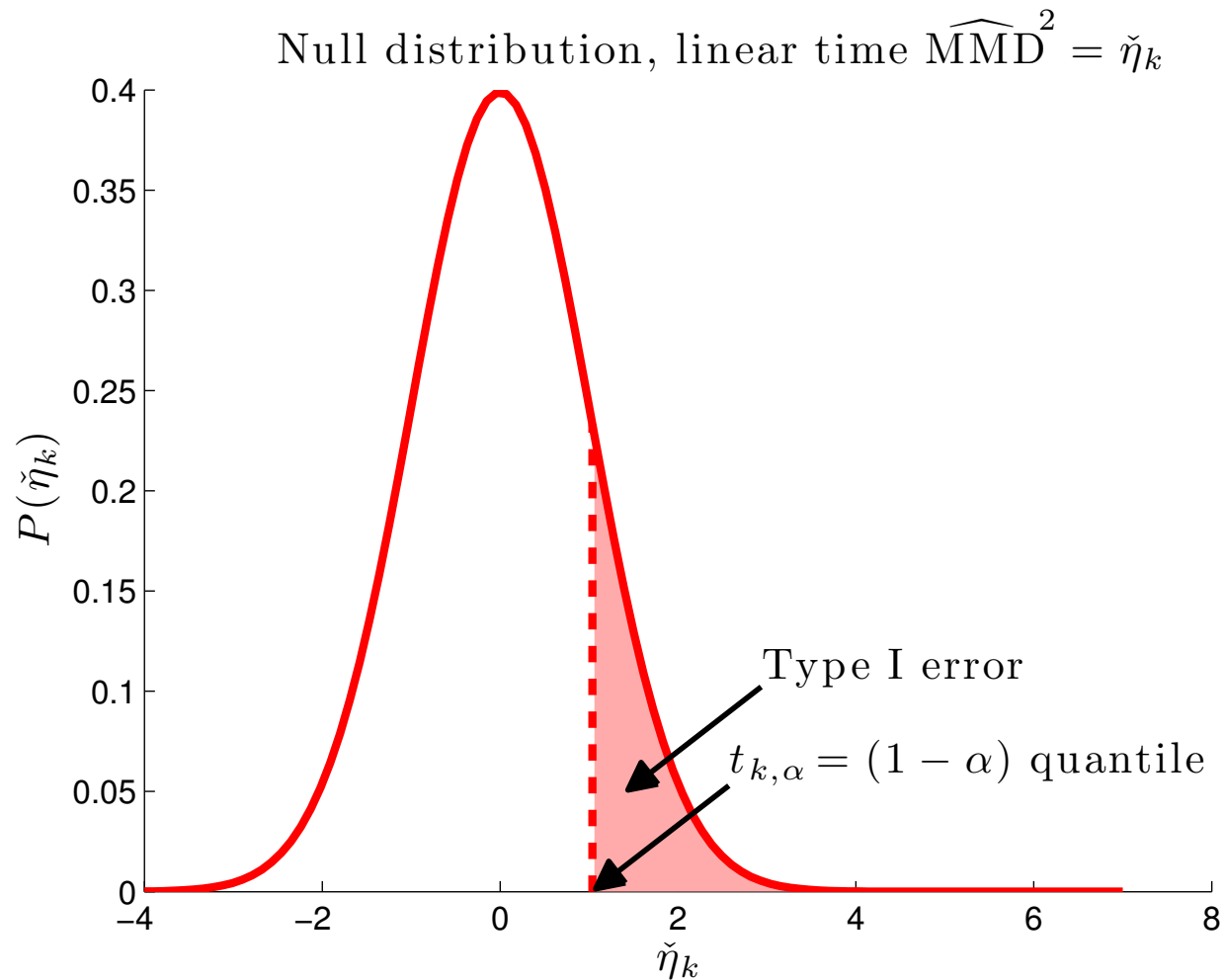
$$m^{1/2} (\check{\eta}_k - \eta_k(p, q)) \xrightarrow{D} \mathcal{N}(0, 2\sigma_k^2)$$

- assuming $0 < \mathbb{E}(h_k^2) < \infty$ (true for bounded k)
- $\sigma_k^2 = \mathbb{E}_v h_k^2(v) - [\mathbb{E}_v(h_k(v))]^2$.

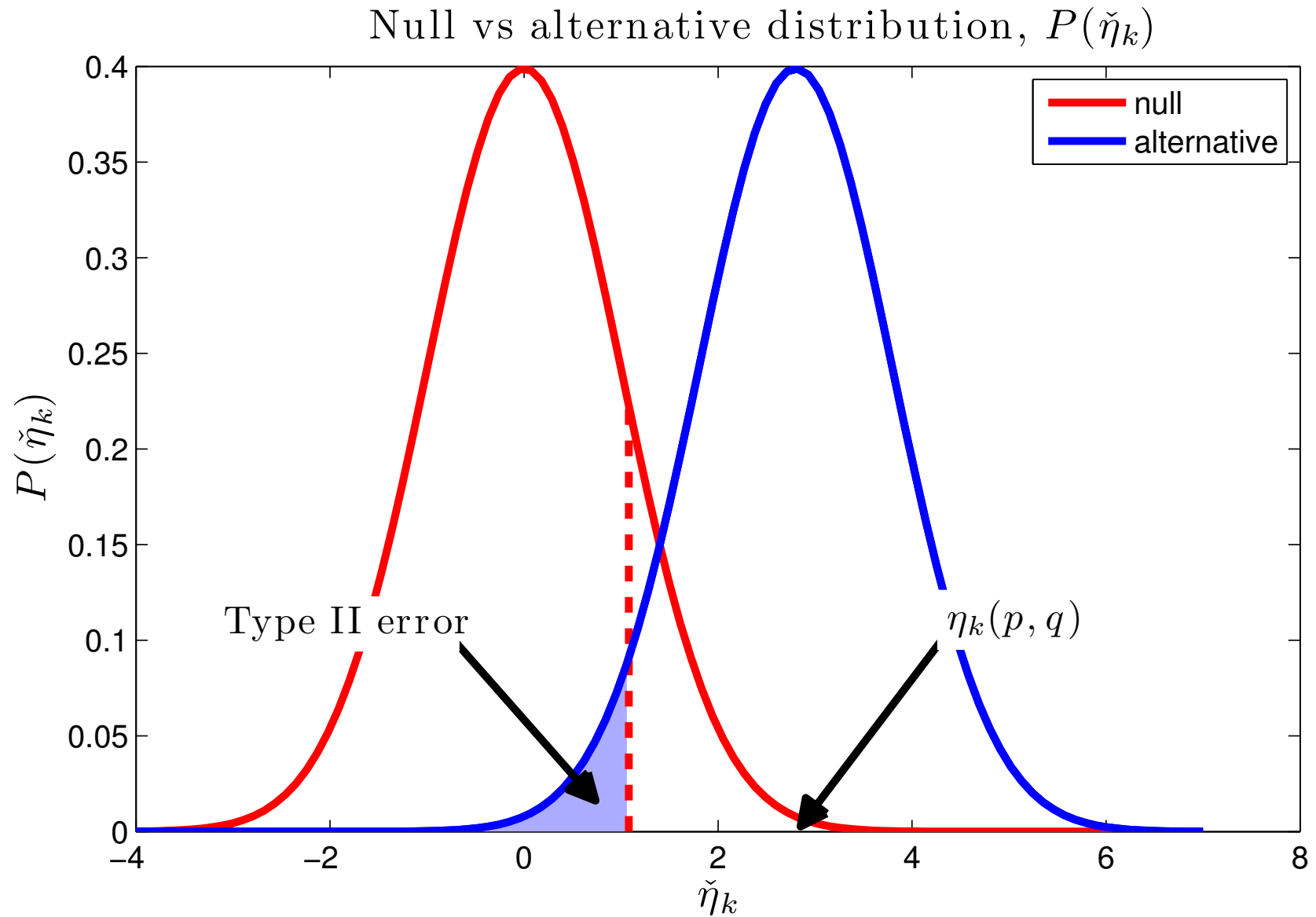
Hypothesis test

Hypothesis test of asymptotic level α :

$$t_{k,\alpha} = m^{-1/2} \sigma_k \sqrt{2} \Phi^{-1}(1 - \alpha) \quad \text{where } \Phi^{-1} \text{ is inverse CDF of } \mathcal{N}(0, 1).$$



Type II error



The best kernel: minimizes Type II error

Type II error: $\check{\eta}_k$ falls below the threshold $t_{k,\alpha}$ and $\eta_k(p, q) > 0$.

Prob. of a Type II error:

$$P(\check{\eta}_k < t_{k,\alpha}) = \Phi \left(\Phi^{-1}(1 - \alpha) - \frac{\eta_k(p, q) \sqrt{m}}{\sigma_k \sqrt{2}} \right)$$

where Φ is a Normal CDF.

The best kernel: minimizes Type II error

Type II error: $\check{\eta}_k$ falls below the threshold $t_{k,\alpha}$ and $\eta_k(p, q) > 0$.

Prob. of a Type II error:

$$P(\check{\eta}_k < t_{k,\alpha}) = \Phi \left(\Phi^{-1}(1 - \alpha) - \frac{\eta_k(p, q) \sqrt{m}}{\sigma_k \sqrt{2}} \right)$$

where Φ is a Normal CDF.

Since Φ monotonic, **best kernel choice to minimize Type II error prob.** is:

$$k_* = \arg \max_{k \in \mathcal{K}} \eta_k(p, q) \sigma_k^{-1},$$

where \mathcal{K} is the family of kernels under consideration.

Learning the best kernel in a family

Define the family of kernels as follows:

$$\mathcal{K} := \left\{ k : k = \sum_{u=1}^d \beta_u k_u, \|\beta\|_1 = D, \beta_u \geq 0, \forall u \in \{1, \dots, d\} \right\}.$$

Properties: if at least one $\beta_u > 0$

- all $k \in \mathcal{K}$ are valid kernels,
- If all k_u characteristic then k characteristic

Test statistic

The squared MMD becomes

$$\eta_k(p, q) = \|\mu_k(p) - \mu_k(q)\|_{\mathcal{F}_k}^2 = \sum_{u=1}^d \beta_u \eta_u(p, q),$$

where $\eta_u(p, q) := \mathbb{E}_v h_u(v)$.

Test statistic

The squared MMD becomes

$$\eta_k(p, q) = \|\mu_k(p) - \mu_k(q)\|_{\mathcal{F}_k}^2 = \sum_{u=1}^d \beta_u \eta_u(p, q),$$

where $\eta_u(p, q) := \mathbb{E}_v h_u(v)$.

Denote:

- $\beta = (\beta_1, \beta_2, \dots, \beta_d)^\top \in \mathbb{R}^d$,
- $h = (h_1, h_2, \dots, h_d)^\top \in \mathbb{R}^d$,
 - $h_u(x, x', y, y') = k_u(x, x') + k_u(y, y') - k_u(x, y') - k_u(x', y)$
- $\eta = \mathbb{E}_v(h) = (\eta_1, \eta_2, \dots, \eta_d)^\top \in \mathbb{R}^d$.

Quantities for test:

$$\eta_k(p, q) = \mathbb{E}(\beta^\top h) = \beta^\top \eta \quad \sigma_k^2 := \beta^\top \text{cov}(h) \beta.$$

Optimization of ratio $\eta_k(p, q)\sigma_k^{-1}$

Empirical test parameters:

$$\hat{\eta}_k = \beta^\top \hat{\eta} \qquad \hat{\sigma}_{k,\lambda} = \sqrt{\beta^\top (\hat{Q} + \lambda_m I) \beta},$$

\hat{Q} is empirical estimate of $\text{cov}(h)$.

Note: $\hat{\eta}_k, \hat{\sigma}_{k,\lambda}$ computed on training data, vs $\check{\eta}_k, \check{\sigma}_k$ on data to be tested
(why?)

Optimization of ratio $\eta_k(p, q)\sigma_k^{-1}$

Empirical test parameters:

$$\hat{\eta}_k = \beta^\top \hat{\eta} \quad \hat{\sigma}_{k,\lambda} = \sqrt{\beta^\top \left(\hat{Q} + \lambda_m I \right) \beta},$$

\hat{Q} is empirical estimate of $\text{cov}(h)$.

Note: $\hat{\eta}_k, \hat{\sigma}_{k,\lambda}$ computed on training data, vs $\check{\eta}_k, \check{\sigma}_k$ on data to be tested
(why?)

Objective:

$$\begin{aligned} \hat{\beta}^* &= \arg \max_{\beta \succeq 0} \hat{\eta}_k(p, q) \hat{\sigma}_{k,\lambda}^{-1} \\ &= \arg \max_{\beta \succeq 0} \left(\beta^\top \hat{\eta} \right) \left(\beta^\top \left(\hat{Q} + \lambda_m I \right) \beta \right)^{-1/2} \\ &=: \alpha(\beta; \hat{\eta}, \hat{Q}) \end{aligned}$$

Optimization of ratio $\eta_k(p, q)\sigma_k^{-1}$

Assume: $\hat{\eta}$ has at least one positive entry

Then there exists $\beta \succeq 0$ s.t. $\alpha(\beta; \hat{\eta}, \hat{Q}) > 0$.

Thus: $\alpha(\hat{\beta}^*; \hat{\eta}, \hat{Q}) > 0$

Optimization of ratio $\eta_k(p, q)\sigma_k^{-1}$

Assume: $\hat{\eta}$ has at least one positive entry

Then there exists $\beta \succeq 0$ s.t. $\alpha(\beta; \hat{\eta}, \hat{Q}) > 0$.

Thus: $\alpha(\hat{\beta}^*; \hat{\eta}, \hat{Q}) > 0$

Solve easier problem: $\hat{\beta}^* = \arg \max_{\beta \succeq 0} \alpha^2(\beta; \hat{\eta}, \hat{Q})$.

Quadratic program:

$$\min\{\beta^\top (\hat{Q} + \lambda_m I) \beta : \beta^\top \hat{\eta} = 1, \beta \succeq 0\}$$

Optimization of ratio $\eta_k(p, q)\sigma_k^{-1}$

Assume: $\hat{\eta}$ has at least one positive entry

Then there exists $\beta \succeq 0$ s.t. $\alpha(\beta; \hat{\eta}, \hat{Q}) > 0$.

Thus: $\alpha(\hat{\beta}^*; \hat{\eta}, \hat{Q}) > 0$

Solve easier problem: $\hat{\beta}^* = \arg \max_{\beta \succeq 0} \alpha^2(\beta; \hat{\eta}, \hat{Q})$.

Quadratic program:

$$\min\{\beta^\top (\hat{Q} + \lambda_m I) \beta : \beta^\top \hat{\eta} = 1, \beta \succeq 0\}$$

What if $\hat{\eta}$ has no positive entries?

Test procedure

1. Split the data into **testing** and **training**.
2. On the **training** data:
 - (a) Compute $\hat{\eta}_u$ for all $k_u \in \mathcal{K}$
 - (b) If at least one $\hat{\eta}_u > 0$, solve the QP to get β^* , else choose random kernel from \mathcal{K}
3. On the **test** data:
 - (a) Compute $\check{\eta}_{k^*}$ using $k^* = \sum_{u=1}^d \beta^* k_u$
 - (b) Compute test threshold \check{t}_{α, k^*} using $\check{\sigma}_{k^*}$
4. Reject null if $\check{\eta}_{k^*} > \check{t}_{\alpha, k^*}$

Convergence bounds

Assume bounded kernel, σ_k , bounded away from 0.

If $\lambda_m = \Theta(m^{-1/3})$ then

$$\left| \sup_{k \in \mathcal{K}} \hat{\eta}_k \hat{\sigma}_{k,\lambda}^{-1} - \sup_{k \in \mathcal{K}} \eta_k \sigma_k^{-1} \right| = O_P \left(m^{-1/3} \right).$$

Convergence bounds

Assume bounded kernel, σ_k , bounded away from 0.

If $\lambda_m = \Theta(m^{-1/3})$ then

$$\left| \sup_{k \in \mathcal{K}} \hat{\eta}_k \hat{\sigma}_{k,\lambda}^{-1} - \sup_{k \in \mathcal{K}} \eta_k \sigma_k^{-1} \right| = O_P \left(m^{-1/3} \right).$$

Idea:

$$\begin{aligned} & \left| \sup_{k \in \mathcal{K}} \hat{\eta}_k \hat{\sigma}_{k,\lambda}^{-1} - \sup_{k \in \mathcal{K}} \eta_k \sigma_k^{-1} \right| \\ & \leq \sup_{k \in \mathcal{K}} \left| \hat{\eta}_k \hat{\sigma}_{k,\lambda}^{-1} - \eta_k \sigma_{k,\lambda}^{-1} \right| + \sup_{k \in \mathcal{K}} \left| \eta_k \sigma_{k,\lambda}^{-1} - \eta_k \sigma_k^{-1} \right| \\ & \leq \frac{\sqrt{d}}{D\sqrt{\lambda_m}} \left(C_1 \sup_{k \in \mathcal{K}} |\hat{\eta}_k - \eta_k| + C_2 \sup_{k \in \mathcal{K}} |\hat{\sigma}_{k,\lambda} - \sigma_{k,\lambda}| \right) + C_3 D^2 \lambda_m, \end{aligned}$$

Experiments

Competing approaches

- Median heuristic
- Max. MMD: choose $k_u \in \mathcal{K}$ with the largest $\hat{\eta}_u$
 - same as maximizing $\beta^\top \hat{\eta}$ subject to $\|\beta\|_1 \leq 1$
- ℓ_2 statistic: maximize $\beta^\top \hat{\eta}$ subject to $\|\beta\|_2 \leq 1$
- Cross validation on training set

Also compare with:

- **Single kernel** that maximizes ratio $\eta_k(p, q) \sigma_k^{-1}$

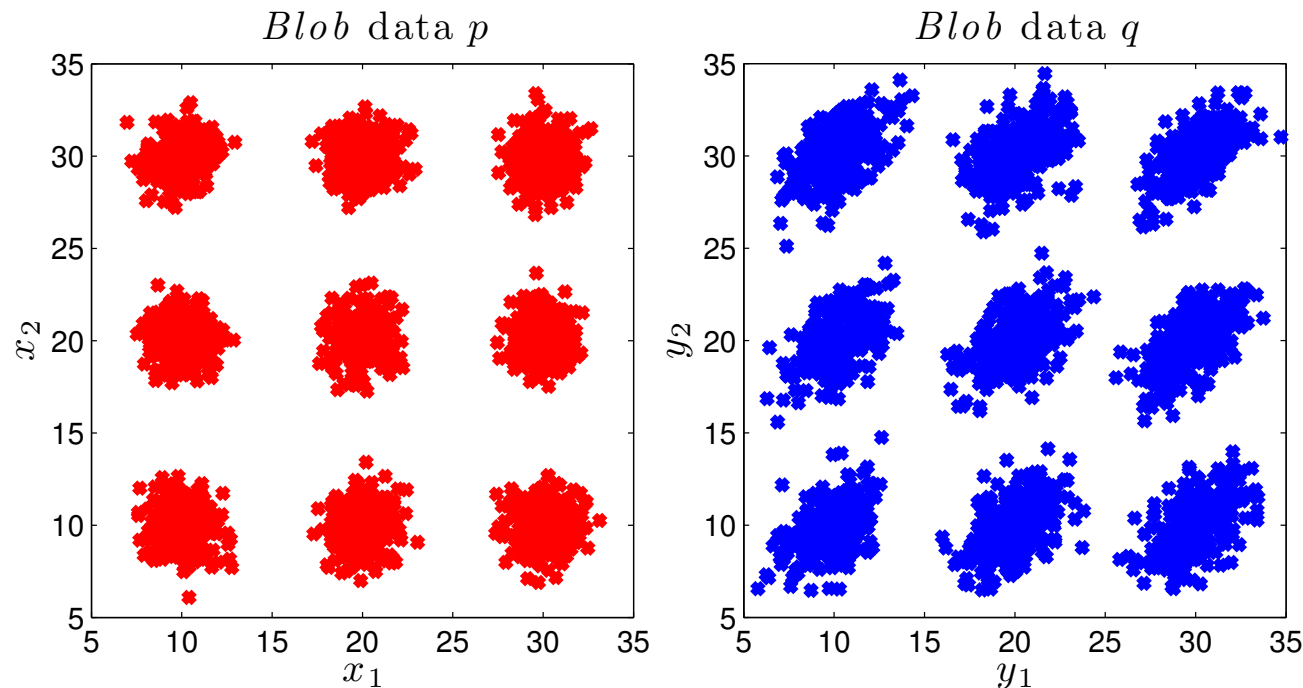
Blobs: data

Difficult problems: lengthscale of the *difference* in distributions not the same as that of the distributions.

Blobs: data

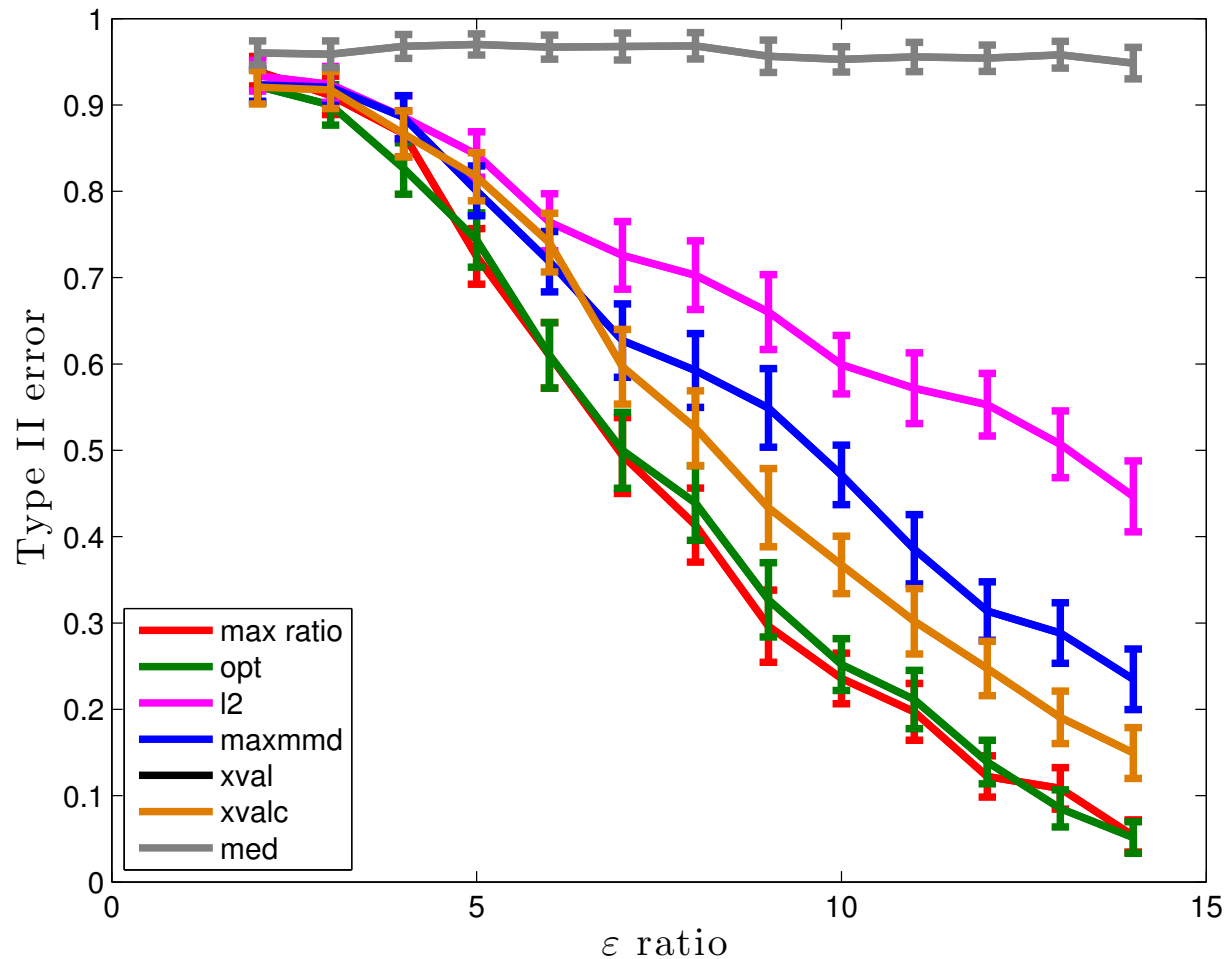
Difficult problems: lengthscale of the *difference* in distributions not the same as that of the distributions.

We distinguish a field of Gaussian blobs with different covariances.



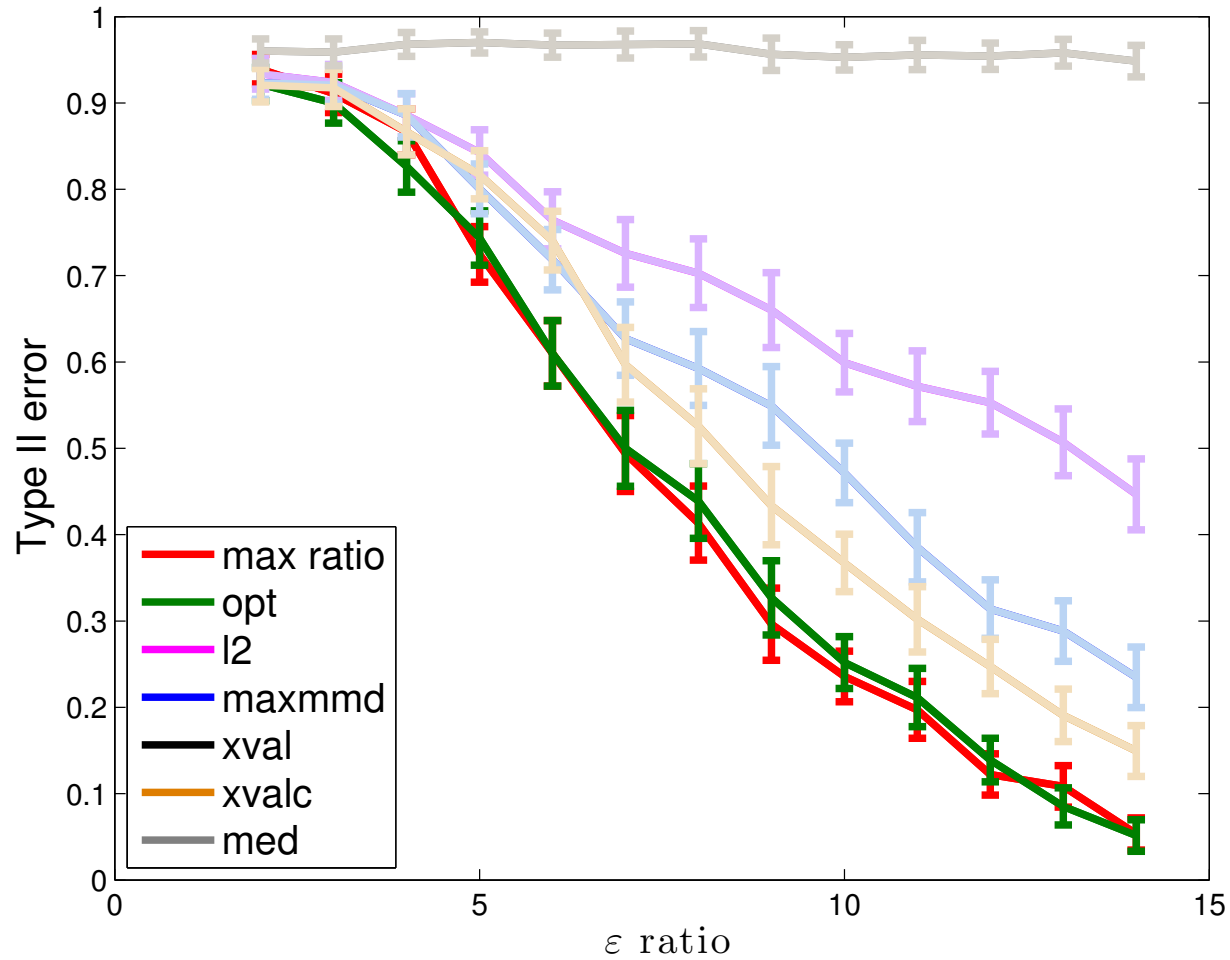
Ratio $\varepsilon = 3.2$ of largest to smallest eigenvalues of blobs in q .

Blobs: results



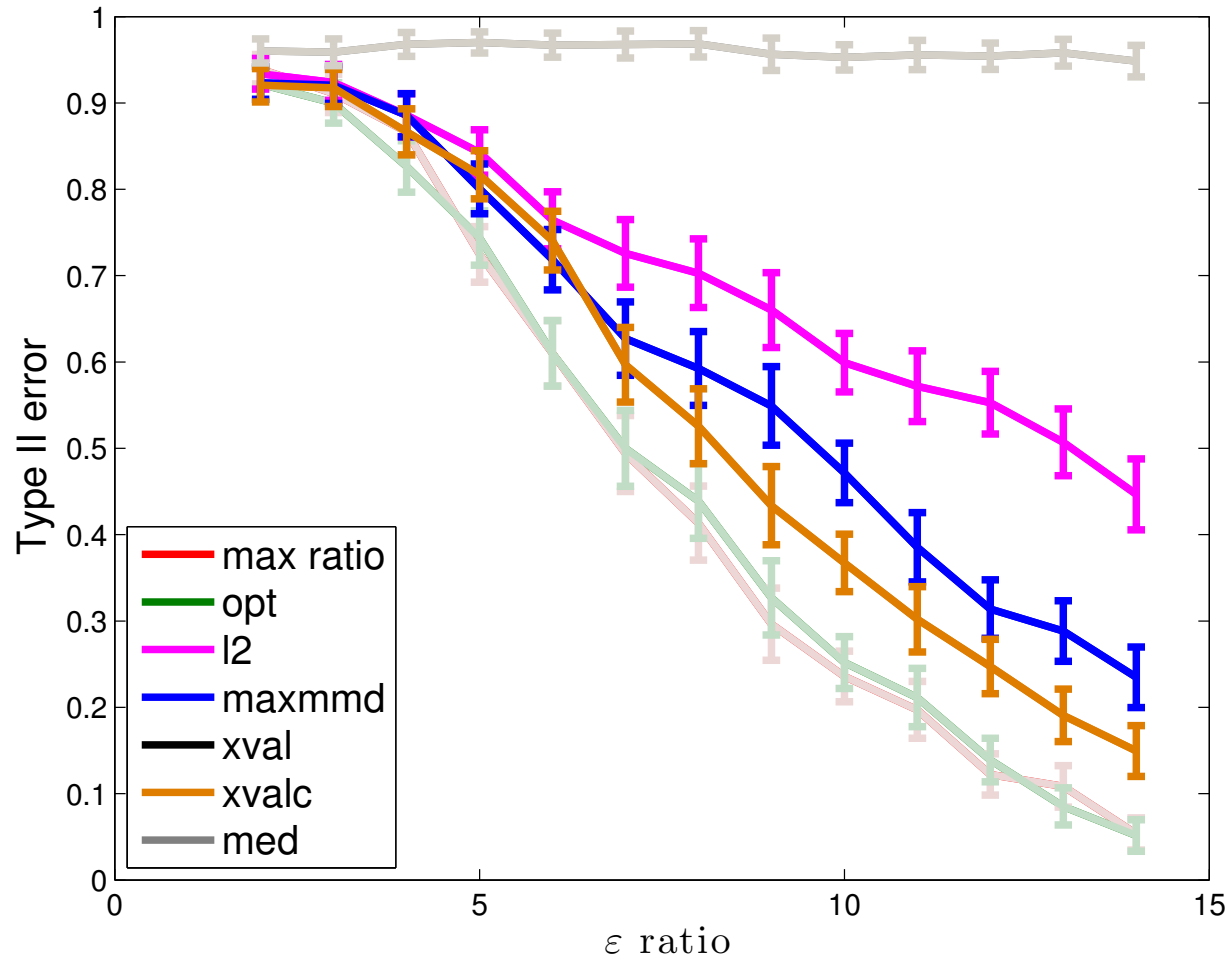
Parameters: $m = 10,000$ (for training and test). **Ratio ϵ** of largest to smallest eigenvalues of blobs in q . Results are average over 617 trials.

Blobs: results



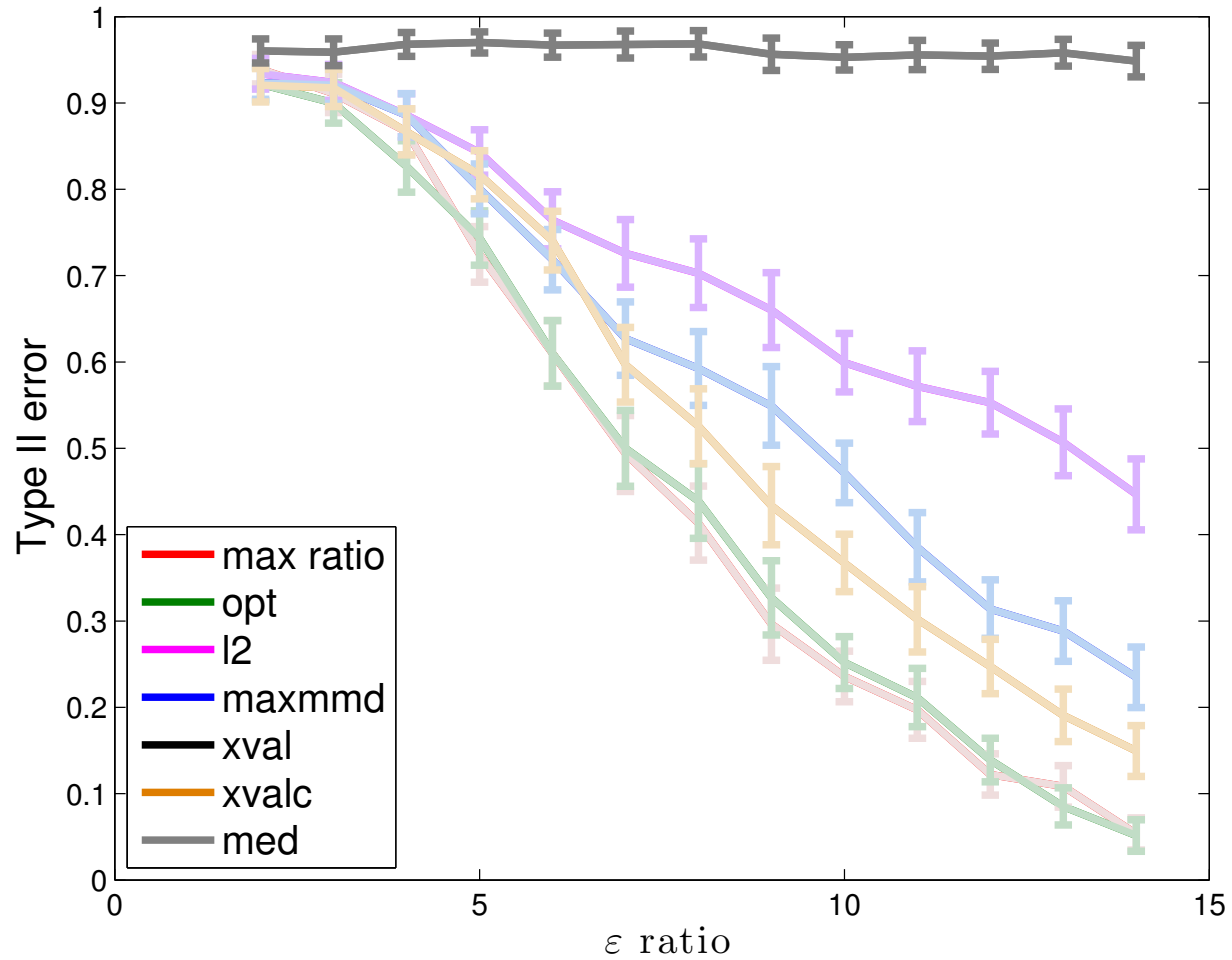
Optimize ratio $\eta_k(p, q)\sigma_k^{-1}$

Blobs: results



Maximize $\eta_k(p, q)$ with β constraint

Blobs: results



Median heuristic

Feature selection: data

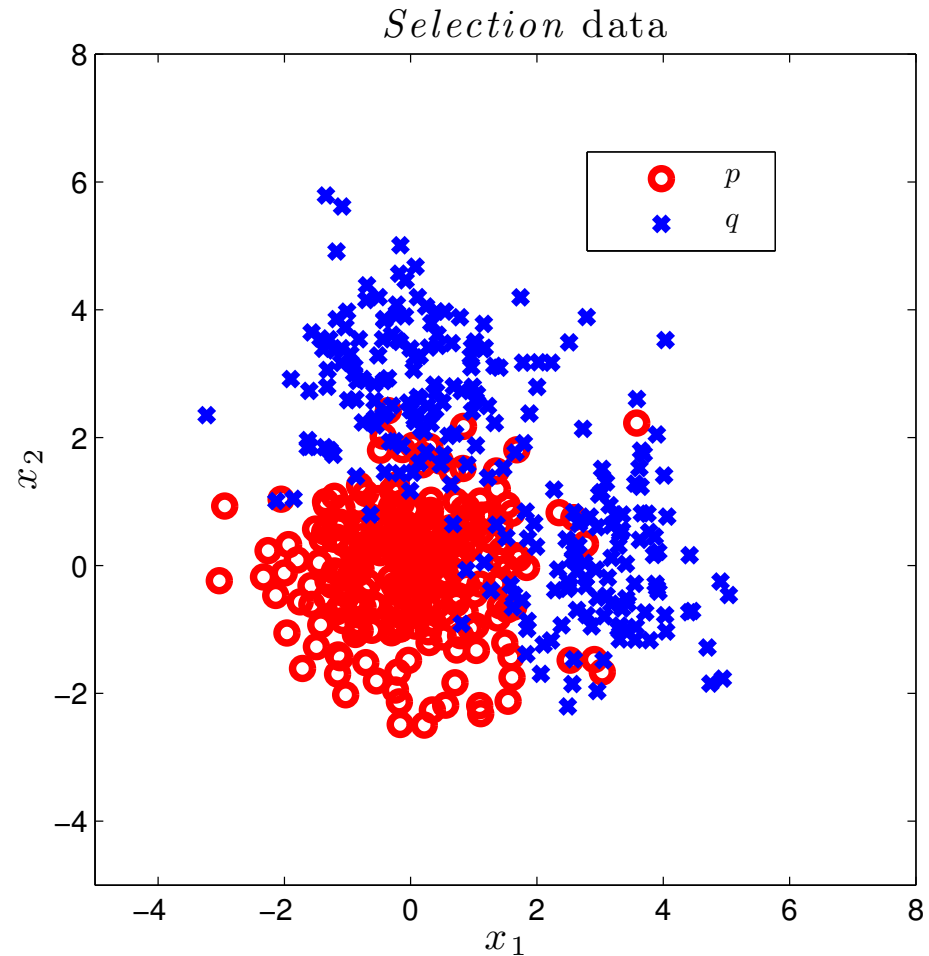
Idea: no single best kernel.

Each of the k_u are univariate (along a single coordinate)

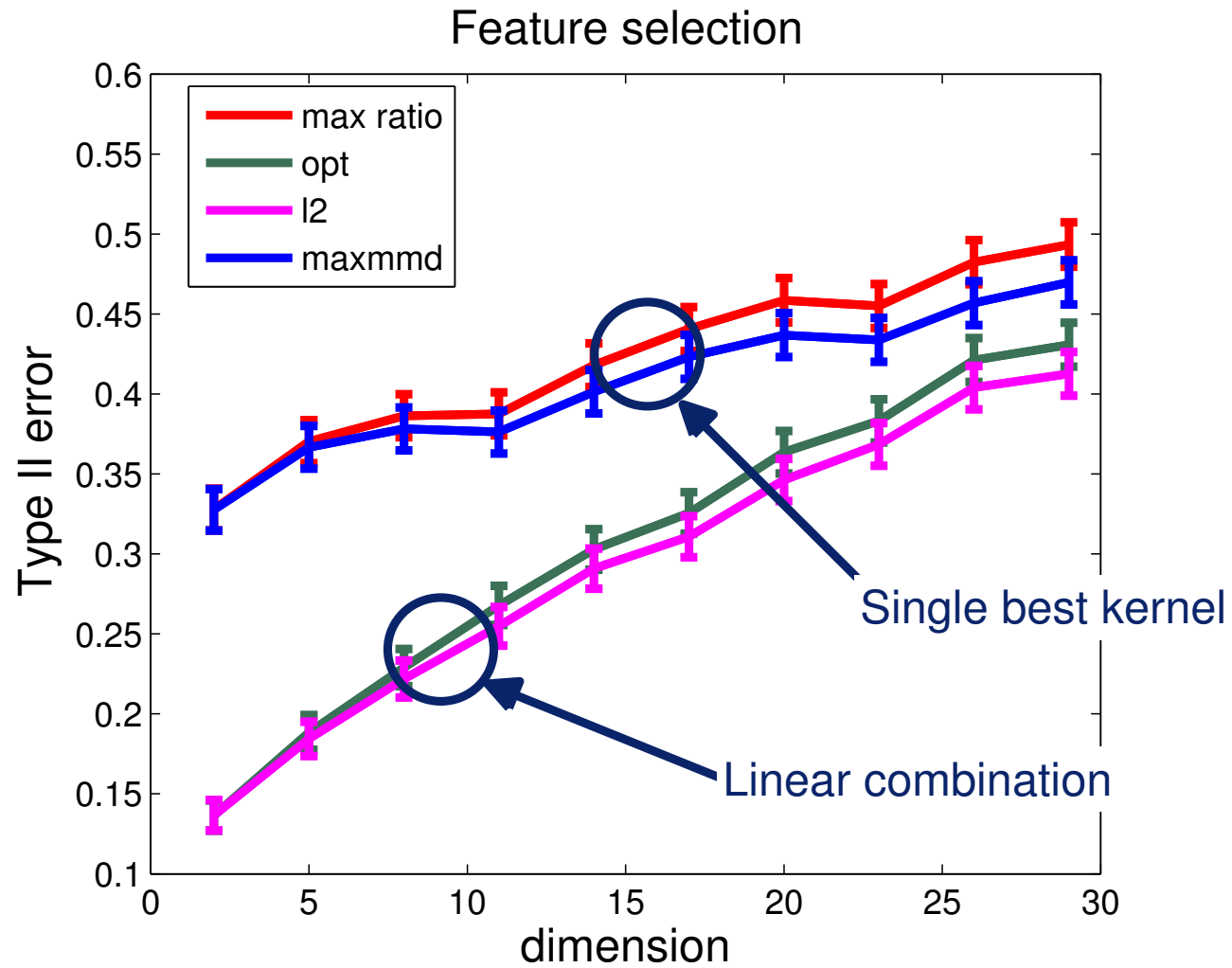
Feature selection: data

Idea: no single best kernel.

Each of the k_u are univariate (along a single coordinate)



Feature selection: results



$m = 10,000$, average over 5000 trials

Amplitude modulated signals

Given an audio signal $s(t)$, an amplitude modulated signal can be defined

$$u(t) = \sin(\omega_c t) [a s(t) + l]$$

- ω_c : carrier frequency
- $a = 0.2$ is signal scaling, $l = 2$ is offset

Amplitude modulated signals

Given an audio signal $s(t)$, an amplitude modulated signal can be defined

$$u(t) = \sin(\omega_c t) [a s(t) + l]$$

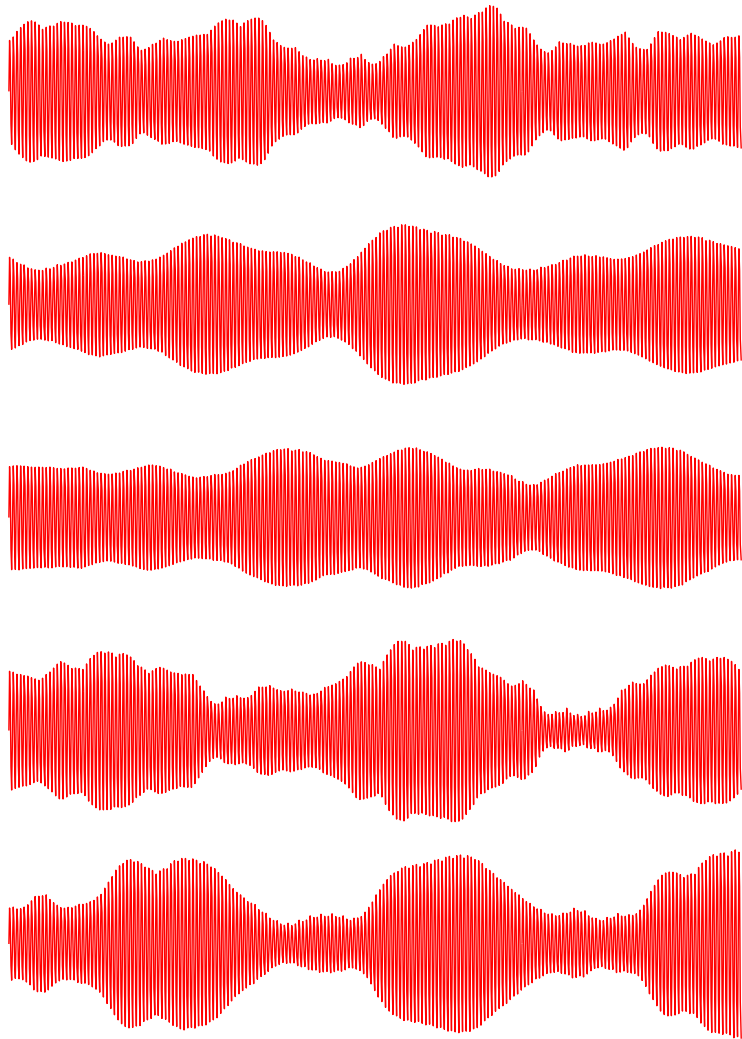
- ω_c : carrier frequency
- $a = 0.2$ is signal scaling, $l = 2$ is offset

Two amplitude modulated signals from same artist (in this case, Magnetic Fields).

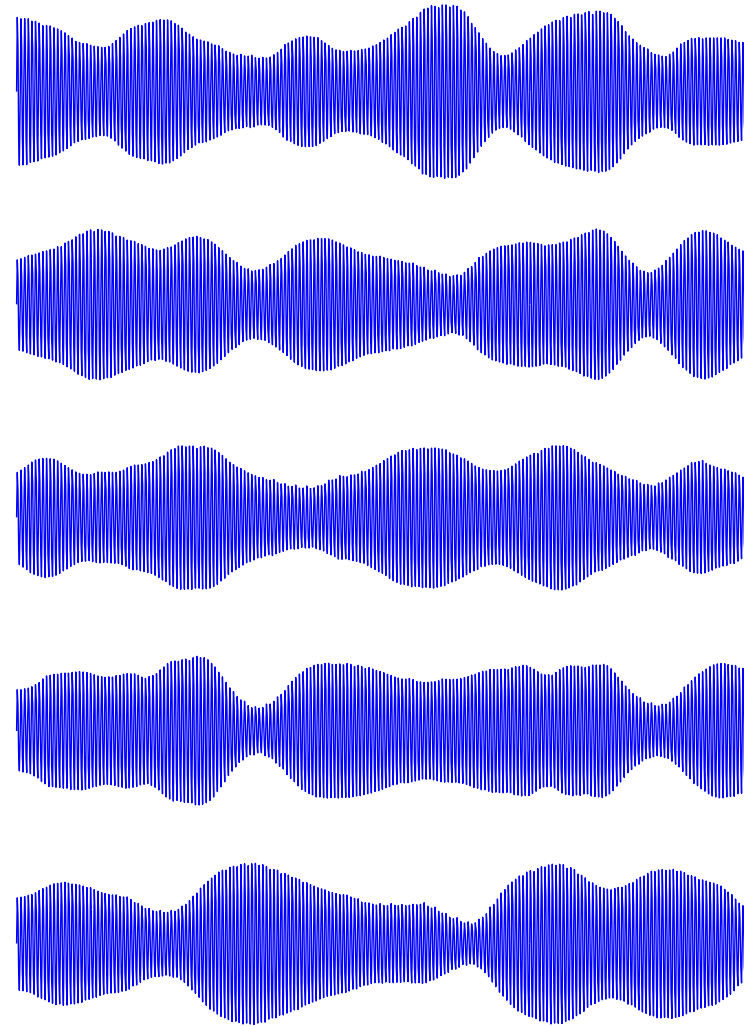
- Music sampled at 8KHz (**very low**)
- Carrier frequency is 24kHz
- AM signal observed at 120kHz
- Samples are extracts of length $N = 1000$, approx. 0.01 sec (**very short**).
- Total dataset size is 30,000 samples from each of p, q .

Amplitude modulated signals

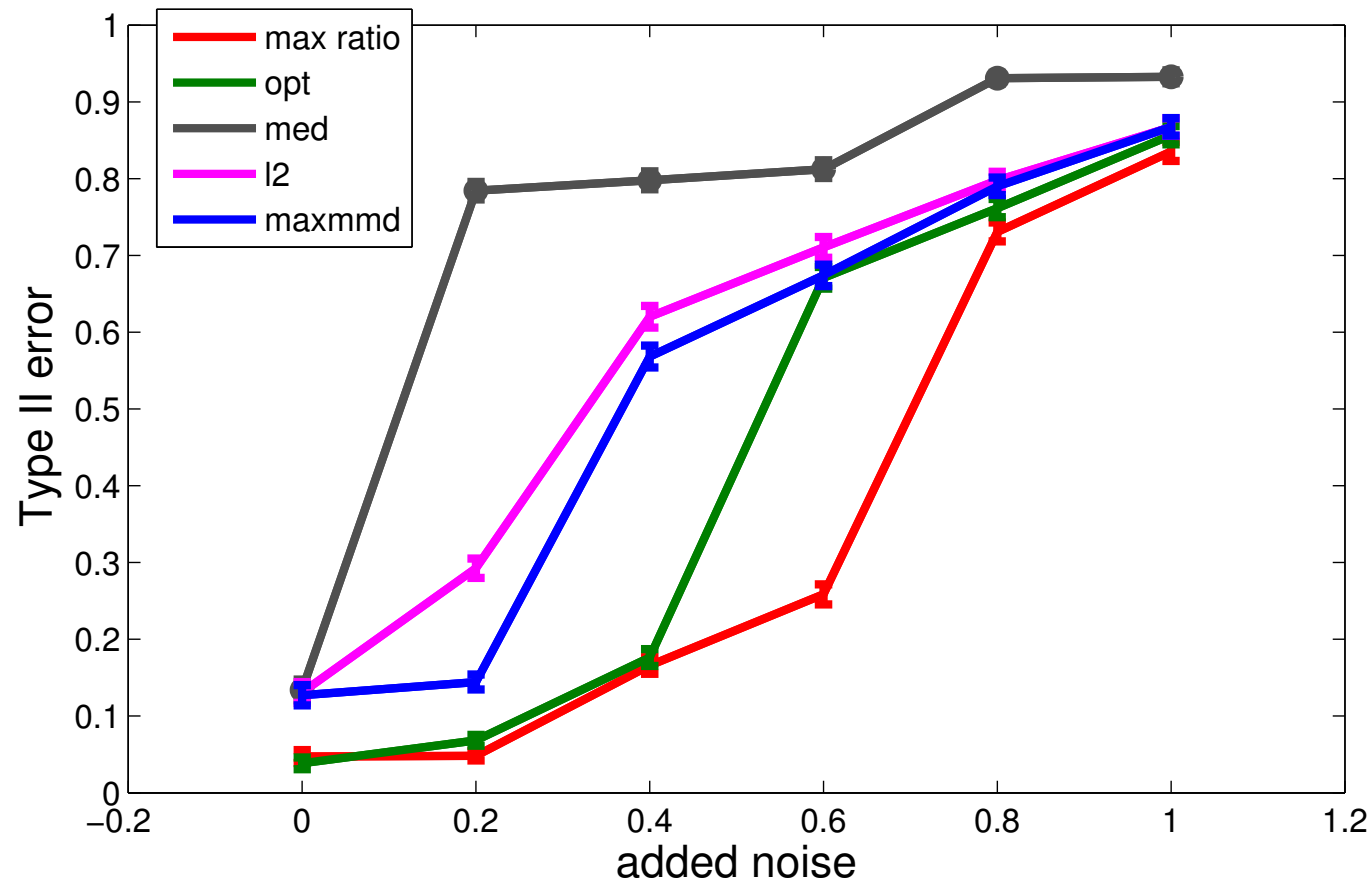
Samples from P



Samples from Q



Results: AM signals



$m = 10,000$ (for training and test) and scaling $a = 0.5$. Average over 4124 trials. Gaussian noise added.

Observations on kernel choice

- It is possible to choose the best kernel for a kernel two-sample test
- Kernel choice matters for “difficult” problems, where the distributions differ on a lengthscale different to that of the data.
- Ongoing work:
 - quadratic time statistic
 - avoid training/test split

Summary

- **MMD** a distance between distributions [ISMB06, NIPS06a, JMLR10, JMLR12a]
 - high dimensionality
 - non-euclidean data (strings, graphs)
 - Nonparametric hypothesis tests
- Measure and test **independence** [ALT05, NIPS07a, NIPS07b, ALT08, JMLR10, JMLR12a]
- **Characteristic RKHS**: MMD a metric [NIPS07b, COLT08, NIPS08a]
 - Easy to check: does spectrum cover \mathbb{R}^d

Co-authors

- **From UCL:**

- Luca Baldassarre
- Steffen Grunewalder
- Guy Lever
- Sam Patterson
- Massimiliano Pontil
- Dino Sejdinovic

- **External:**

- Karsten Borgwardt, MPI
- Wicher Bergsma, LSE
- Kenji Fukumizu, ISM
- Zaid Harchaoui, INRIA
- Bernhard Schoelkopf, MPI
- Alex Smola, CMU/Google
- Le Song, Georgia Tech
- Bharath Sriperumbudur, Cambridge



Selected references

Characteristic kernels and mean embeddings:

- Smola, A., Gretton, A., Song, L., Schoelkopf, B. (2007). A hilbert space embedding for distributions. ALT.
- Sriperumbudur, B., Gretton, A., Fukumizu, K., Schoelkopf, B., Lanckriet, G. (2010). Hilbert space embeddings and metrics on probability measures. JMLR.
- Gretton, A., Borgwardt, K., Rasch, M., Schoelkopf, B., Smola, A. (2012). A kernel two- sample test. JMLR.

Two-sample, independence, conditional independence tests:

- Gretton, A., Fukumizu, K., Teo, C., Song, L., Schoelkopf, B., Smola, A. (2008). A kernel statistical test of independence. NIPS
- Fukumizu, K., Gretton, A., Sun, X., Schoelkopf, B. (2008). Kernel measures of conditional dependence.
- Gretton, A., Fukumizu, K., Harchaoui, Z., Sriperumbudur, B. (2009). A fast, consistent kernel two-sample test. NIPS.
- Gretton, A., Borgwardt, K., Rasch, M., Schoelkopf, B., Smola, A. (2012). A kernel two- sample test. JMLR

Energy distance, relation to kernel distances

- Sejdinovic, D., Sriperumbudur, B., Gretton, A., Fukumizu, K., (2013). Equivalence of distance-based and rkhs-based statistics in hypothesis testing. Annals of Statistics.

Three way interaction

- Sejdinovic, D., Gretton, A., and Bergsma, W. (2013). A Kernel Test for Three-Variable Interactions. NIPS.

Selected references (continued)

Conditional mean embedding, RKHS-valued regression:

- Weston, J., Chapelle, O., Elisseeff, A., Schölkopf, B., and Vapnik, V., (2003). Kernel Dependency Estimation, NIPS.
- Micchelli, C., and Pontil, M., (2005). On Learning Vector-Valued Functions. Neural Computation.
- Caponnetto, A., and De Vito, E. (2007). Optimal Rates for the Regularized Least-Squares Algorithm. Foundations of Computational Mathematics.
- Song, L., and Huang, J., and Smola, A., Fukumizu, K., (2009). Hilbert Space Embeddings of Conditional Distributions. ICML.
- Grunewalder, S., Lever, G., Baldassarre, L., Patterson, S., Gretton, A., Pontil, M. (2012). Conditional mean embeddings as regressors. ICML.
- Grunewalder, S., Gretton, A., Shawe-Taylor, J. (2013). Smooth operators. ICML.

Kernel Bayes rule:

- Song, L., Fukumizu, K., Gretton, A. (2013). Kernel embeddings of conditional distributions: A unified kernel framework for nonparametric inference in graphical models. IEEE Signal Processing Magazine.
- Fukumizu, K., Song, L., Gretton, A. (2013). Kernel Bayes rule: Bayesian inference with positive definite kernels, JMLR

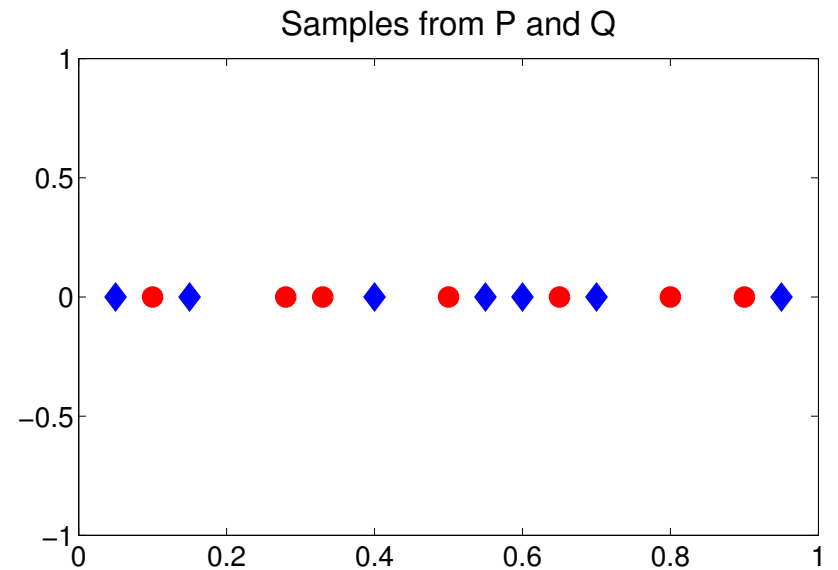
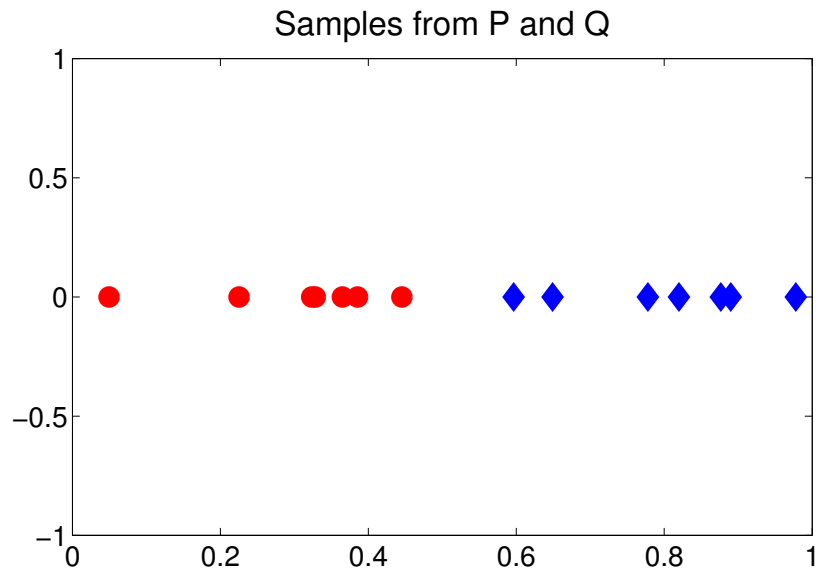
Local departures from the null

What is a hard testing problem?

Local departures from the null

What is a hard testing problem?

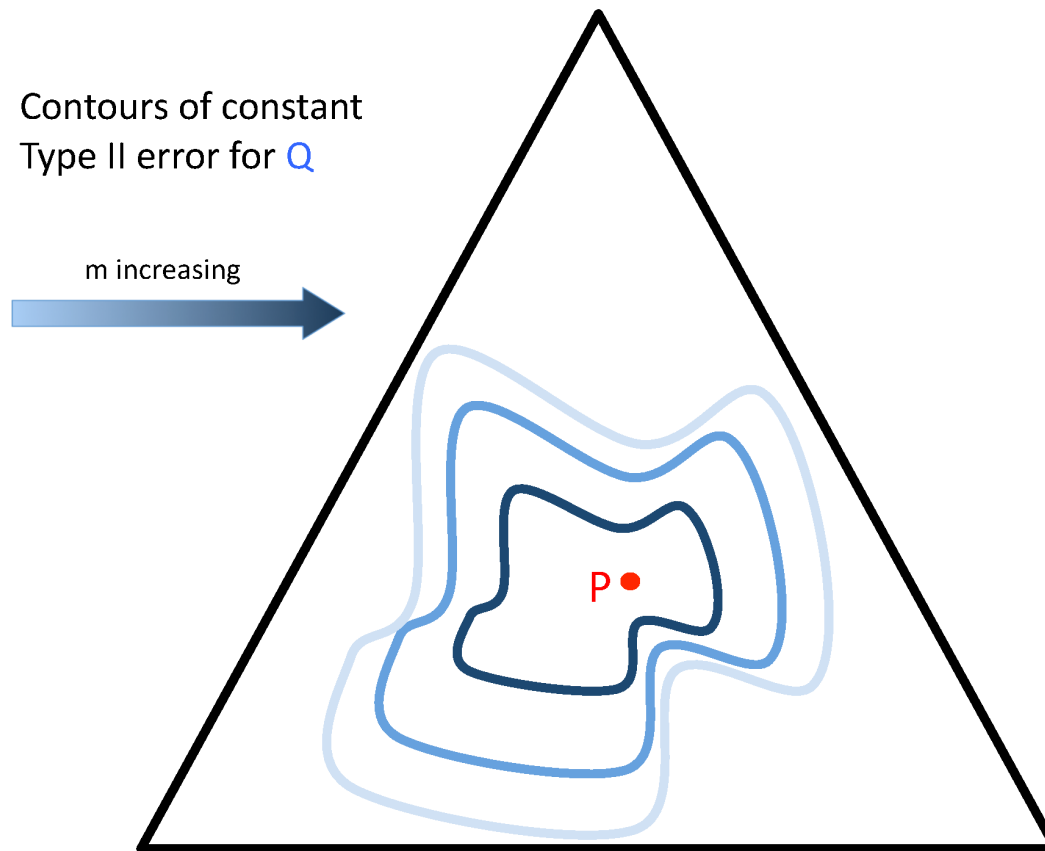
- **First version:** for fixed m , “closer” **P** and **Q** have **higher** Type II error



Local departures from the null

What is a hard testing problem?

- As m increases, distinguish “closer” P and Q with fixed Type II error



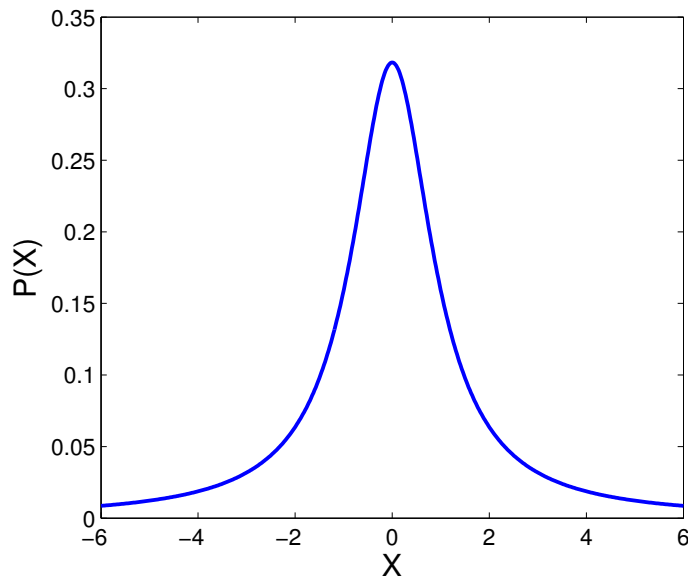
Local departures from the null

What is a hard testing problem?

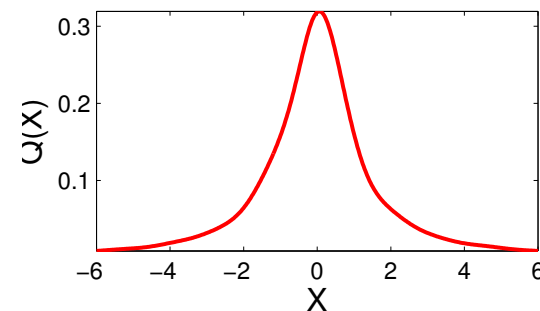
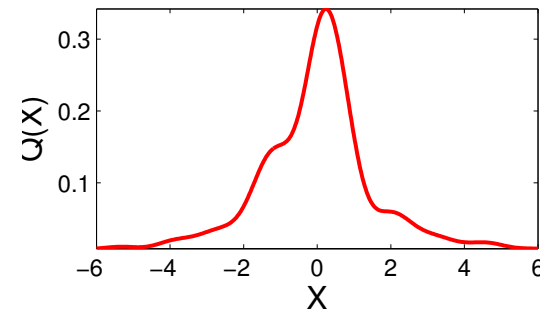
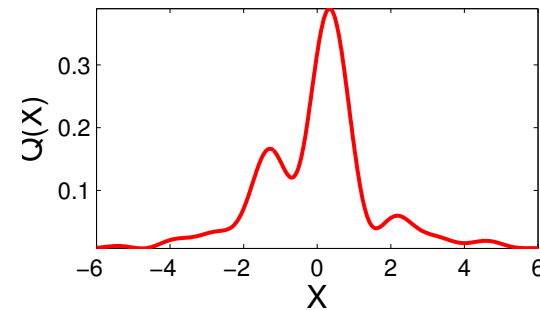
- As m increases, distinguish “closer” \mathbf{P} and \mathbf{Q} with fixed Type II error
- **Example:** $f_{\mathbf{P}}$ and $f_{\mathbf{Q}}$ probability densities, $f_{\mathbf{Q}} = f_{\mathbf{P}} + \delta g$, where $\delta \in \mathbb{R}$, g some *fixed* function such that $f_{\mathbf{Q}}$ is a valid density
 - If $\delta \sim m^{-1/2}$, Type II error approaches a constant

More general local departures from null

- **Example:** $f_{\mathbf{P}}$ and $f_{\mathbf{Q}}$ probability densities, $f_{\mathbf{Q}} = f_{\mathbf{P}} + \delta g$, where $\delta \in \mathbb{R}$, g some *fixed* function such that $f_{\mathbf{Q}}$ is a valid density



VS



Local departures from the null

What is a hard testing problem?

- As we see more samples m , distinguish “closer” \mathbf{P} and \mathbf{Q} with same Type II error
- **Example:** $f_{\mathbf{P}}$ and $f_{\mathbf{Q}}$ probability densities, $f_{\mathbf{Q}} = f_{\mathbf{P}} + \delta g$, where $\delta \in \mathbb{R}$, g some *fixed* function such that $f_{\mathbf{Q}}$ is a valid density
 - If $\delta \sim m^{-1/2}$, Type II error approaches a constant
- ...but **other choices also possible** – how to characterize them all?

Local departures from the null

What is a hard testing problem?

- As we see more samples m , distinguish “closer” \mathbf{P} and \mathbf{Q} with same Type II error
- **Example:** $f_{\mathbf{P}}$ and $f_{\mathbf{Q}}$ probability densities, $f_{\mathbf{Q}} = f_{\mathbf{P}} + \delta g$, where $\delta \in \mathbb{R}$, g some *fixed* function such that $f_{\mathbf{Q}}$ is a valid density
 - If $\delta \sim m^{-1/2}$, Type II error approaches a constant
- ...but **other choices also possible** – how to characterize them all?

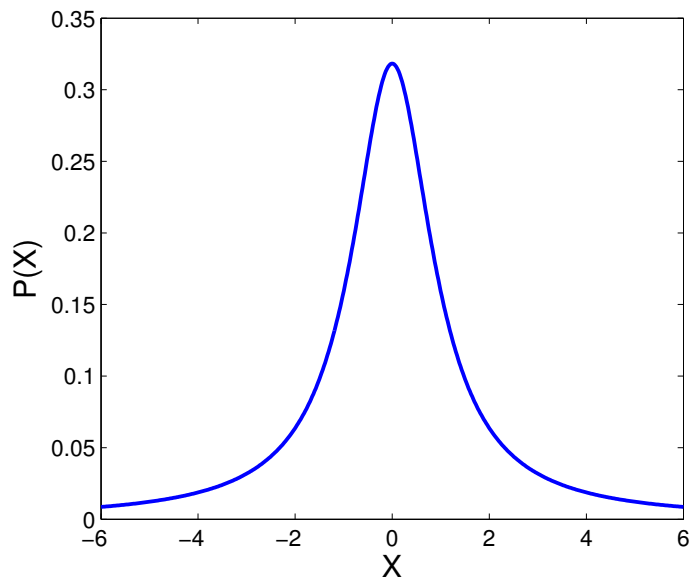
General characterization of local departures from \mathcal{H}_0 :

- Write $\mu_{\mathbf{Q}} = \mu_{\mathbf{P}} + g_m$, where $g_m \in \mathcal{F}$ chosen such that $\mu_{\mathbf{P}} + g_m$ a valid distribution embedding
- Minimum distinguishable distance [JMLR12]

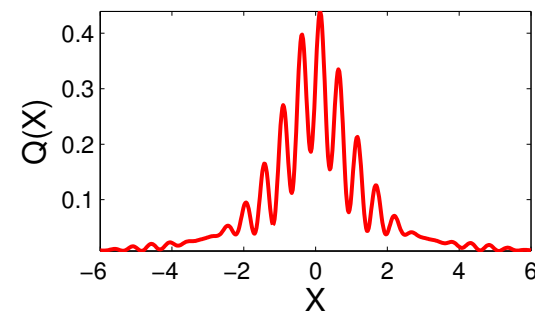
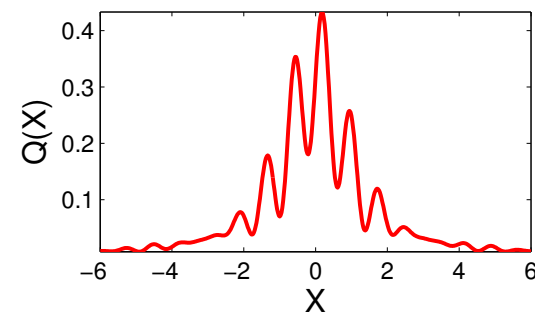
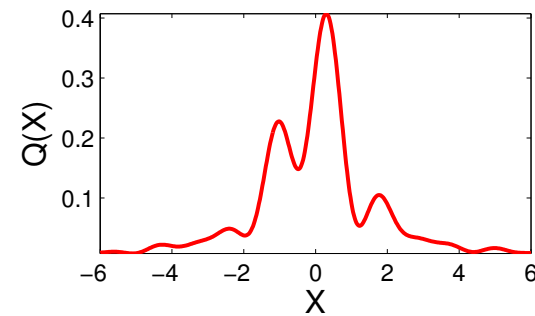
$$\|g_m\|_{\mathcal{F}} = cm^{-1/2}$$

More general local departures from null

- **More advanced example** of a local departure from the null
- Recall: $\mu_{\mathbf{Q}} = \mu_{\mathbf{P}} + g_m$, and $\|g_m\|_{\mathcal{F}} = cm^{-1/2}$



VS



Kernels vs kernels

- How does MMD relate to [Parzen density estimate](#)? [Anderson et al., 1994]

$$\hat{f}_{\mathbf{P}}(x) = \frac{1}{m} \sum_{i=1}^m \kappa(x_i - x), \text{ where } \kappa \text{ satisfies } \int_{\mathcal{X}} \kappa(x) dx = 1 \text{ and } \kappa(x) \geq 0.$$

Kernels vs kernels

- How does MMD relate to [Parzen density estimate](#)? [Anderson et al., 1994]

$$\hat{f}_{\mathbf{P}}(x) = \frac{1}{m} \sum_{i=1}^m \kappa(x_i - x), \text{ where } \kappa \text{ satisfies } \int_{\mathcal{X}} \kappa(x) dx = 1 \text{ and } \kappa(x) \geq 0.$$

- [L₂ distance](#) between Parzen density estimates:

$$\begin{aligned} D_2(\hat{f}_{\mathbf{P}}, \hat{f}_{\mathbf{Q}})^2 &= \int \left[\frac{1}{m} \sum_{i=1}^m \kappa(x_i - z) - \frac{1}{m} \sum_{i=1}^m \kappa(y_i - z) \right]^2 dz \\ &= \frac{1}{m^2} \sum_{i,j=1}^m k(x_i - x_j) + \frac{1}{m^2} \sum_{i,j=1}^m k(y_i - y_j) - \frac{2}{m^2} \sum_{i,j=1}^m k(x_i - y_j), \end{aligned}$$

where $k(x - y) = \int \kappa(x - z)\kappa(y - z)dz$

Kernels vs kernels

- How does MMD relate to **Parzen density estimate**? [Anderson et al., 1994]

$$\hat{f}_{\mathbf{P}}(x) = \frac{1}{m} \sum_{i=1}^m \kappa(x_i - x), \text{ where } \kappa \text{ satisfies } \int_{\mathcal{X}} \kappa(x) dx = 1 \text{ and } \kappa(x) \geq 0.$$

- **L_2 distance** between Parzen density estimates:

$$\begin{aligned} D_2(\hat{f}_{\mathbf{P}}, \hat{f}_{\mathbf{Q}})^2 &= \int \left[\frac{1}{m} \sum_{i=1}^m \kappa(x_i - z) - \frac{1}{m} \sum_{i=1}^m \kappa(y_i - z) \right]^2 dz \\ &= \frac{1}{m^2} \sum_{i,j=1}^m k(x_i - x_j) + \frac{1}{m^2} \sum_{i,j=1}^m k(y_i - y_j) - \frac{2}{m^2} \sum_{i,j=1}^m k(x_i - y_j), \end{aligned}$$

where $k(x - y) = \int \kappa(x - z)\kappa(y - z)dz$

- $f_{\mathbf{Q}} = f_{\mathbf{P}} + \delta g$, **minimum distance** to discriminate $f_{\mathbf{P}}$ from $f_{\mathbf{Q}}$ is $\delta = (m)^{-1/2} h_m^{-d/2}$, where h_m is width of κ .

Characteristic Kernels (via universality)

Characteristic: MMD a metric (MMD = 0 iff $\mathbf{P} = \mathbf{Q}$) [NIPS07b, COLT08]

Characteristic Kernels (via universality)

Characteristic: MMD a **metric** (MMD = 0 iff **P** = **Q**) [NIPS07b, COLT08]

Classical result: **P** = **Q** if and only if $\mathbb{E}_{\mathbf{P}}(f(x)) = \mathbb{E}_{\mathbf{Q}}(f(y))$ for all $f \in C(\mathcal{X})$,
the space of bounded continuous functions on \mathcal{X} [Dudley, 2002]

Characteristic Kernels (via universality)

Characteristic: MMD a **metric** (MMD = 0 iff $\mathbf{P} = \mathbf{Q}$) [NIPS07b, COLT08]

Classical result: $\mathbf{P} = \mathbf{Q}$ if and only if $\mathbb{E}_{\mathbf{P}}(f(x)) = \mathbb{E}_{\mathbf{Q}}(f(y))$ for all $f \in C(\mathcal{X})$, the space of bounded continuous functions on \mathcal{X} [Dudley, 2002]

Universal RKHS: $k(x, x')$ continuous, \mathcal{X} compact, and \mathcal{F} dense in $C(\mathcal{X})$ with respect to L_{∞} [Steinwart, 2001]

Characteristic Kernels (via universality)

Characteristic: MMD a **metric** (MMD = 0 iff $\mathbf{P} = \mathbf{Q}$) [NIPS07b, COLT08]

Classical result: $\mathbf{P} = \mathbf{Q}$ if and only if $\mathbb{E}_{\mathbf{P}}(f(x)) = \mathbb{E}_{\mathbf{Q}}(f(y))$ for all $f \in C(\mathcal{X})$, the space of bounded continuous functions on \mathcal{X} [Dudley, 2002]

Universal RKHS: $k(x, x')$ continuous, \mathcal{X} compact, and \mathcal{F} dense in $C(\mathcal{X})$ with respect to L_{∞} [Steinwart, 2001]

If \mathcal{F} **universal**, then $\text{MMD} \{\mathbf{P}, \mathbf{Q}; F\} = 0$ iff $\mathbf{P} = \mathbf{Q}$

Characteristic Kernels (via universality)

Proof:

First, it is clear that $\mathbf{P} = \mathbf{Q}$ implies $\text{MMD} \{ \mathbf{P}, \mathbf{Q}; F \}$ is zero.

Converse: by the universality of \mathcal{F} , for any given $\epsilon > 0$ and $f \in C(\mathcal{X}) \exists g \in \mathcal{F}$

$$\|f - g\|_{\infty} \leq \epsilon.$$

Characteristic Kernels (via universality)

Proof:

First, it is clear that $\mathbf{P} = \mathbf{Q}$ implies $\text{MMD}\{\mathbf{P}, \mathbf{Q}; F\}$ is zero.

Converse: by the universality of \mathcal{F} , for any given $\epsilon > 0$ and $f \in C(\mathcal{X}) \exists g \in \mathcal{F}$

$$\|f - g\|_{\infty} \leq \epsilon.$$

We next make the expansion

$$|\mathbf{E}_{\mathbf{P}} f(x) - \mathbf{E}_{\mathbf{Q}} f(y)| \leq |\mathbf{E}_{\mathbf{P}} f(x) - \mathbf{E}_{\mathbf{P}} g(x)| + |\mathbf{E}_{\mathbf{P}} g(x) - \mathbf{E}_{\mathbf{Q}} g(y)| + |\mathbf{E}_{\mathbf{Q}} g(y) - \mathbf{E}_{\mathbf{Q}} f(y)|.$$

The first and third terms satisfy

$$|\mathbf{E}_{\mathbf{P}} f(x) - \mathbf{E}_{\mathbf{P}} g(x)| \leq \mathbf{E}_{\mathbf{P}} |f(x) - g(x)| \leq \epsilon.$$

Characteristic Kernels (via universality)

Proof (continued):

Next, write

$$\mathbf{E}_{\mathbf{P}}g(x) - \mathbf{E}_{\mathbf{Q}}g(y) = \langle g(\cdot), \mu_{\mathbf{P}} - \mu_{\mathbf{Q}} \rangle_{\mathcal{F}} = 0,$$

since $\text{MMD}\{\mathbf{P}, \mathbf{Q}; F\} = 0$ implies $\mu_{\mathbf{P}} = \mu_{\mathbf{Q}}$. Hence

$$|\mathbf{E}_{\mathbf{P}}f(x) - \mathbf{E}_{\mathbf{Q}}f(y)| \leq 2\epsilon$$

for all $f \in C(\mathcal{X})$ and $\epsilon > 0$, which implies $\mathbf{P} = \mathbf{Q}$.

References

- V. Alba Fernández, M. Jiménez-Gamero, and J. Muñoz García. A test for the two-sample problem based on empirical characteristic functions. *Comput. Stat. Data An.*, 52:3730–3748, 2008.
- N. Anderson, P. Hall, and D. Titterton. Two-sample test statistics for measuring discrepancies between two multivariate probability density functions using kernel-based density estimates. *Journal of Multivariate Analysis*, 50:41–54, 1994.
- M. Arcones and E. Giné. On the bootstrap of u and v statistics. *The Annals of Statistics*, 20(2):655–674, 1992.
- R. M. Dudley. *Real analysis and probability*. Cambridge University Press, Cambridge, UK, 2002.
- Andrey Feuerverger. A consistent test for bivariate dependence. *International Statistical Review*, 61(3):419–433, 1993.
- K. Fukumizu, B. Sriperumbudur, A. Gretton, and B. Schoelkopf. Characteristic kernels on groups and semigroups. In *Advances in Neural Information Processing Systems 21*, pages 473–480, Red Hook, NY, 2009. Curran Associates Inc.
- Z. Harchaoui, F. Bach, and E. Moulines. Testing for homogeneity with kernel Fisher discriminant analysis. In *Advances in Neural Information Processing Systems 20*, pages 609–616. MIT Press, Cambridge, MA, 2008.
- W. Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58:13–30, 1963.
- Wassily Hoeffding. A class of statistics with asymptotically normal distribution. *The Annals of Mathematical Statistics*, 19(3):293–325, 1948.
- N. L. Johnson, S. Kotz, and N. Balakrishnan. *Continuous Univariate Distributions. Volume 1*. John Wiley and Sons, 2nd edition, 1994.
- A. Kankainen. *Consistent Testing of Total Independence Based on the Empirical Characteristic Function*. PhD thesis, University of Jyväskylä, 1995.
- S. Mallat. *A Wavelet Tour of Signal Processing*. Academic Press, 2nd edition, 1999.
- C. McDiarmid. On the method of bounded differences. In *Survey in Combinatorics*, pages 148–188. Cambridge University Press, 1989.

- A. Müller. Integral probability metrics and their generating classes of functions. *Advances in Applied Probability*, 29(2):429–443, 1997.
- R. Serfling. *Approximation Theorems of Mathematical Statistics*. Wiley, New York, 1980.
- B. Sriperumbudur, A. Gretton, K. Fukumizu, G. Lanckriet, and B. Schölkopf. Hilbert space embeddings and metrics on probability measures. *Journal of Machine Learning Research*, 11:1517–1561, 2010.
- I. Steinwart. On the influence of the kernel on the consistency of support vector machines. *Journal of Machine Learning Research*, 2:67–93, 2001.
- Ingo Steinwart and Andreas Christmann. *Support Vector Machines*. Information Science and Statistics. Springer, 2008.
- G. Székely and M. Rizzo. Brownian distance covariance. *Annals of Applied Statistics*, 4(3):1233–1303, 2009.
- G. Székely, M. Rizzo, and N. Bakirov. Measuring and testing dependence by correlation of distances. *Ann. Stat.*, 35(6):2769–2794, 2007.
- S. K. Zhou and R. Chellappa. From sample similarity to ensemble similarity: Probabilistic distance measures in reproducing kernel hilbert space. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28(6):917–929, 2006.