

# Representing and comparing probabilities with kernels: Part 2

Arthur Gretton

Gatsby Computational Neuroscience Unit,  
University College London

MLSS Tuebingen, 2020

## Comparing two samples

- Given: Samples from unknown distributions  $P$  and  $Q$ .
- Goal: do  $P$  and  $Q$  differ?



$\sim P$



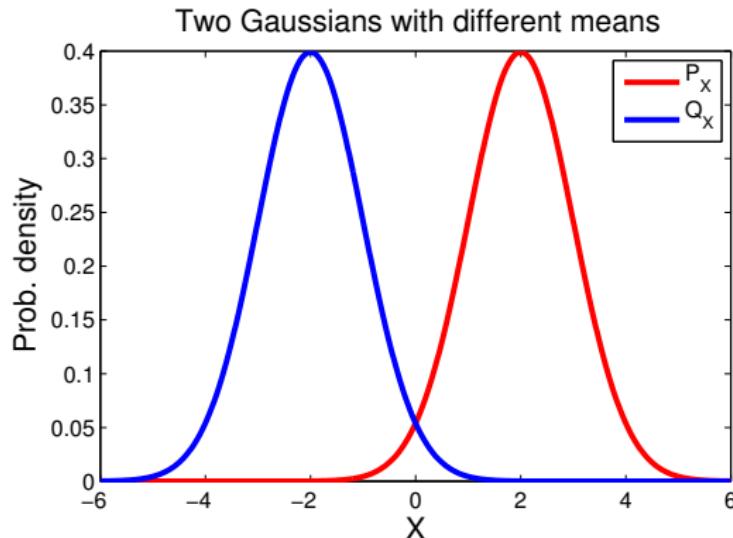
$\sim Q$

# Outline

- Maximum Mean Discrepancy (MMD)...
  - ...as a difference in feature means
  - ...as an integral probability metric (not just a technicality!)
- A statistical test based on the MMD
- Next slides: training generative adversarial networks with MMD
  - Gradient regularisation and data adaptivity

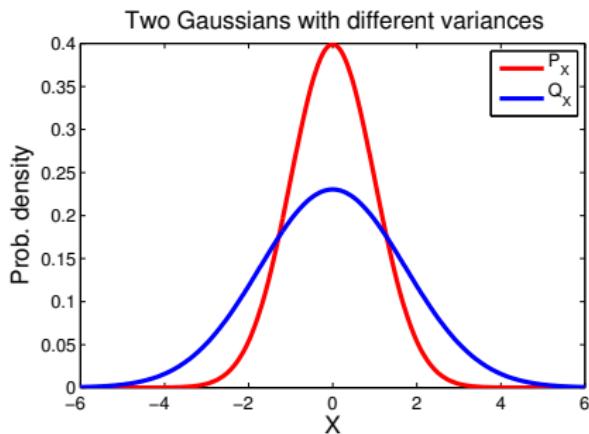
## Feature mean difference

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



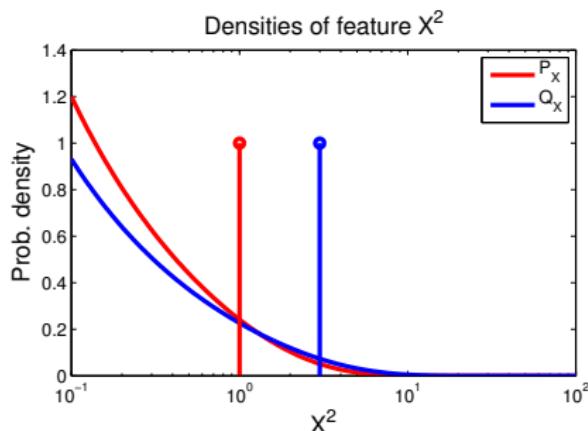
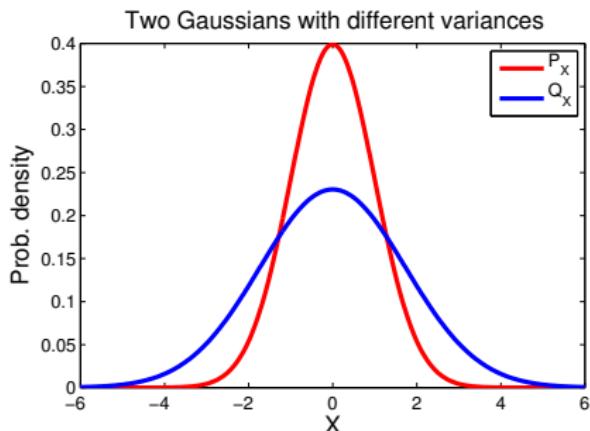
## Feature mean difference

- 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  $\varphi(x) = x^2$



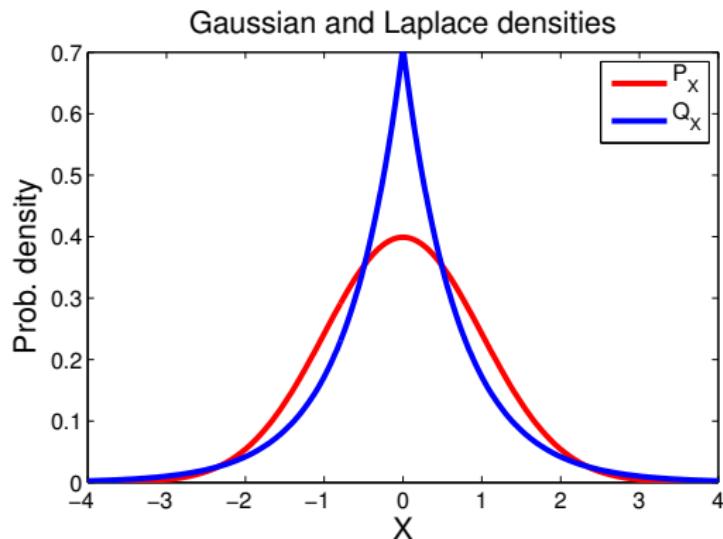
## Feature mean difference

- 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



## Infinitely many features using kernels

Kernels: dot products  
of features

Exponentiated quadratic kernel

$$k(x, x') = \exp(-\gamma \|x - x'\|^2)$$

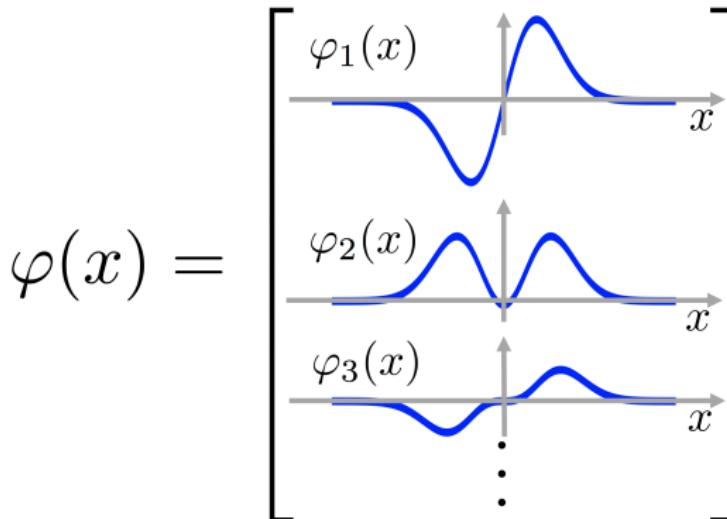
Feature map  $\varphi(x) \in \mathcal{F}$ ,

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

For positive definite  $k$ ,

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

Infinitely many features  
 $\varphi(x)$ , dot product in  
closed form!



## Infinitely many features of *distributions*

Given  $P$  a Borel **probability measure** on  $\mathcal{X}$ , define feature map of probability  $P$ ,

$$\mu_P = [\dots \mathbf{E}_P [\varphi_i(X)] \dots]$$

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

$$\langle \mu_P, \mu_Q \rangle_{\mathcal{F}} = \mathbf{E}_{P,Q} k(\textcolor{blue}{x}, \textcolor{red}{y})$$

for  $x \sim P$  and  $y \sim Q$ .

Fine print: feature map  $\varphi(x)$  must be Bochner integrable for all probability measures considered.  
Always true if kernel bounded.

## Infinitely many features of *distributions*

Given  $P$  a Borel **probability measure** on  $\mathcal{X}$ , define feature map of probability  $P$ ,

$$\mu_P = [\dots \mathbf{E}_P [\varphi_i(X)] \dots]$$

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

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

for  $x \sim P$  and  $y \sim Q$ .

**Fine print:** feature map  $\varphi(x)$  must be Bochner integrable for all probability measures considered.  
Always true if kernel bounded.

## The maximum mean discrepancy

The maximum mean discrepancy is the distance between **feature means**:

$$\begin{aligned} MMD^2(P, Q) &= \|\mu_P - \mu_Q\|_{\mathcal{F}}^2 \\ &= \langle \mu_P, \mu_P \rangle_{\mathcal{F}} + \langle \mu_Q, \mu_Q \rangle_{\mathcal{F}} - 2 \langle \mu_P, \mu_Q \rangle_{\mathcal{F}} \\ &= \underbrace{\mathbf{E}_P k(X, X')}_{(a)} + \underbrace{\mathbf{E}_Q k(Y, Y')}_{(a)} - 2 \underbrace{\mathbf{E}_{P,Q} k(X, Y)}_{(b)} \end{aligned}$$

## The maximum mean discrepancy

The maximum mean discrepancy is the distance between **feature means**:

$$\begin{aligned} MMD^2(P, Q) &= \|\mu_P - \mu_Q\|_{\mathcal{F}}^2 \\ &= \langle \mu_P, \mu_P \rangle_{\mathcal{F}} + \langle \mu_Q, \mu_Q \rangle_{\mathcal{F}} - 2 \langle \mu_P, \mu_Q \rangle_{\mathcal{F}} \\ &= \underbrace{\mathbf{E}_P k(X, X')}_{(a)} + \underbrace{\mathbf{E}_Q k(Y, Y')}_{(a)} - 2 \underbrace{\mathbf{E}_{P, Q} k(X, Y)}_{(b)} \end{aligned}$$

## The maximum mean discrepancy

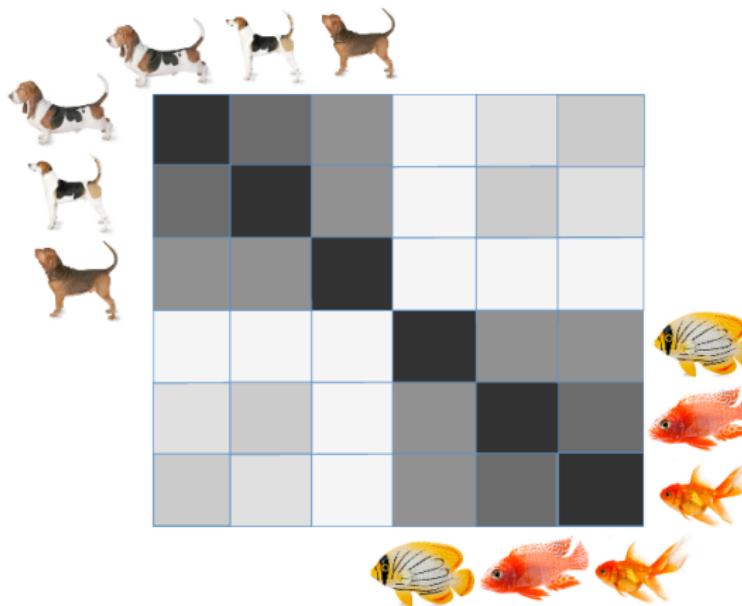
The maximum mean discrepancy is the distance between feature means:

$$\begin{aligned} MMD^2(P, Q) &= \|\mu_P - \mu_Q\|_{\mathcal{F}}^2 \\ &= \langle \mu_P, \mu_P \rangle_{\mathcal{F}} + \langle \mu_Q, \mu_Q \rangle_{\mathcal{F}} - 2 \langle \mu_P, \mu_Q \rangle_{\mathcal{F}} \\ &= \underbrace{\mathbf{E}_P k(X, X')}_{(a)} + \underbrace{\mathbf{E}_Q k(Y, Y')}_{(a)} - 2 \underbrace{\mathbf{E}_{P,Q} k(X, Y)}_{(b)} \end{aligned}$$

(a)= within distrib. similarity, (b)= cross-distrib. similarity.

## Illustration of MMD

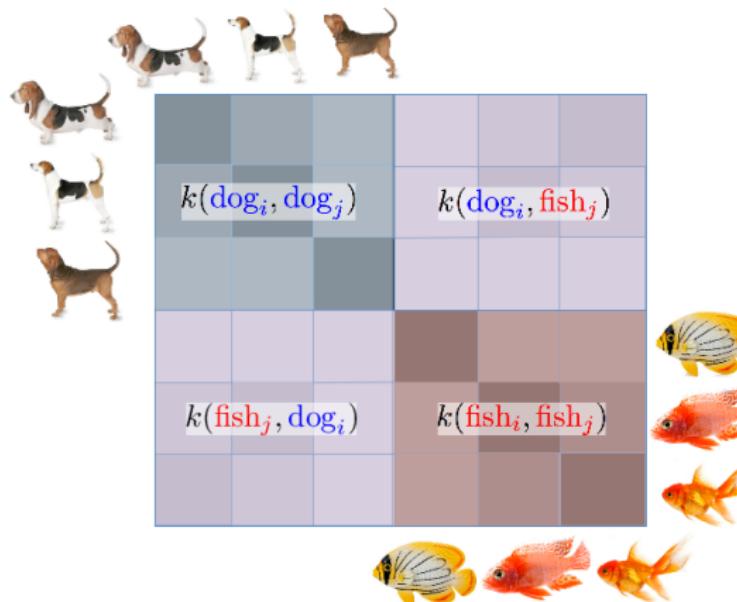
- Dogs ( $= P$ ) and fish ( $= Q$ ) example revisited
- Each entry is one of  $k(\text{dog}_i, \text{dog}_j)$ ,  $k(\text{dog}_i, \text{fish}_j)$ , or  $k(\text{fish}_i, \text{fish}_j)$



## Illustration of MMD

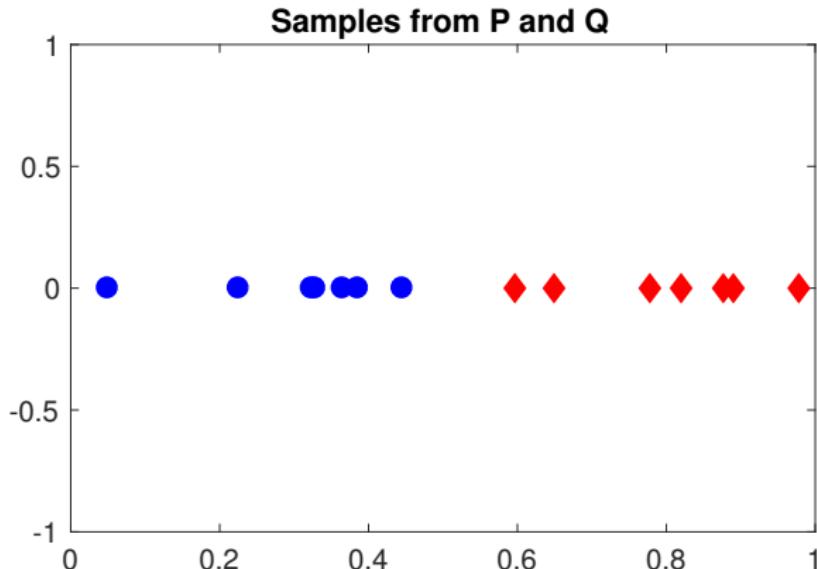
The maximum mean discrepancy:

$$\widehat{MMD}^2 = \frac{1}{n(n-1)} \sum_{i \neq j} k(\text{dog}_i, \text{dog}_j) + \frac{1}{n(n-1)} \sum_{i \neq j} k(\text{fish}_i, \text{fish}_j) - \frac{2}{n^2} \sum_{i,j} k(\text{dog}_i, \text{fish}_j)$$



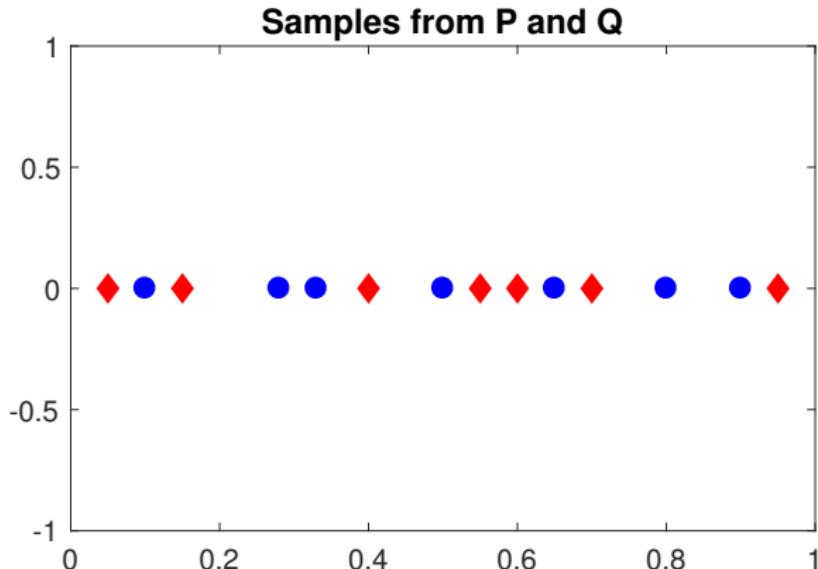
## MMD as an integral probability metric

Are  $P$  and  $Q$  different?



## MMD as an integral probability metric

Are  $P$  and  $Q$  different?

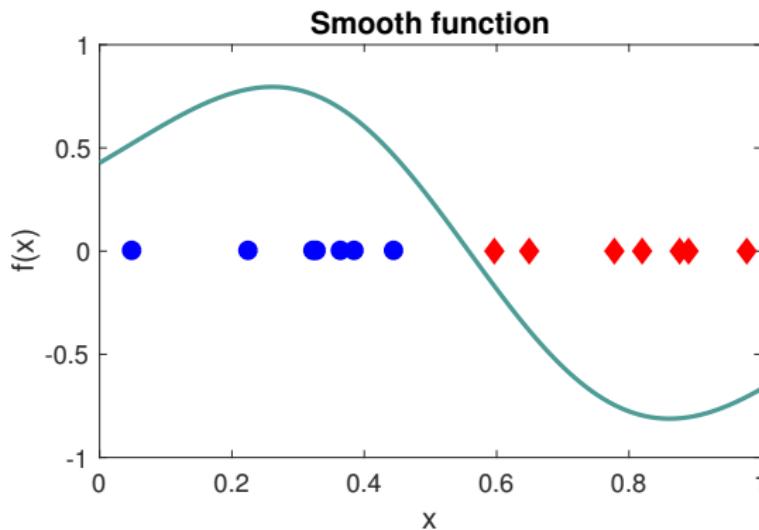


## MMD as an integral probability metric

Integral probability metric:

Find a "well behaved function"  $f(x)$  to maximize

$$\mathbf{E}_P f(\mathcal{X}) - \mathbf{E}_Q f(\mathcal{Y})$$

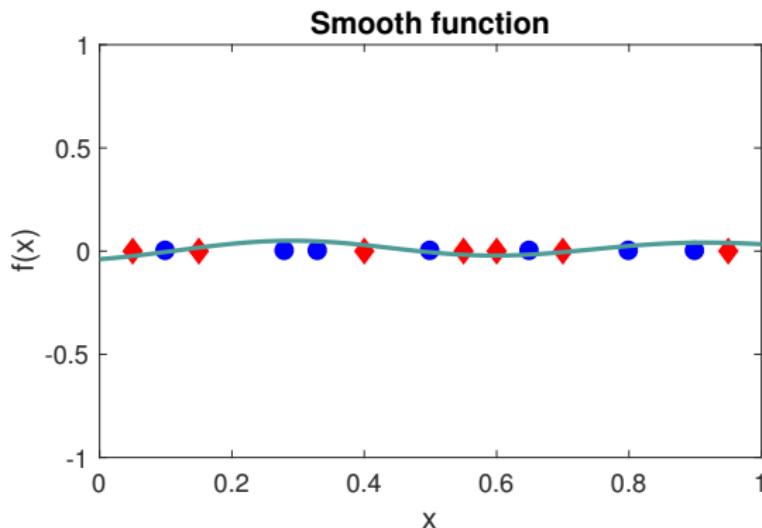


## MMD as an integral probability metric

Integral probability metric:

Find a "well behaved function"  $f(x)$  to maximize

$$\mathbf{E}_P f(\mathcal{X}) - \mathbf{E}_Q f(\mathcal{Y})$$



## MMD as an integral probability metric

Maximum mean discrepancy: smooth function for  $P$  vs  $Q$

$$MMD(P, Q; \mathcal{F}) := \sup_{\|f\|_{\mathcal{F}} \leq 1} [\mathbf{E}_P f(X) - \mathbf{E}_Q f(Y)]$$

$(\mathcal{F} = \text{unit ball in RKHS } \mathcal{F})$

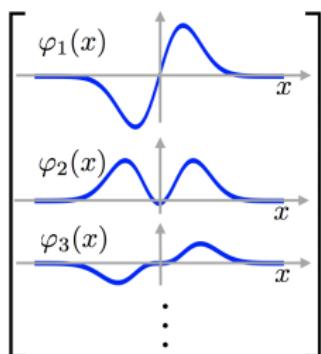
## MMD as an integral probability metric

Maximum mean discrepancy: smooth function for  $P$  vs  $Q$

$$MMD(P, Q; \mathcal{F}) := \sup_{\|\mathbf{f}\|_{\mathcal{F}} \leq 1} [\mathbf{E}_{P\mathbf{f}}(X) - \mathbf{E}_{Q\mathbf{f}}(Y)]$$

$(\mathcal{F} = \text{unit ball in RKHS } \mathcal{F})$

Functions are linear combinations of features:

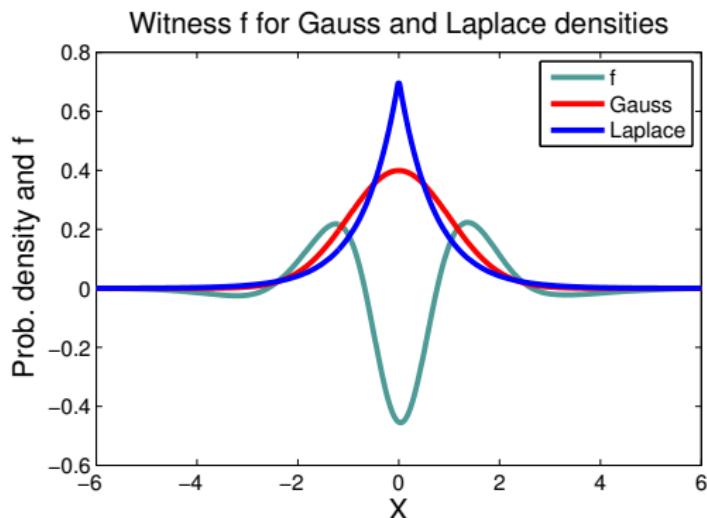
$$f(x) = \langle \mathbf{f}, \varphi(x) \rangle_{\mathcal{F}} = \sum_{\ell=1}^{\infty} f_{\ell} \varphi_{\ell}(x) = \begin{bmatrix} f_1 \\ f_2 \\ f_3 \\ \vdots \end{bmatrix}^{\top}$$

$$\|\mathbf{f}\|_{\mathcal{F}}^2 := \sum_{i=1}^{\infty} f_i^2 \leq 1$$

## MMD as an integral probability metric

Maximum mean discrepancy: smooth function for  $P$  vs  $Q$

$$MMD(P, Q; \mathcal{F}) := \sup_{\|f\|_{\mathcal{F}} \leq 1} [\mathbf{E}_P f(X) - \mathbf{E}_Q f(Y)]$$

$(\mathcal{F} = \text{unit ball in RKHS } \mathcal{F})$



## MMD as an integral probability metric

Maximum mean discrepancy: smooth function for  $P$  vs  $Q$

$$MMD(P, Q; \mathcal{F}) := \sup_{\|f\|_{\mathcal{F}} \leq 1} [\mathbf{E}_P f(X) - \mathbf{E}_Q f(Y)]$$

$(\mathcal{F} = \text{unit ball in RKHS } \mathcal{F})$

For characteristic RKHS  $\mathcal{F}$ ,  $MMD(P, Q; \mathcal{F}) = 0$  iff  $P = Q$

Other choices for witness function class:

- Bounded continuous [Dudley, 2002]
- Bounded variation 1 (Kolmogorov metric) [Müller, 1997]
- Bounded Lipschitz (Wasserstein distances) [Dudley, 2002]

## MMD as an integral probability metric

Maximum mean discrepancy: smooth function for  $P$  vs  $Q$

$$MMD(P, Q; F) := \sup_{\|f\|_F \leq 1} [\mathbf{E}_P f(X) - \mathbf{E}_Q f(Y)]$$

$(F = \text{unit ball in RKHS } \mathcal{F})$

Expectations of functions are linear combinations  
of expected features

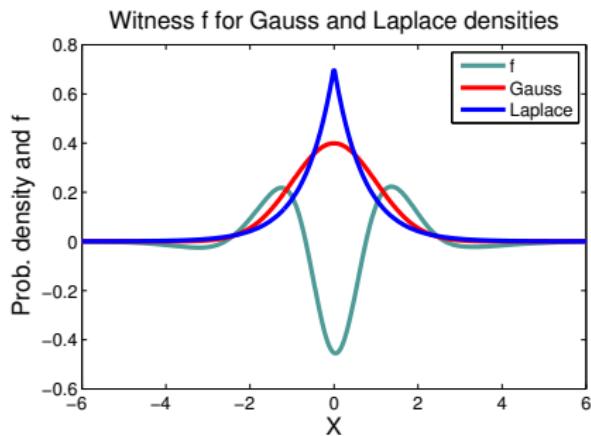
$$\mathbf{E}_P(f(X)) = \langle f, \mathbf{E}_P \varphi(X) \rangle_F = \langle f, \mu_P \rangle_F$$

(always true if kernel is bounded)

## Integral prob. metric vs feature difference

The MMD:

$$\begin{aligned} MMD(P, Q; \mathcal{F}) \\ = \sup_{\|\mathbf{f}\|_{\mathcal{F}} \leq 1} [\mathbf{E}_P f(\mathbf{X}) - \mathbf{E}_Q f(\mathbf{Y})] \end{aligned}$$



## Integral prob. metric vs feature difference

The MMD:

use

$$\begin{aligned} & MMD(P, Q; \mathcal{F}) \\ &= \sup_{\|\mathbf{f}\|_{\mathcal{F}} \leq 1} [\mathbf{E}_P \mathbf{f}(\mathbf{X}) - \mathbf{E}_Q \mathbf{f}(\mathbf{Y})] \\ &= \sup_{\|\mathbf{f}\|_{\mathcal{F}} \leq 1} \langle \mathbf{f}, \boldsymbol{\mu}_P - \boldsymbol{\mu}_Q \rangle_{\mathcal{F}} \end{aligned}$$
$$\mathbf{E}_P \mathbf{f}(\mathbf{X}) = \langle \boldsymbol{\mu}_P, \mathbf{f} \rangle_{\mathcal{F}}$$

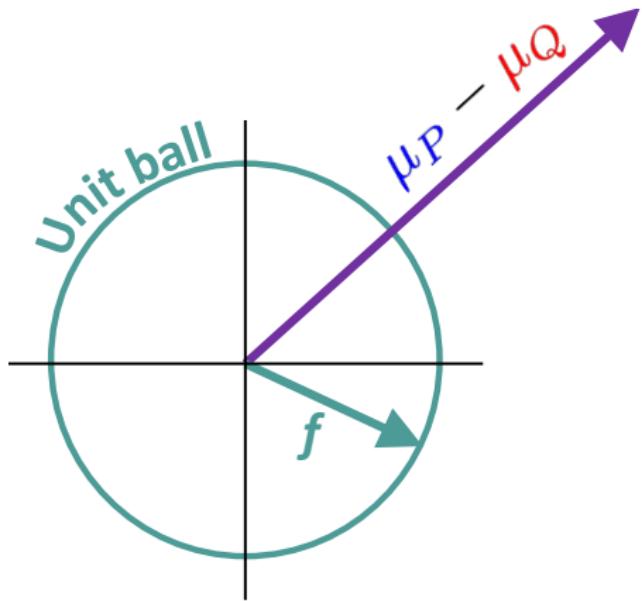
## Integral prob. metric vs feature difference

The MMD:

$$MMD(P, Q; \mathcal{F})$$

$$= \sup_{\|f\|_{\mathcal{F}} \leq 1} [\mathbf{E}_P f(X) - \mathbf{E}_Q f(Y)]$$

$$= \sup_{\|f\|_{\mathcal{F}} \leq 1} \langle f, \mu_P - \mu_Q \rangle_{\mathcal{F}}$$



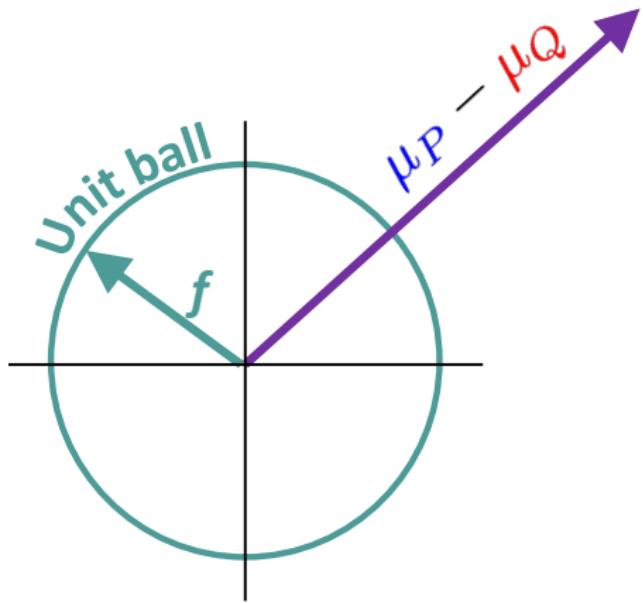
## Integral prob. metric vs feature difference

The MMD:

$$MMD(P, Q; \mathcal{F})$$

$$= \sup_{\|f\|_{\mathcal{F}} \leq 1} [\mathbf{E}_P f(X) - \mathbf{E}_Q f(Y)]$$

$$= \sup_{\|f\|_{\mathcal{F}} \leq 1} \langle f, \mu_P - \mu_Q \rangle_{\mathcal{F}}$$



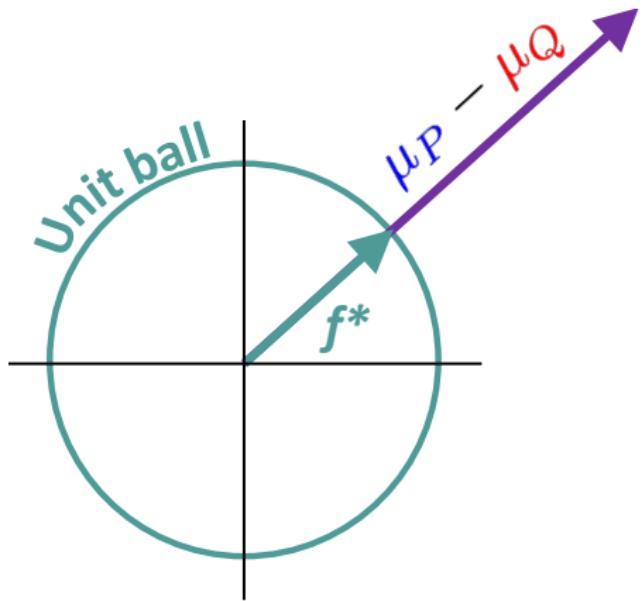
## Integral prob. metric vs feature difference

The MMD:

$$MMD(P, Q; \mathcal{F})$$

$$= \sup_{\|f\|_{\mathcal{F}} \leq 1} [\mathbf{E}_P f(X) - \mathbf{E}_Q f(Y)]$$

$$= \sup_{\|f\|_{\mathcal{F}} \leq 1} \langle f, \mu_P - \mu_Q \rangle_{\mathcal{F}}$$



$$f^* = \frac{\mu_P - \mu_Q}{\|\mu_P - \mu_Q\|}$$

## Integral prob. metric vs feature difference

The MMD:

$$\begin{aligned}MMD(P, Q; F) &= \sup_{\|f\|_F \leq 1} [\mathbf{E}_P f(X) - \mathbf{E}_Q f(Y)] \\&= \sup_{\|f\|_F \leq 1} \langle f, \mu_P - \mu_Q \rangle_F \\&= \|\mu_P - \mu_Q\|\end{aligned}$$

Function view and feature view equivalent

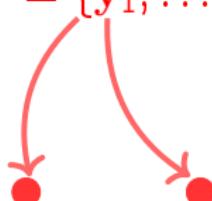
## Construction of MMD witness

Construction of empirical **witness function** (proof: next slide!)

Observe  $X = \{x_1, \dots, x_n\} \sim P$

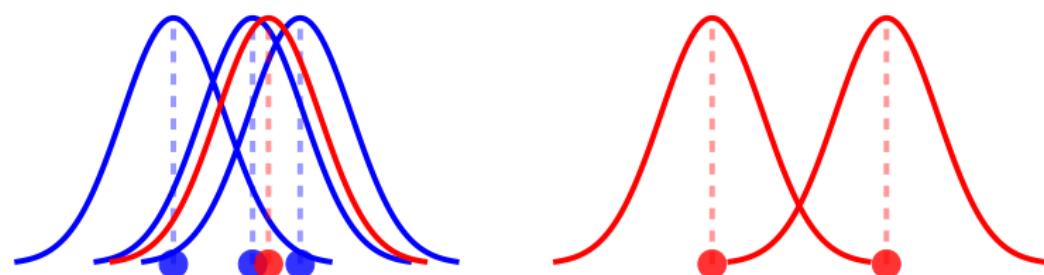


Observe  $Y = \{y_1, \dots, y_n\} \sim Q$



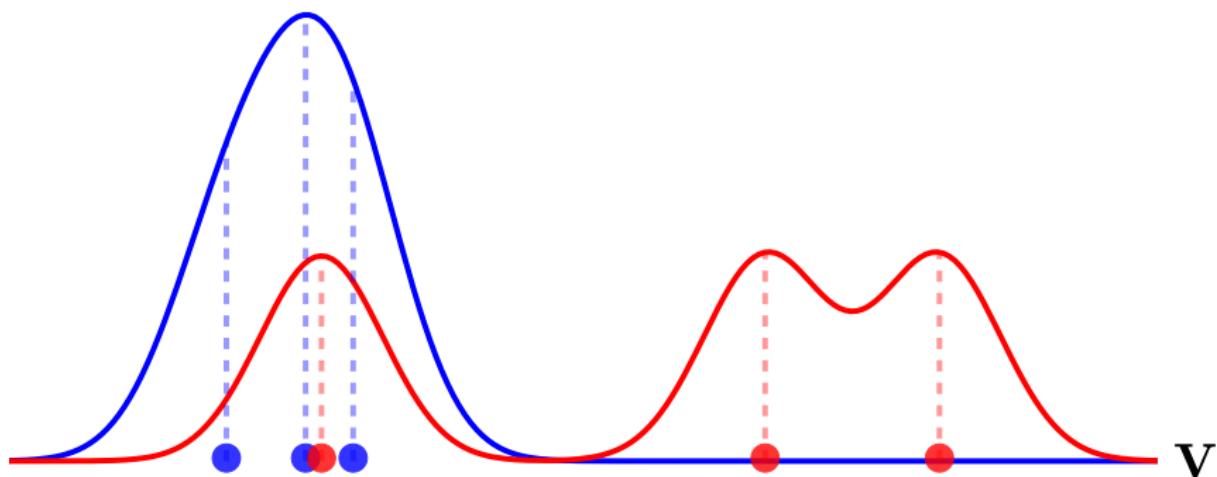
## Construction of MMD witness

Construction of empirical **witness function** (proof: next slide!)



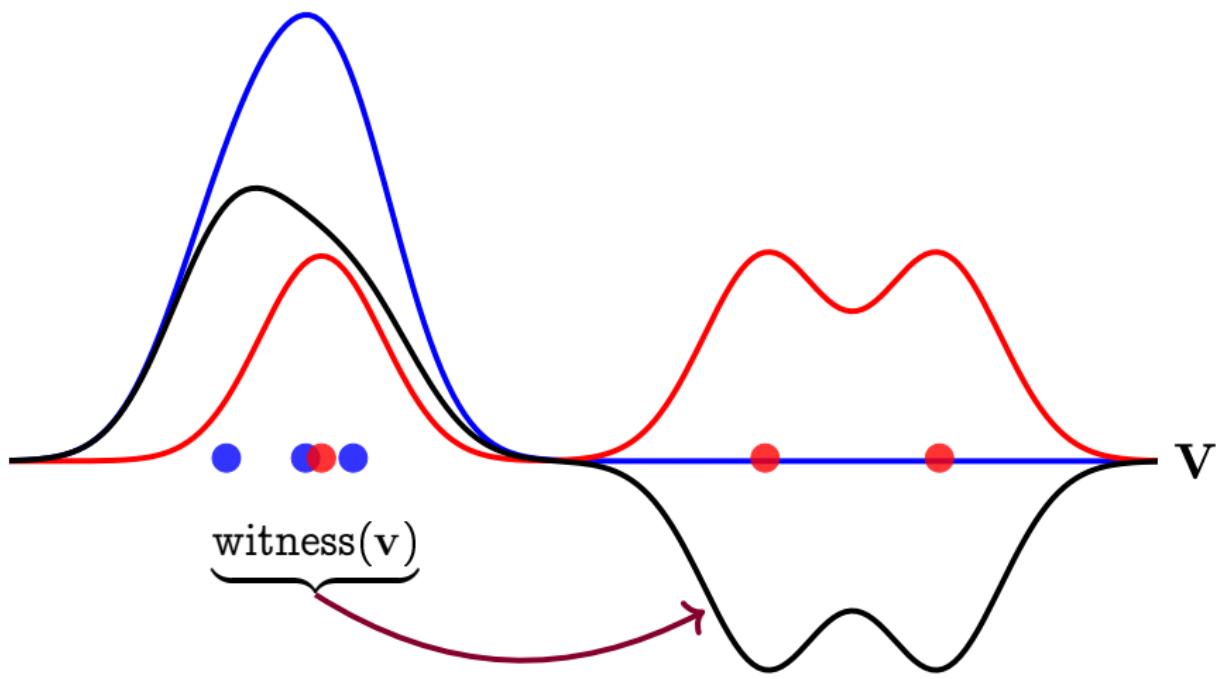
## Construction of MMD witness

Construction of empirical **witness function** (proof: next slide!)



## Construction of MMD witness

Construction of empirical **witness function** (proof: next slide!)



## Derivation of empirical witness function

Recall the witness function expression

$$f^* \propto \mu_P - \mu_Q$$

## Derivation of empirical witness function

Recall the **witness function** expression

$$f^* \propto \mu_P - \mu_Q$$

The empirical feature mean for  $P$

$$\hat{\mu}_P := \frac{1}{n} \sum_{i=1}^n \varphi(x_i)$$

## Derivation of empirical witness function

Recall the **witness function** expression

$$\textcolor{teal}{f}^* \propto \mu_P - \mu_Q$$

The empirical feature mean for  $P$

$$\widehat{\mu}_P := \frac{1}{n} \sum_{i=1}^n \varphi(x_i)$$

The empirical witness function at  $v$

$$\textcolor{teal}{f}^*(v) = \langle \textcolor{teal}{f}^*, \varphi(v) \rangle_{\mathcal{F}}$$

## Derivation of empirical witness function

Recall the **witness function** expression

$$\textcolor{teal}{f}^* \propto \mu_P - \mu_Q$$

The empirical feature mean for  $P$

$$\widehat{\mu}_P := \frac{1}{n} \sum_{i=1}^n \varphi(x_i)$$

The empirical witness function at  $v$

$$\begin{aligned}\textcolor{teal}{f}^*(v) &= \langle \textcolor{teal}{f}^*, \varphi(v) \rangle_{\mathcal{F}} \\ &\propto \langle \widehat{\mu}_P - \widehat{\mu}_Q, \varphi(v) \rangle_{\mathcal{F}}\end{aligned}$$

## Derivation of empirical witness function

Recall the witness function expression

$$\textcolor{teal}{f}^* \propto \mu_P - \mu_Q$$

The empirical feature mean for  $P$

$$\widehat{\mu}_P := \frac{1}{n} \sum_{i=1}^n \varphi(x_i)$$

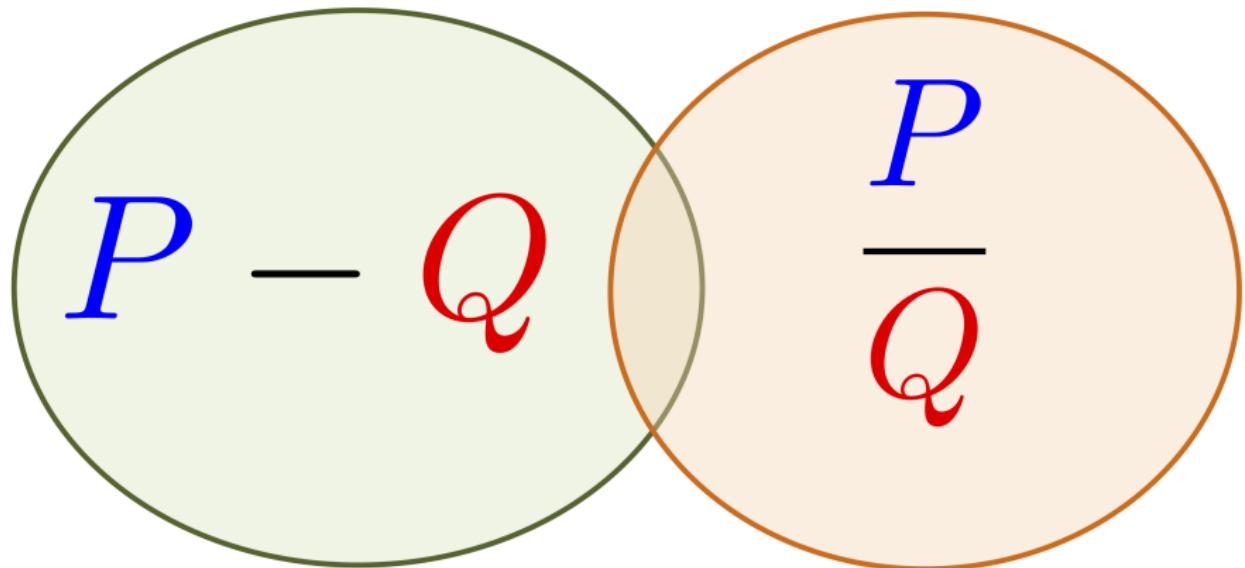
The empirical witness function at  $v$

$$\begin{aligned}\textcolor{teal}{f}^*(v) &= \langle \textcolor{teal}{f}^*, \varphi(v) \rangle_{\mathcal{F}} \\ &\propto \langle \widehat{\mu}_P - \widehat{\mu}_Q, \varphi(v) \rangle_{\mathcal{F}} \\ &= \frac{1}{n} \sum_{i=1}^n k(\textcolor{blue}{x}_i, v) - \frac{1}{n} \sum_{i=1}^n k(\textcolor{red}{y}_i, v)\end{aligned}$$

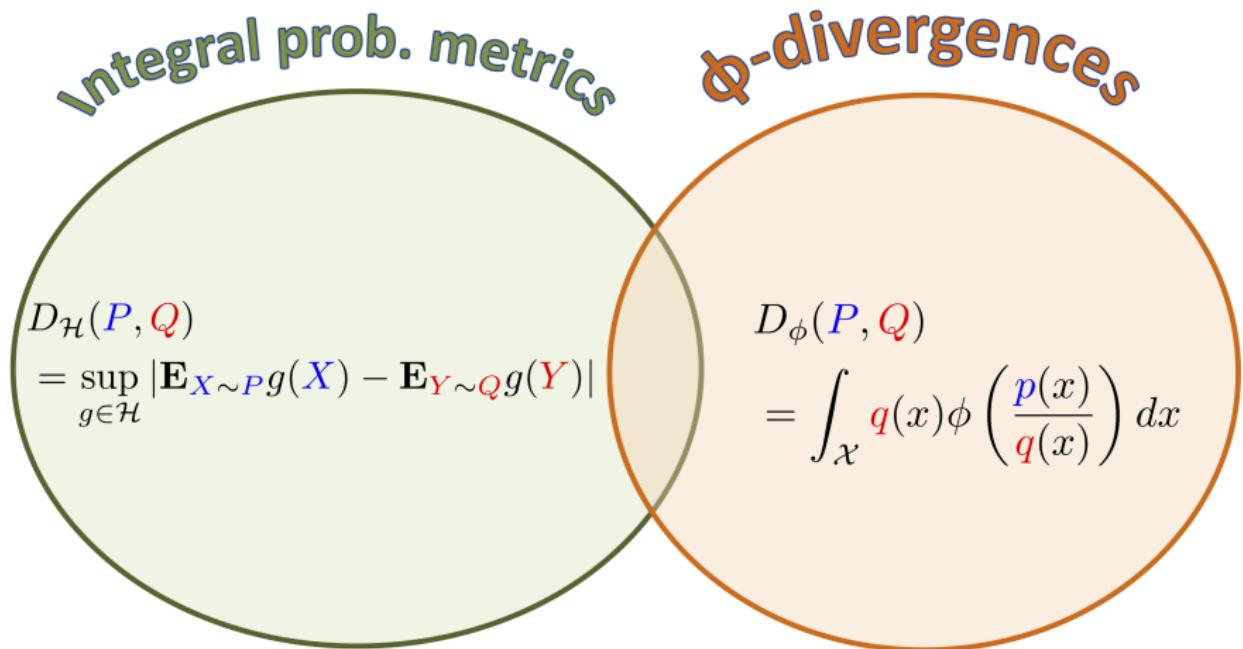
Don't need explicit feature coefficients  $f^* := [ f_1^* \ f_2^* \ \dots ]$

# Interlude: divergence measures

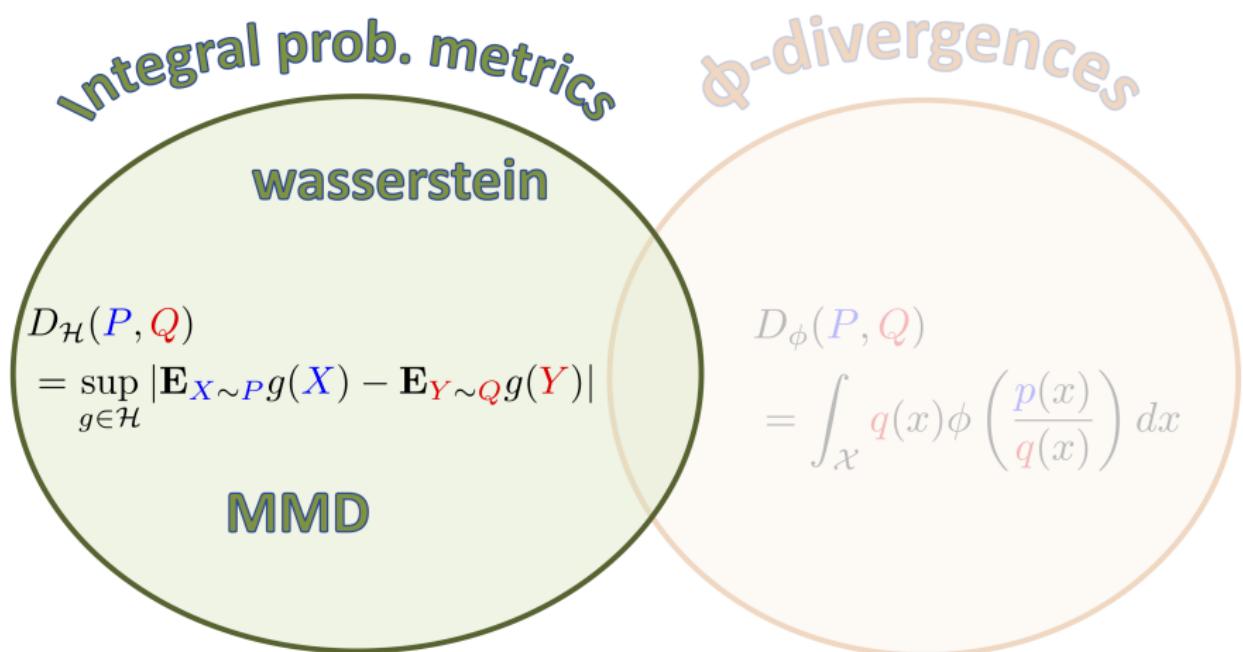
## Divergences



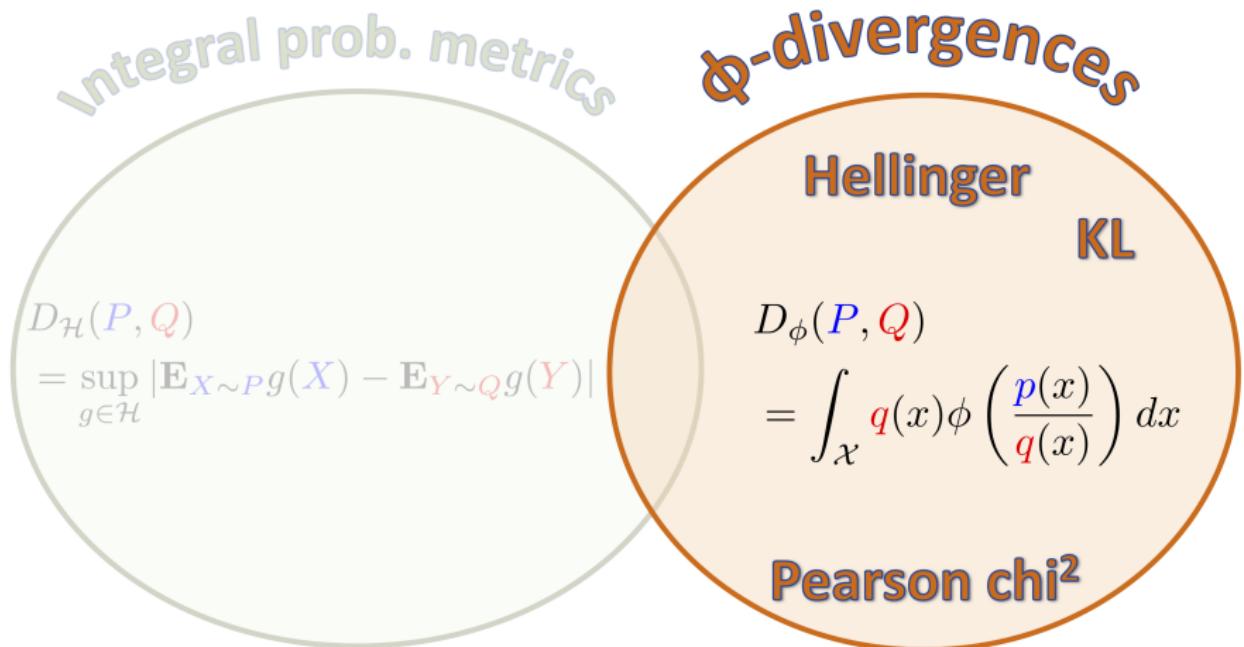
## Divergences



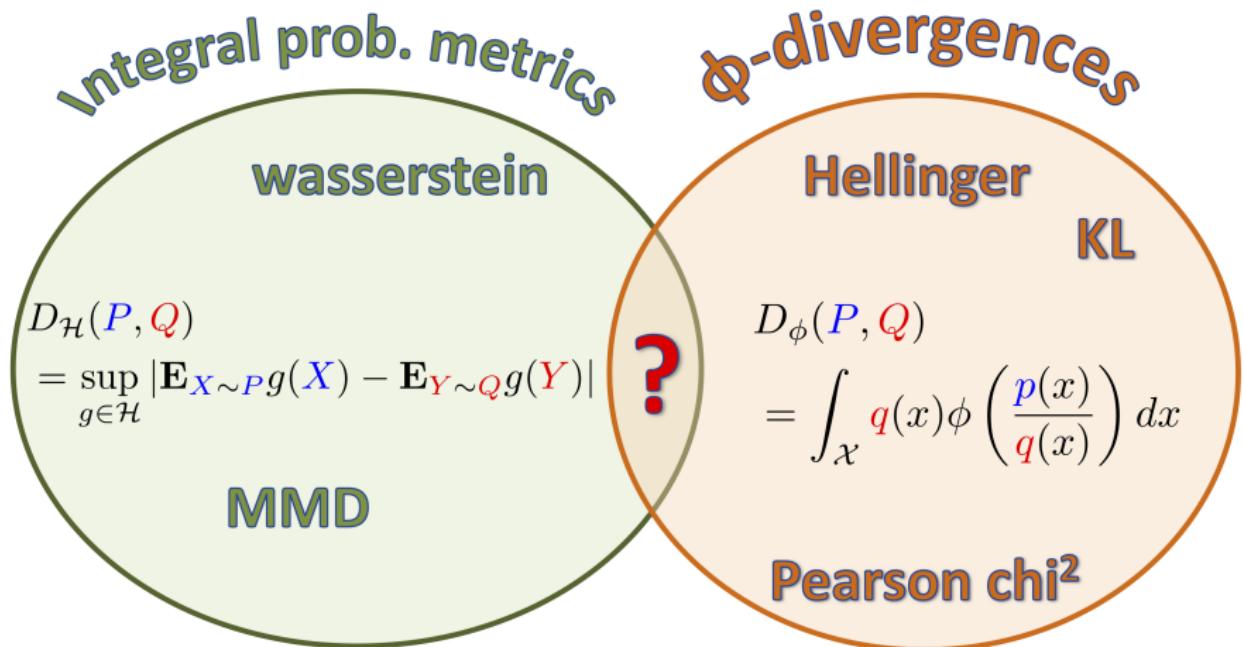
## The integral probability metrics



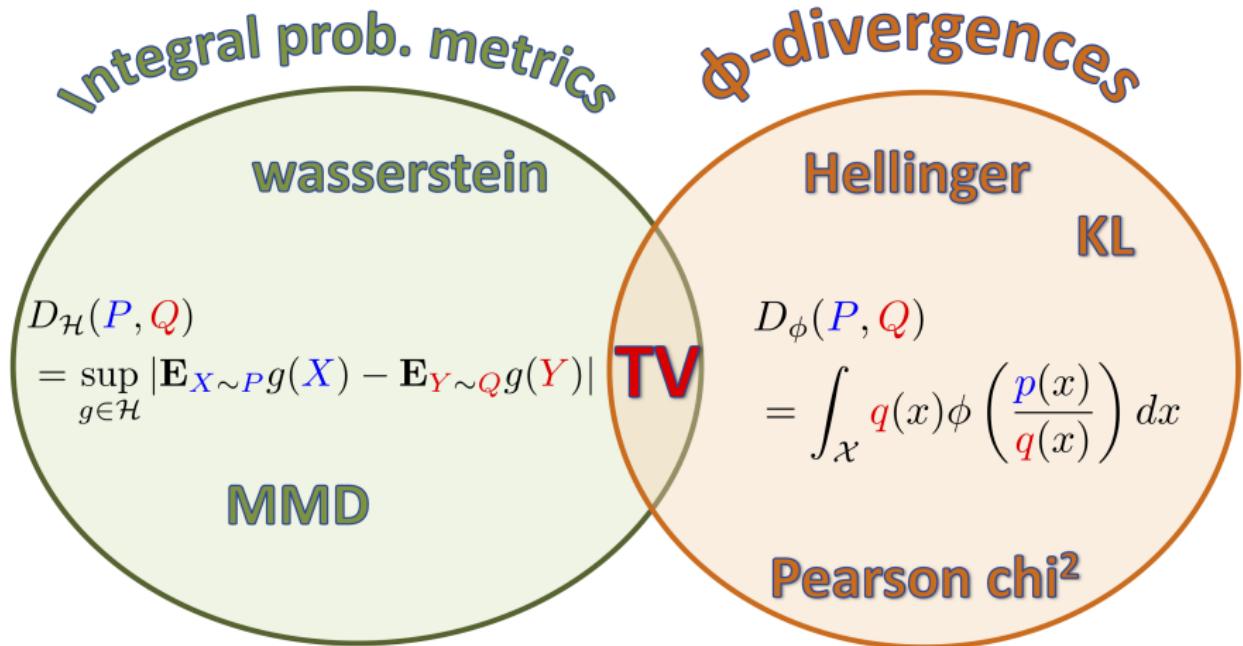
## The $\phi$ -divergences



## Divergences



## Divergences



Sriperumbudur, Fukumizu, G, Schoelkopf, Lanckriet, EJS (2012)

# Two-Sample Testing with MMD

## A statistical test using MMD

The empirical MMD:

$$\widehat{MMD}^2 = \frac{1}{n(n-1)} \sum_{i \neq j} k(\mathbf{x}_i, \mathbf{x}_j) + \frac{1}{n(n-1)} \sum_{i \neq j} k(\mathbf{y}_i, \mathbf{y}_j) - \frac{2}{n^2} \sum_{i,j} k(\mathbf{x}_i, \mathbf{y}_j)$$

How does this help decide whether  $P = Q$ ?

## A statistical test using MMD

The empirical MMD:

$$\widehat{MMD}^2 = \frac{1}{n(n-1)} \sum_{i \neq j} k(\mathbf{x}_i, \mathbf{x}_j) + \frac{1}{n(n-1)} \sum_{i \neq j} k(\mathbf{y}_i, \mathbf{y}_j) - \frac{2}{n^2} \sum_{i,j} k(\mathbf{x}_i, \mathbf{y}_j)$$

Perspective from [statistical hypothesis testing](#):

- Null hypothesis  $\mathcal{H}_0$  when  $P = Q$ 
  - should see  $\widehat{MMD}^2$  “close to zero”.
- Alternative hypothesis  $\mathcal{H}_1$  when  $P \neq Q$ 
  - should see  $\widehat{MMD}^2$  “far from zero”

## A statistical test using MMD

The empirical MMD:

$$\widehat{MMD}^2 = \frac{1}{n(n-1)} \sum_{i \neq j} k(\mathbf{x}_i, \mathbf{x}_j) + \frac{1}{n(n-1)} \sum_{i \neq j} k(\mathbf{y}_i, \mathbf{y}_j) - \frac{2}{n^2} \sum_{i,j} k(\mathbf{x}_i, \mathbf{y}_j)$$

Perspective from [statistical hypothesis testing](#):

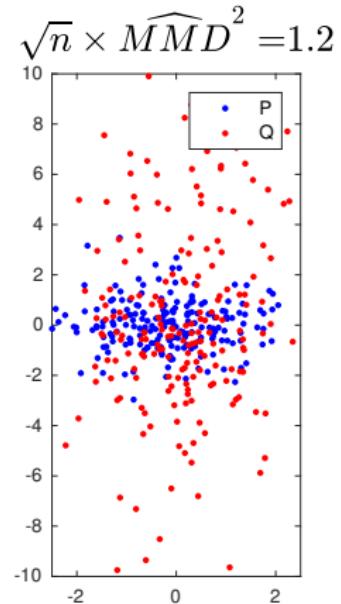
- Null hypothesis  $\mathcal{H}_0$  when  $P = Q$ 
  - should see  $\widehat{MMD}^2$  “close to zero”.
- Alternative hypothesis  $\mathcal{H}_1$  when  $P \neq Q$ 
  - should see  $\widehat{MMD}^2$  “far from zero”

Want [Threshold](#)  $c_\alpha$  for  $\widehat{MMD}^2$  to get [false positive rate](#)  $\alpha$

## Behaviour of $\widehat{MMD}^2$ when $P \neq Q$

Draw  $n = 200$  i.i.d samples from  $P$  and  $Q$

- Laplace with different y-variance.
- $\sqrt{n} \times \widehat{MMD}^2 = 1.2$

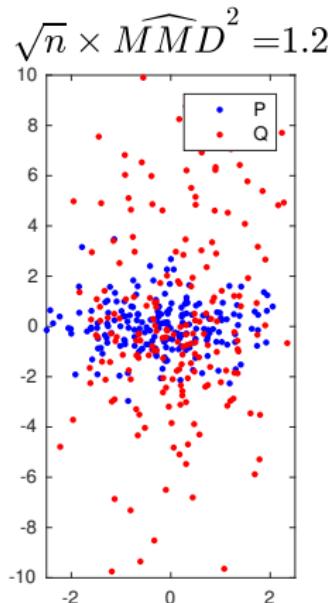
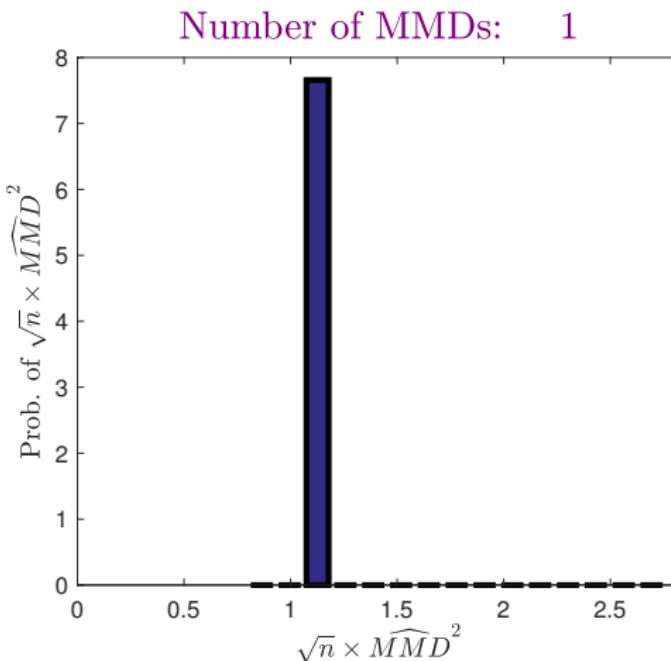


## Behaviour of $\widehat{MMD}^2$ when $P \neq Q$

Draw  $n = 200$  i.i.d samples from  $P$  and  $Q$

- Laplace with different y-variance.

- $\sqrt{n} \times \widehat{MMD}^2 = 1.2$

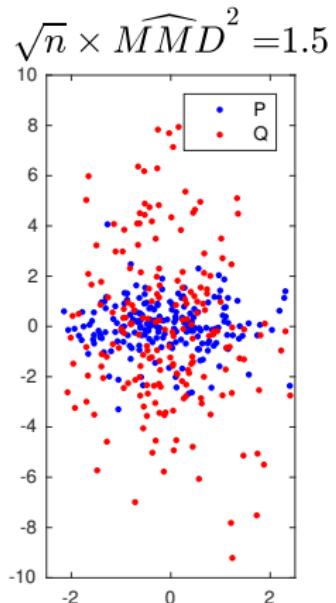
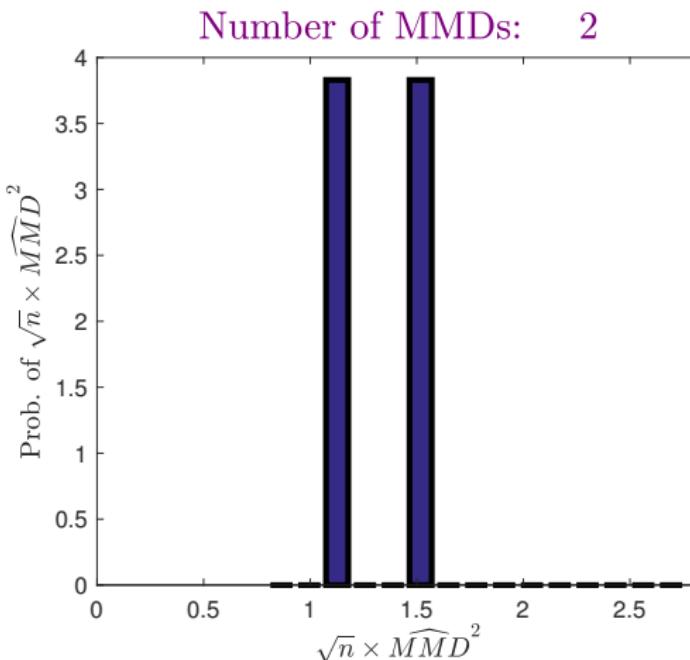


## Behaviour of $\widehat{MMD}^2$ when $P \neq Q$

Draw  $n = 200$  new samples from  $P$  and  $Q$

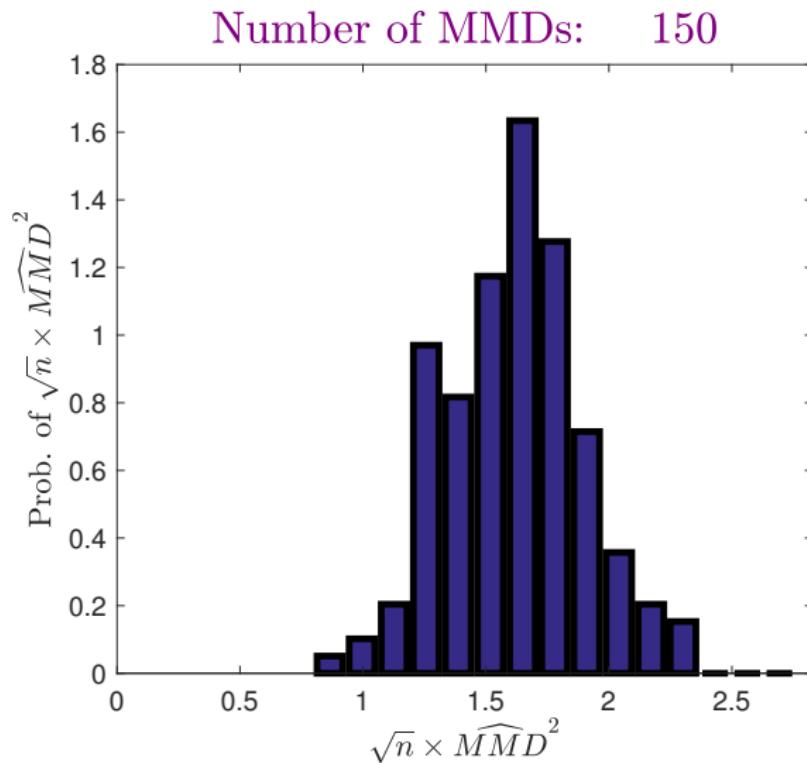
- Laplace with different y-variance.

- $\sqrt{n} \times \widehat{MMD}^2 = 1.5$



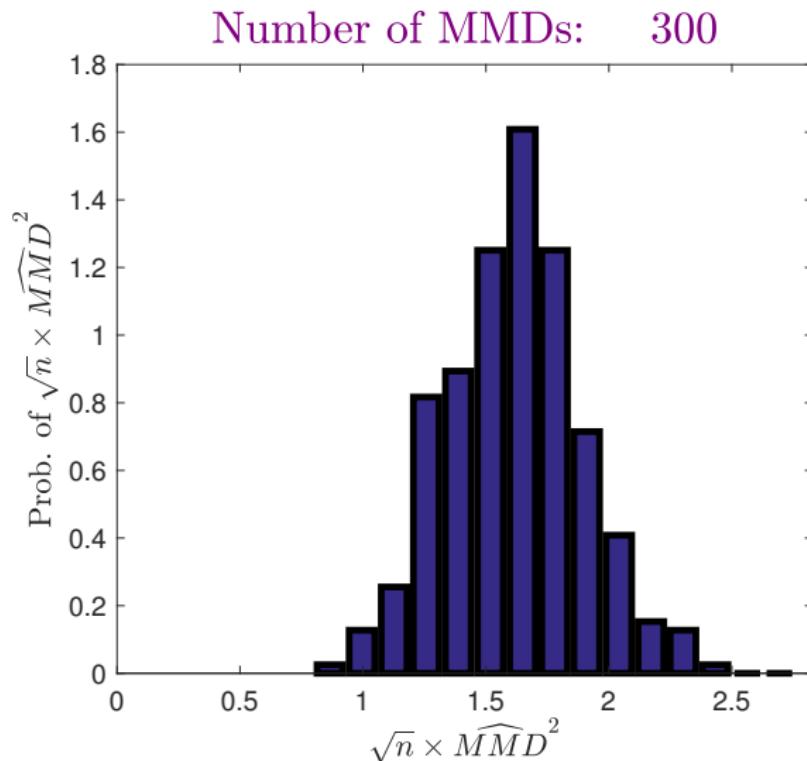
## Behaviour of $\widehat{MMD}^2$ when $P \neq Q$

Repeat this 150 times ...



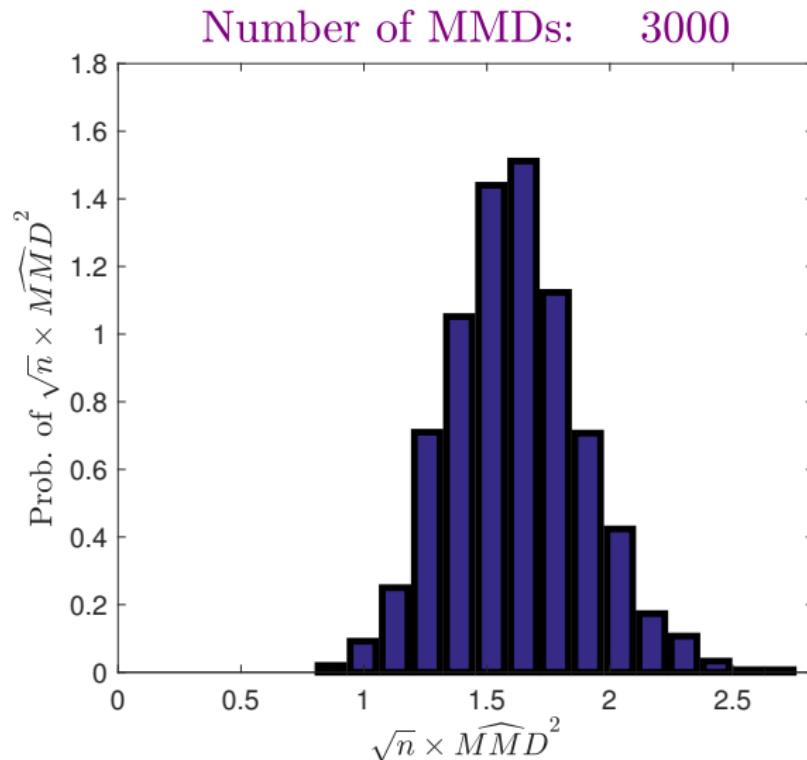
## Behaviour of $\widehat{MMD}^2$ when $P \neq Q$

Repeat this 300 times ...



## Behaviour of $\widehat{MMD}^2$ when $P \neq Q$

Repeat this 3000 times . . .



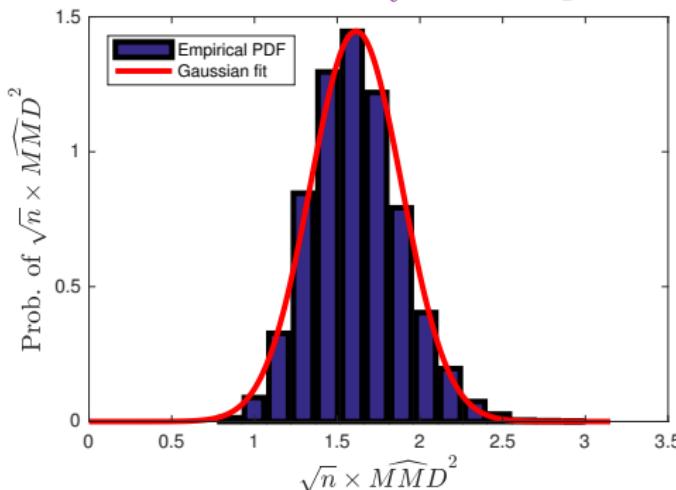
# Asymptotics of $\widehat{MMD}^2$ when $P \neq Q$

When  $P \neq Q$ , statistic is asymptotically normal,

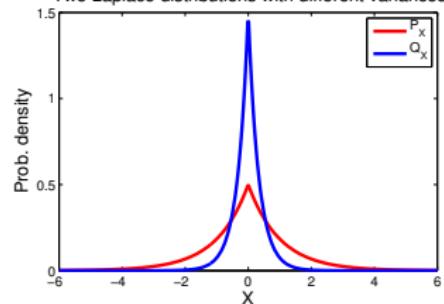
$$\frac{\widehat{MMD}^2 - MMD^2(P, Q)}{\sqrt{V_n(P, Q)}} \xrightarrow{D} \mathcal{N}(0, 1),$$

where variance  $V_n(P, Q) = O(n^{-1})$ .

MMD density under  $\mathcal{H}_1$



Two Laplace distributions with different variances

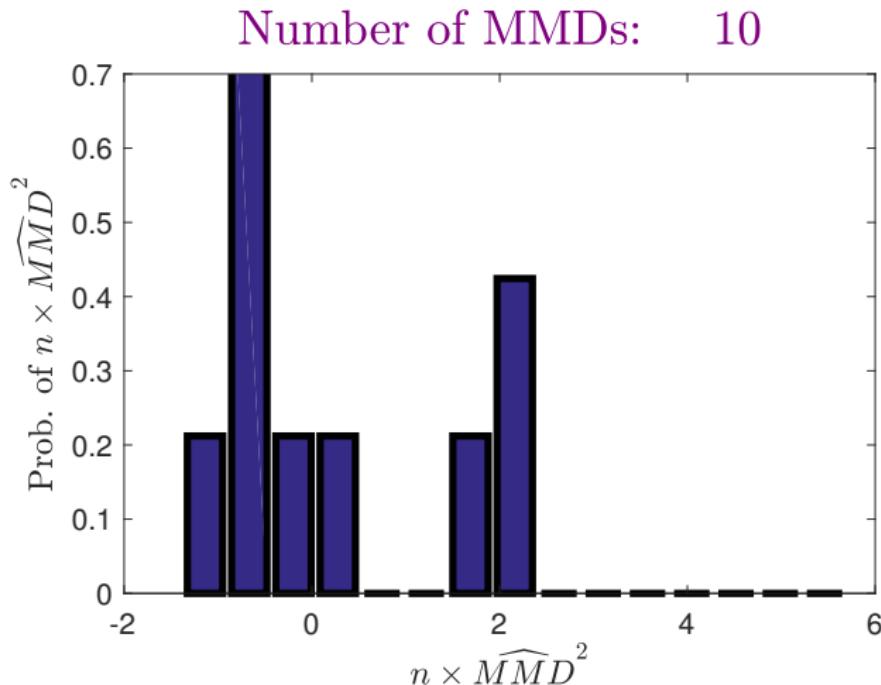


## Behaviour of $\widehat{MMD}^2$ when $P = Q$

What happens when  $P$  and  $Q$  are the same?

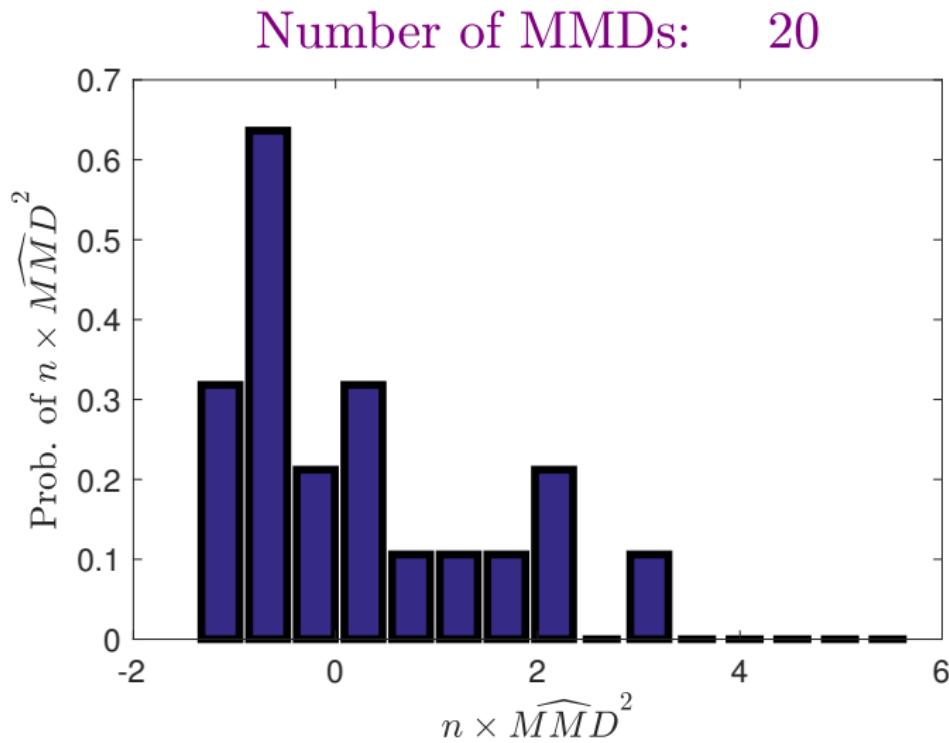
## Behaviour of $\widehat{MMD}^2$ when $P = Q$

- Case of  $P = Q = \mathcal{N}(0, 1)$



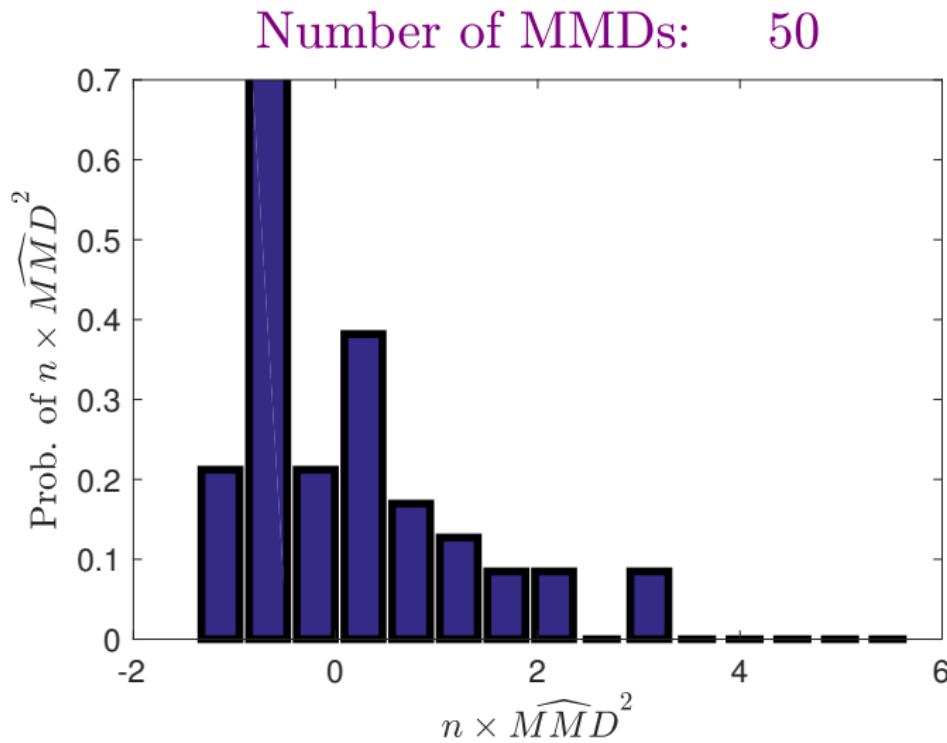
## Behaviour of $\widehat{MMD}^2$ when $P = Q$

- Case of  $P = Q = \mathcal{N}(0, 1)$



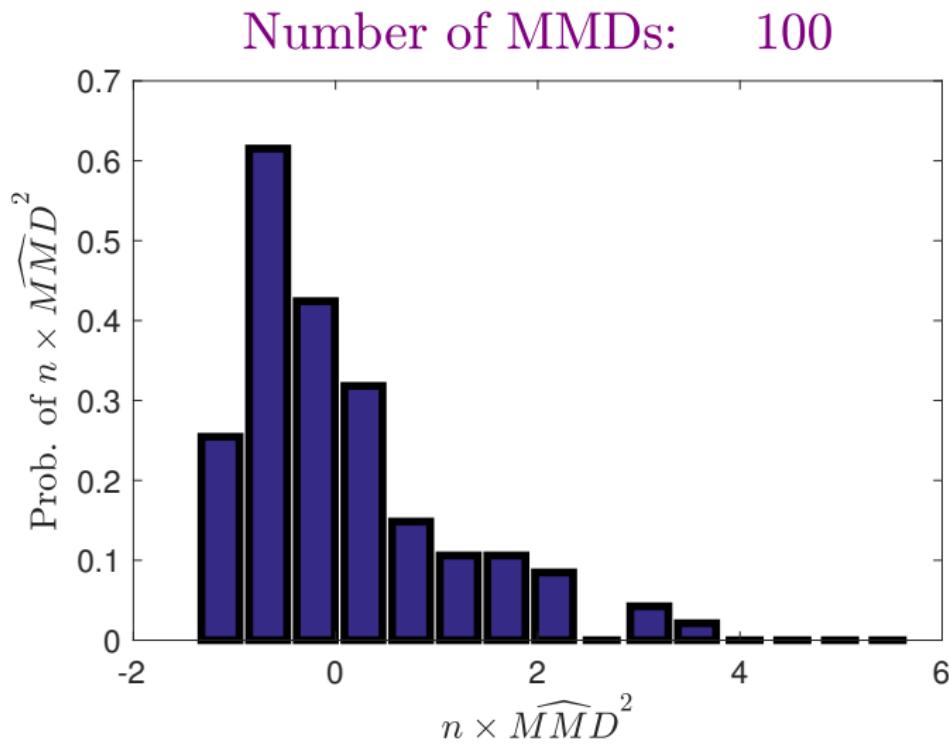
## Behaviour of $\widehat{MMD}^2$ when $P = Q$

- Case of  $P = Q = \mathcal{N}(0, 1)$



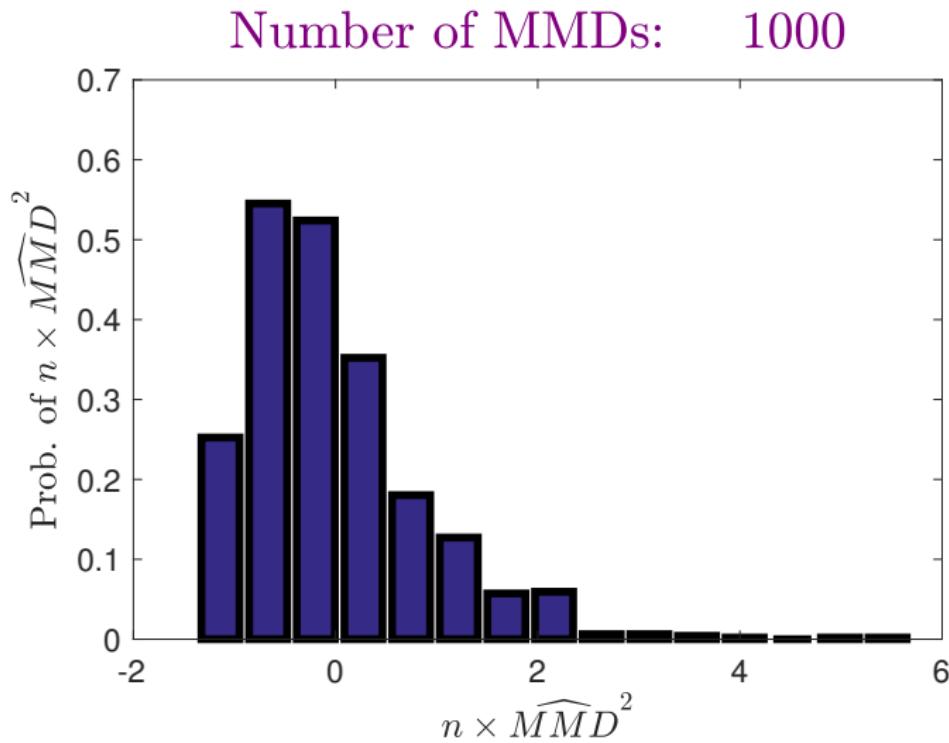
## Behaviour of $\widehat{MMD}^2$ when $P = Q$

- Case of  $P = Q = \mathcal{N}(0, 1)$



## Behaviour of $\widehat{MMD}^2$ when $P = Q$

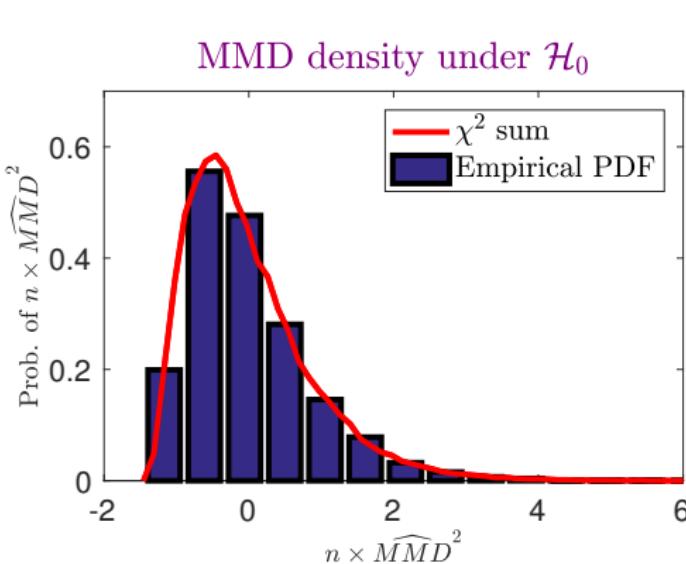
- Case of  $P = Q = \mathcal{N}(0, 1)$



## Asymptotics of $\widehat{MMD}^2$ when $P = Q$

Where  $P = Q$ , statistic has asymptotic distribution

$$n\widehat{MMD}^2 \sim \sum_{l=1}^{\infty} \lambda_l [z_l^2 - 2]$$



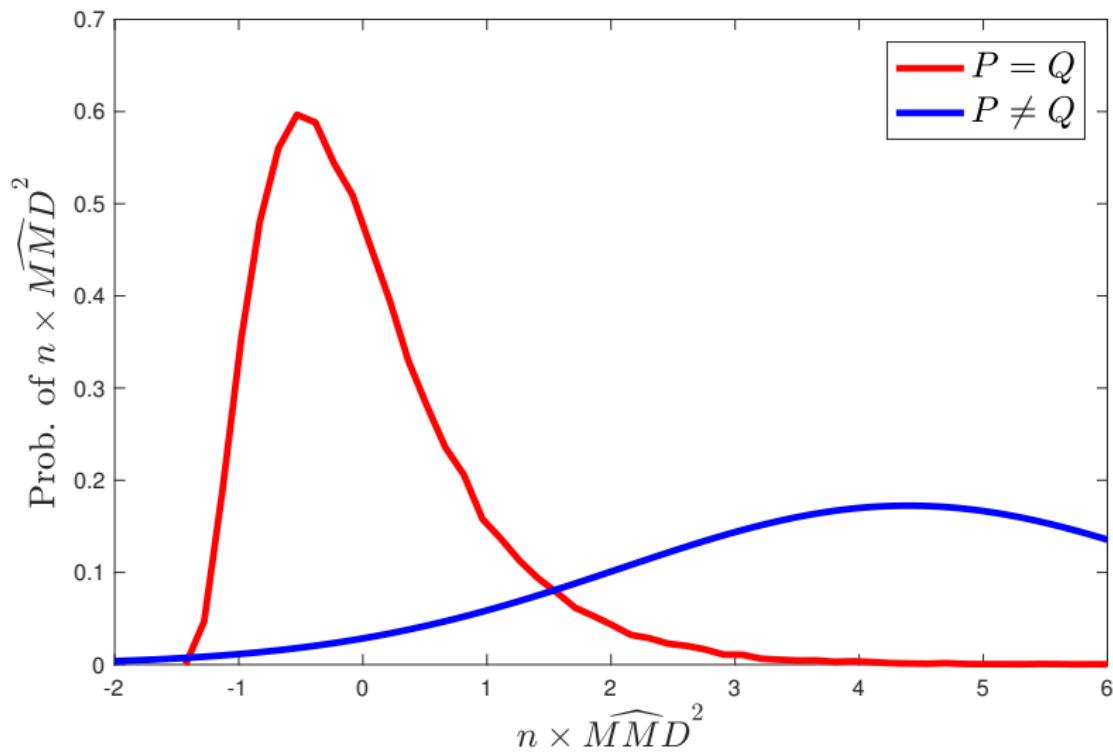
where

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

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

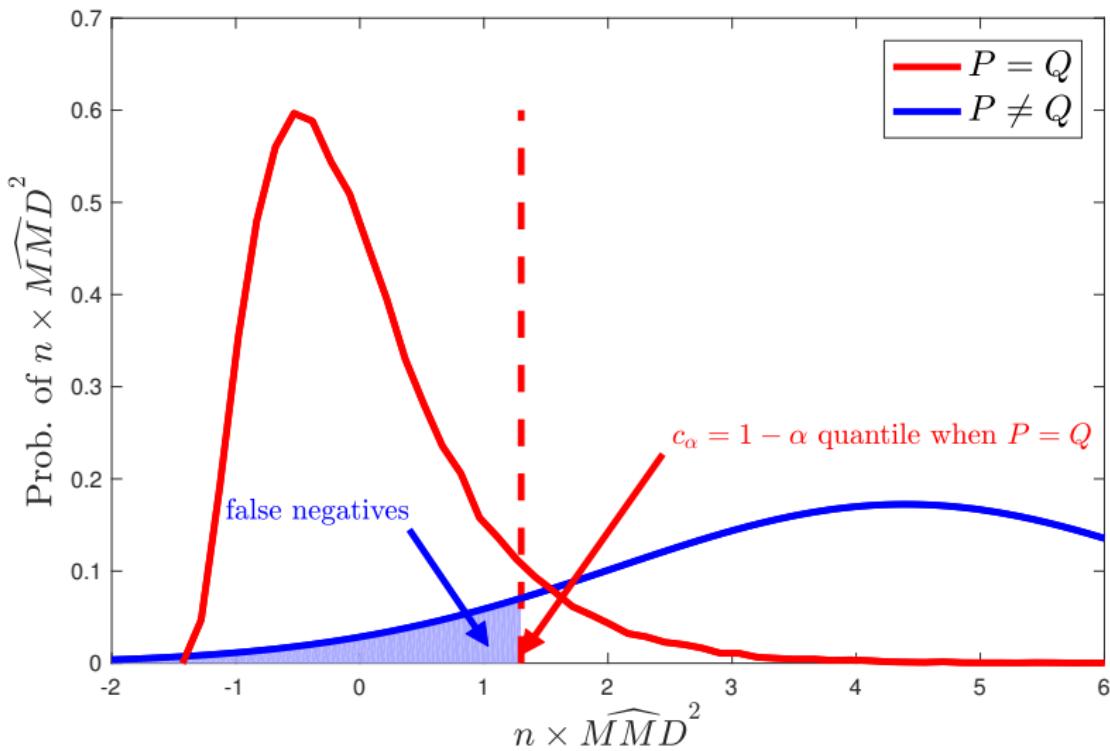
## A statistical test

A summary of the asymptotics:



# A statistical test

**Test construction:** (G., Borgwardt, Rasch, Schoelkopf, and Smola, JMLR 2012)



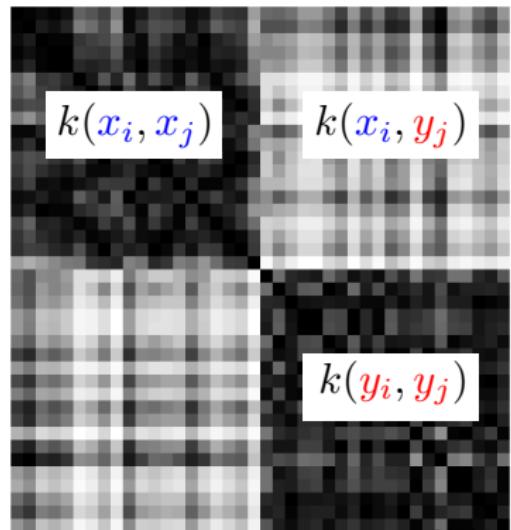
## How do we get test threshold $c_\alpha$ ?

Original empirical MMD for dogs and fish:

$$X = \begin{bmatrix} \text{Basset Hound} & \text{Beagle} & \text{Basset Hound} & \dots \end{bmatrix}$$

$$Y = \begin{bmatrix} \text{Butterfly Fish} & \text{Coral Fish} & \text{Goldfish} & \dots \end{bmatrix}$$

$$\widehat{\text{MMD}}^2 = \frac{1}{n(n-1)} \sum_{i \neq j} k(\mathbf{x}_i, \mathbf{x}_j) + \frac{1}{n(n-1)} \sum_{i \neq j} k(\mathbf{y}_i, \mathbf{y}_j) - \frac{2}{n^2} \sum_{i,j} k(\mathbf{x}_i, \mathbf{y}_j)$$

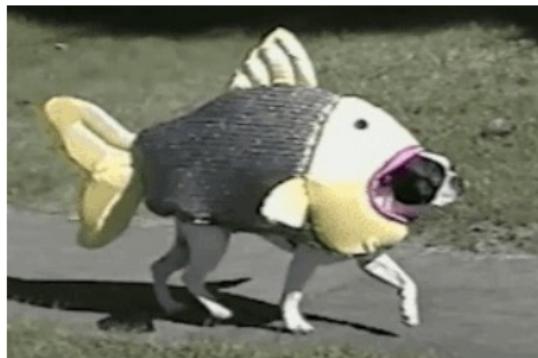


## How do we get test threshold $c_\alpha$ ?

Permuted dog and fish samples (**merdogs**):

$$\tilde{X} = \begin{bmatrix} \text{fish} & \text{dog} & \text{fish} & \dots \end{bmatrix}$$

$$\tilde{Y} = \begin{bmatrix} \text{dog} & \text{fish} & \text{dog} & \dots \end{bmatrix}$$



## How do we get test threshold $c_\alpha$ ?

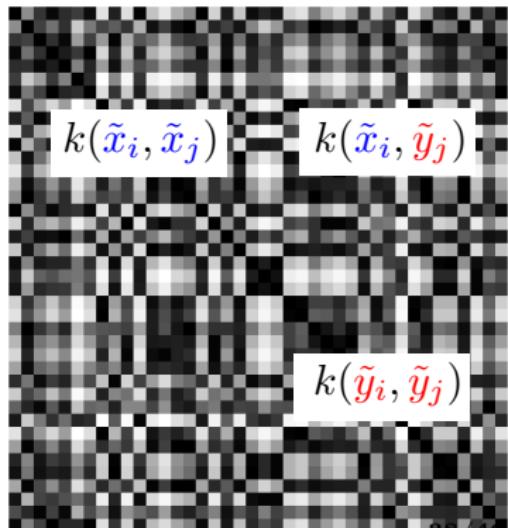
Permuted dog and fish samples (**merdogs**):

$$\tilde{X} = [\text{fish emoji} \quad \text{dog emoji} \quad \text{fish emoji} \quad \dots]$$

$$\tilde{Y} = [\text{dog emoji} \quad \text{fish emoji} \quad \text{dog emoji} \quad \dots]$$

$$\begin{aligned}\widehat{MMD}^2 &= \frac{1}{n(n-1)} \sum_{i \neq j} k(\tilde{x}_i, \tilde{x}_j) \\ &\quad + \frac{1}{n(n-1)} \sum_{i \neq j} k(\tilde{y}_i, \tilde{y}_j) \\ &\quad - \frac{2}{n^2} \sum_{i,j} k(\tilde{x}_i, \tilde{y}_j)\end{aligned}$$

Permutation simulates  
 $P = Q$



How to choose the best kernel:  
optimising the kernel parameters

## The best test for the job

- A test's power depends on  $k(x, x')$ ,  $P$ , and  $Q$  (and  $n$ )
- With characteristic kernel, MMD test has power  $\rightarrow 1$  as  $n \rightarrow \infty$  for any (fixed) problem
  - But, for many  $P$  and  $Q$ , will have terrible power with reasonable  $n$ !

## The best test for the job

- A test's power depends on  $k(x, x')$ ,  $P$ , and  $Q$  (and  $n$ )
- With characteristic kernel, MMD test has power  $\rightarrow 1$  as  $n \rightarrow \infty$  for any (fixed) problem
  - But, for many  $P$  and  $Q$ , will have terrible power with reasonable  $n$ !
- You *can* choose a good kernel for a given problem
- You *can't* get one kernel that has good finite-sample power for all problems
  - No one test can have all that power

## Choosing a kernel for the test

- Simple choice: exponentiated quadratic

$$k(x, y) = \exp\left(-\frac{1}{2\sigma^2}\|x - y\|^2\right)$$

- *Characteristic:* for any  $\sigma$ : for any  $P$  and  $Q$ , power  $\rightarrow 1$  as  $n \rightarrow \infty$

## Choosing a kernel for the test

- Simple choice: exponentiated quadratic

$$k(x, y) = \exp\left(-\frac{1}{2\sigma^2}\|x - y\|^2\right)$$

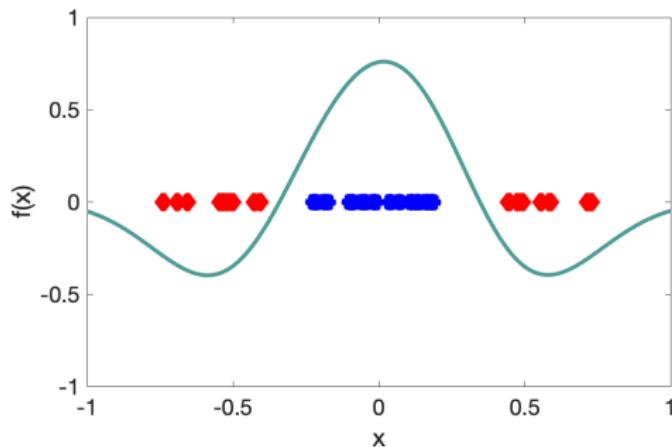
- *Characteristic:* for any  $\sigma$ : for any  $P$  and  $Q$ , power  $\rightarrow 1$  as  $n \rightarrow \infty$
- But choice of  $\sigma$  is very important for finite  $n\dots$

## Choosing a kernel for the test

- Simple choice: exponentiated quadratic

$$k(x, y) = \exp\left(-\frac{1}{2\sigma^2}\|x - y\|^2\right)$$

- *Characteristic:* for any  $\sigma$ : for any  $P$  and  $Q$ , power  $\rightarrow 1$  as  $n \rightarrow \infty$
- But choice of  $\sigma$  is very important for finite  $n\dots$

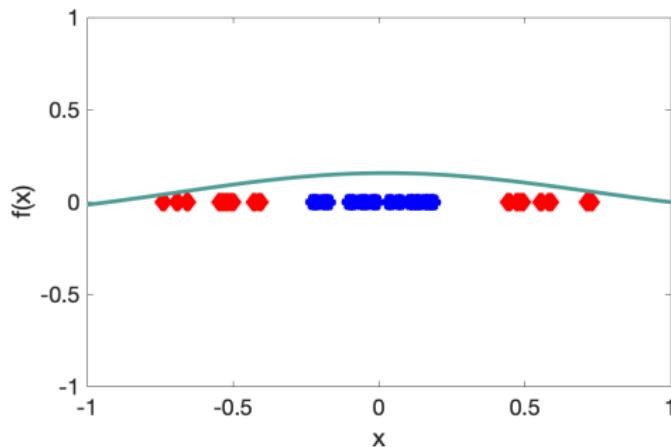


## Choosing a kernel for the test

- Simple choice: exponentiated quadratic

$$k(x, y) = \exp\left(-\frac{1}{2\sigma^2}\|x - y\|^2\right)$$

- *Characteristic:* for any  $\sigma$ : for any  $P$  and  $Q$ , power  $\rightarrow 1$  as  $n \rightarrow \infty$
- But choice of  $\sigma$  is very important for finite  $n\dots$

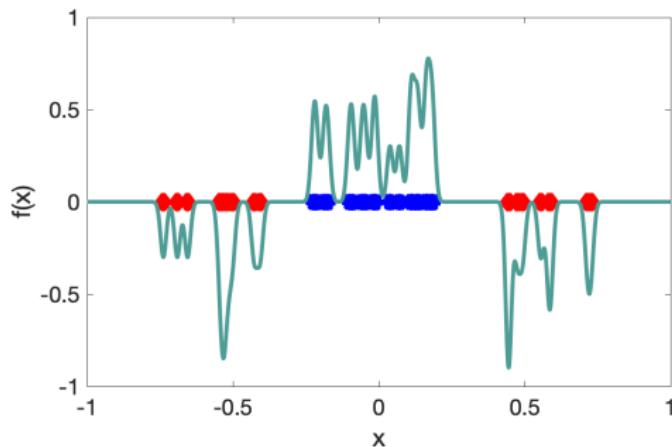


## Choosing a kernel for the test

- Simple choice: exponentiated quadratic

$$k(x, y) = \exp\left(-\frac{1}{2\sigma^2}\|x - y\|^2\right)$$

- *Characteristic:* for any  $\sigma$ : for any  $P$  and  $Q$ , power  $\rightarrow 1$  as  $n \rightarrow \infty$
- But choice of  $\sigma$  is very important for finite  $n\dots$



## Choosing a kernel for the test

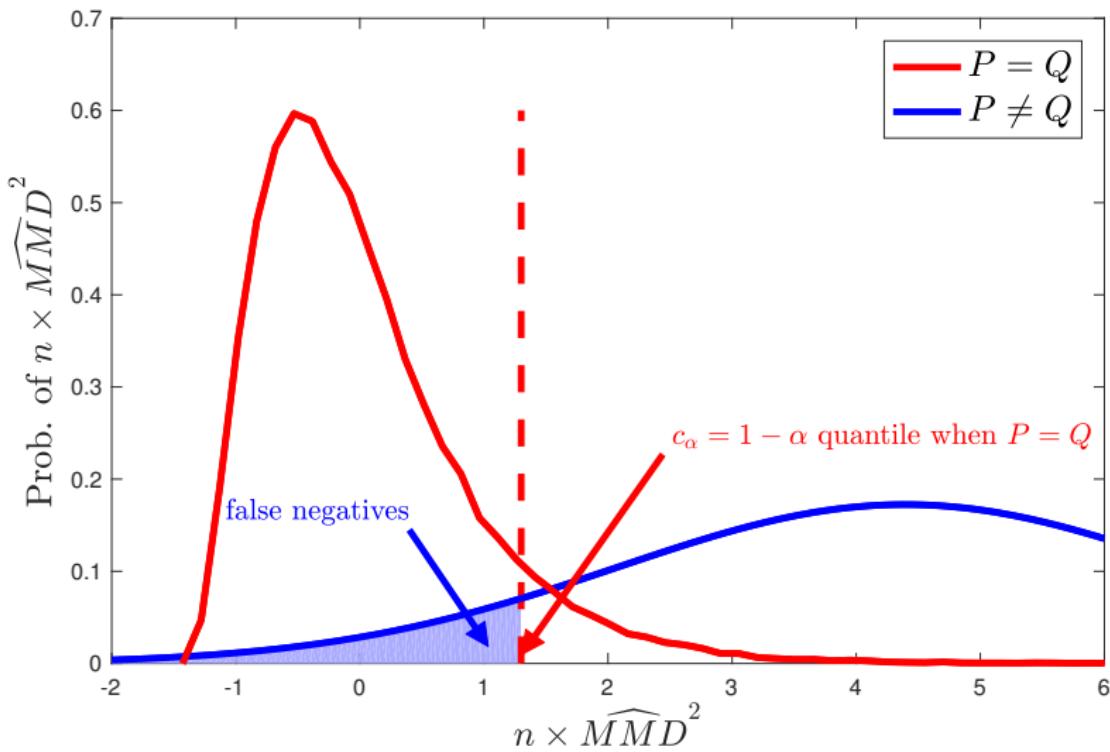
- Simple choice: exponentiated quadratic

$$k(x, y) = \exp\left(-\frac{1}{2\sigma^2}\|x - y\|^2\right)$$

- *Characteristic:* for any  $\sigma$ : for any  $P$  and  $Q$ , power  $\rightarrow 1$  as  $n \rightarrow \infty$
- But choice of  $\sigma$  is very important for finite  $n\dots$
- ... and some problems (e.g. images) might have no good choice for  $\sigma$

## Graphical illustration

- Maximising test power same as minimizing false negatives



## Optimizing kernel for test power

The power of our test ( $\Pr_1$  denotes probability under  $P \neq Q$ ):

$$\Pr_1 \left( n \widehat{\text{MMD}}^2 > \hat{c}_\alpha \right)$$

## Optimizing kernel for test power

The power of our test ( $\Pr_1$  denotes probability under  $P \neq Q$ ):

$$\begin{aligned} & \Pr_1 \left( n \widehat{\text{MMD}}^2 > \hat{c}_\alpha \right) \\ & \rightarrow \Phi \left( \frac{\text{MMD}^2(P, Q)}{\sqrt{V_n(P, Q)}} - \frac{c_\alpha}{n \sqrt{V_n(P, Q)}} \right) \end{aligned}$$

where

- $\Phi$  is the CDF of the standard normal distribution.
- $\hat{c}_\alpha$  is an estimate of  $c_\alpha$  test threshold.

## Optimizing kernel for test power

The power of our test ( $\Pr_1$  denotes probability under  $P \neq Q$ ):

$$\begin{aligned} & \Pr_1 \left( n \widehat{\text{MMD}}^2 > \hat{c}_\alpha \right) \\ & \rightarrow \Phi \left( \underbrace{\frac{\text{MMD}^2(P, Q)}{\sqrt{V_n(P, Q)}}}_{O(n^{1/2})} - \underbrace{\frac{c_\alpha}{n \sqrt{V_n(P, Q)}}}_{O(n^{-1/2})} \right) \end{aligned}$$

For large  $n$ , second term negligible!

## Optimizing kernel for test power

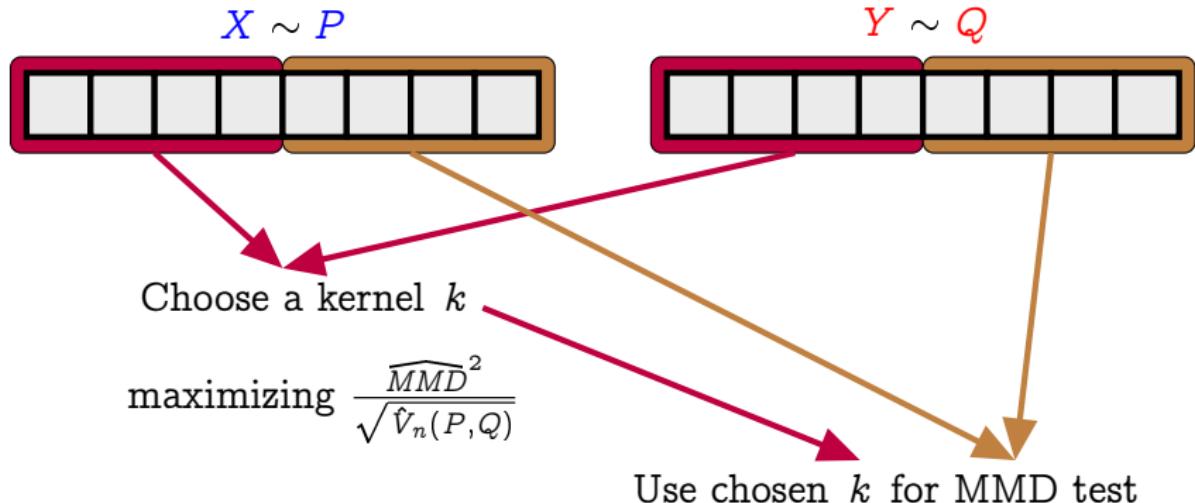
The power of our test ( $\Pr_1$  denotes probability under  $P \neq Q$ ):

$$\begin{aligned} & \Pr_1 \left( n \widehat{\text{MMD}}^2 > \hat{c}_\alpha \right) \\ & \rightarrow \Phi \left( \frac{\text{MMD}^2(P, Q)}{\sqrt{V_n(P, Q)}} - \frac{c_\alpha}{n \sqrt{V_n(P, Q)}} \right) \end{aligned}$$

To maximize test power, maximize

$$\frac{\text{MMD}^2(P, Q)}{\sqrt{V_n(P, Q)}}$$

## Data splitting

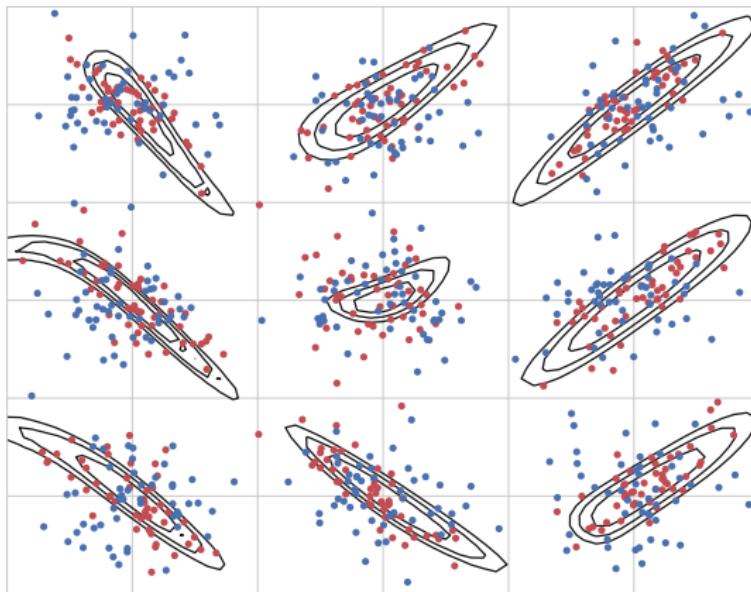


## Learning a kernel helps a lot

Kernel with deep learned features:

$$k_{\theta}(x, y) = [(1 - \epsilon)\kappa(\Phi_{\theta}(x), \Phi_{\theta}(y)) + \epsilon] q(x, y)$$

$\kappa$  and  $q$  are Gaussian kernels



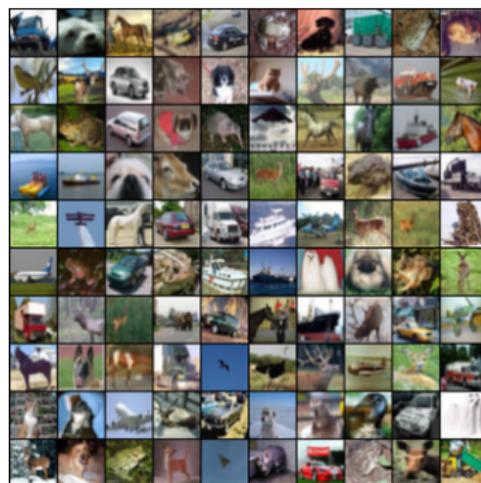
# Learning a kernel helps a lot

Kernel with deep learned features:

$$k_{\theta}(x, y) = [(1 - \epsilon)\kappa(\Phi_{\theta}(x), \Phi_{\theta}(y)) + \epsilon] q(x, y)$$

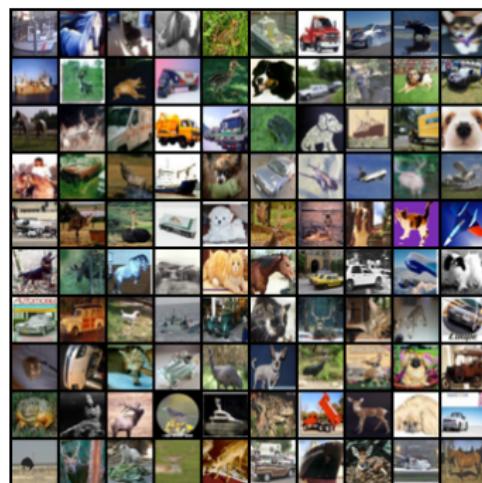
$\kappa$  and  $q$  are Gaussian kernels

- CIFAR-10 vs CIFAR-10.1, **null rejected 75% of time**



CIFAR-10 test set (Krizhevsky 2009)

$$X \sim P$$



CIFAR-10.1 (Recht+ ICML 2019)

$$Y \sim Q$$

# Learning a kernel helps a lot

Kernel with deep learned features:

$$k_{\theta}(x, y) = [(1 - \epsilon)\kappa(\Phi_{\theta}(x), \Phi_{\theta}(y)) + \epsilon] q(x, y)$$

$\kappa$  and  $q$  are Gaussian kernels

- CIFAR-10 vs CIFAR-10.1, **null rejected 75% of time**

arXiv.org > stat > arXiv:2002.09116

Statistics > Machine Learning

[Submitted on 21 Feb 2020]

## Learning Deep Kernels for Non-Parametric Two-Sample Tests

Feng Liu, Wenkai Xu, Jie Lu, Guangquan Zhang, Arthur Gretton, D. J. Sutherland

Accepted to ICML 2020

# Questions?



- A brief introduction to RKHS
- Maximum Mean Discrepancy (MMD)...
  - ...as a difference in feature means
  - ...as an integral probability metric  
*(not just a technicality!)*
- A statistical test based on the MMD