# Controlling Wasserstein distances by maximum mean discrepencies with applications to compressive statistical learning

**Titouan Vayer**
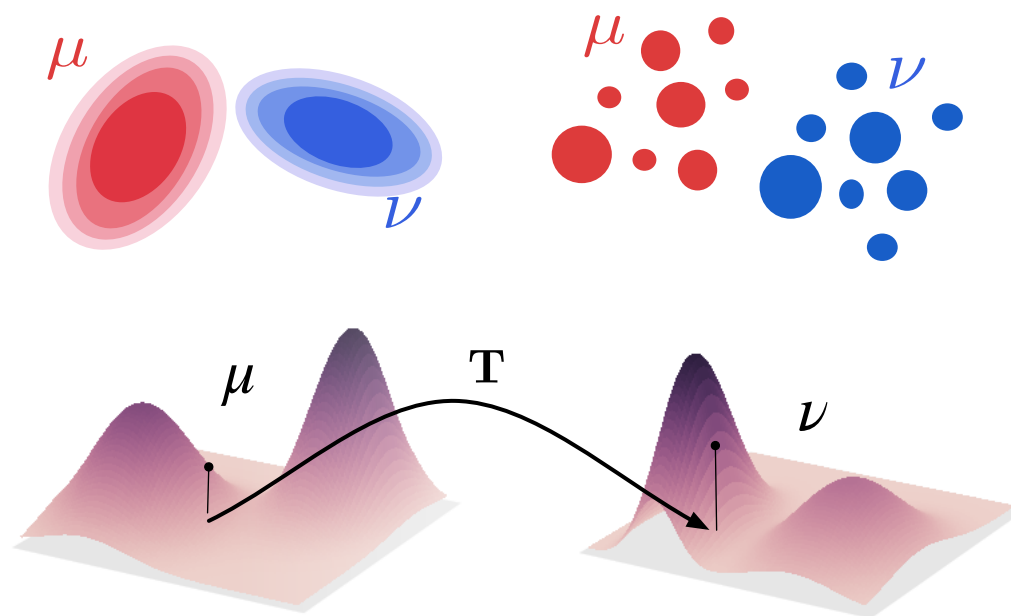
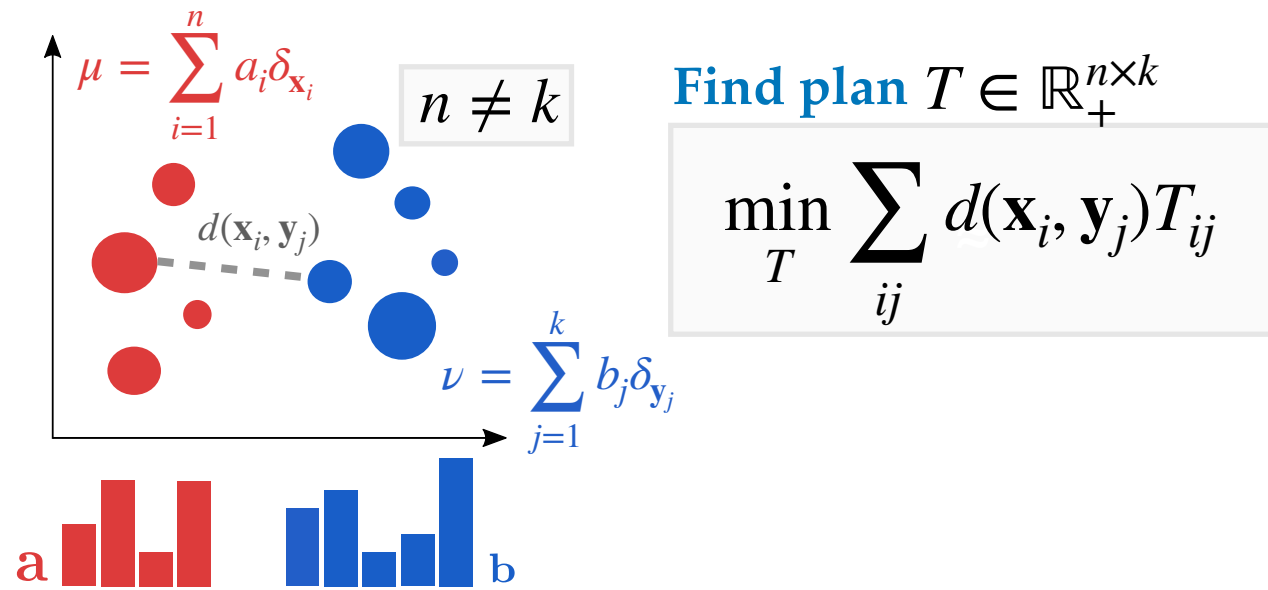**Rémi Gribonval**

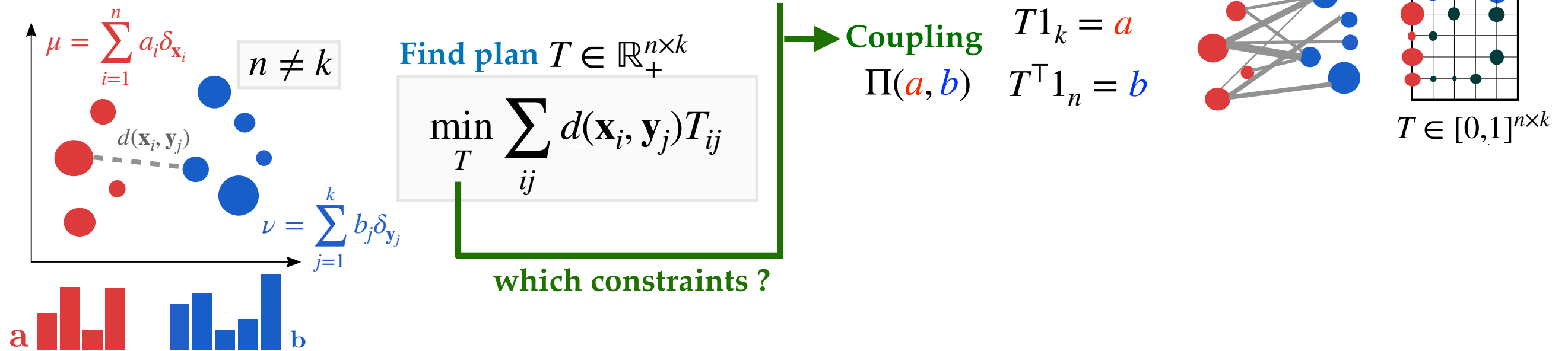# From Optimal Transport to Maximum Mean Discrepancy

# From Wasserstein to MMD

✦ **Classical optimal transport (in a nutshell)**

$\mu = \sum\limits_{i=1}^{n} a_i \delta_{\mathbf{x}_i}$

$n \neq k$

**Find plan** $T \in \mathbb{R}^{n \times k}_{+}$

$$\min_{T} \sum_{ij} d(\mathbf{x}_i, \mathbf{y}_j) T_{ij}$$

$d(\mathbf{x}_i, \mathbf{y}_j)$

$\nu = \sum\limits_{j=1}^{k} b_j \delta_{\mathbf{y}_j}$

$\mathbf{a}$

$\mathbf{b}$

# From Wasserstein to MMD

## ✦ Classical optimal transport (in a nutshell)



$$\mu = \sum_{i=1}^{n} a_i \delta_{\mathbf{x}_i}$$

$$n \neq k$$

$$d(\mathbf{x}_i, \mathbf{y}_j)$$

$$\nu = \sum_{j=1}^{k} b_j \delta_{\mathbf{y}_j}$$

**a**      **b**

**Find plan** $T \in \mathbb{R}_+^{n \times k}$

$$\min_{T} \sum_{ij} d(\mathbf{x}_i, \mathbf{y}_j) T_{ij}$$

**which constraints ?**

**Coupling**    $T 1_k = a$

$\Pi(a, b)$    $T^\top 1_n = b$

$$T \in [0,1]^{n \times k}$$

# From Wasserstein to MMD

✦ **Classical optimal transport (in a nutshell)**



$$\mu = \sum_{i=1}^{n} a_i \delta_{\mathbf{x}_i}$$

$$n \neq k$$

$d(\mathbf{x}_i, \mathbf{y}_j)$

$$\nu = \sum_{j=1}^{k} b_j \delta_{\mathbf{y}_j}$$

**a**     **b**

**Find plan** $T \in \mathbb{R}_+^{n \times k}$

$$\min_{T} \sum_{ij} d(\mathbf{x}_i, \mathbf{y}_j) T_{ij}$$

**which constraints ?**

**Coupling**  $T 1_k = a$

$\Pi(a, b)$   $T^\top 1_n = b$

$T \in [0,1]^{n \times k}$

✦ **Wasserstein distance**  |  $\mu \in \mathscr{P}(X)$
                              |  $\nu \in \mathscr{P}(X)$

$$W_p(\mu, \nu) = \left( \min_{T} \int_{X \times X} \approx d(x, y)^p \, \mathrm{d}T(x, y) \right)^{1/p}$$

✦ It is always **well-defined**

✦ It is a proper distance on $\mathscr{P}(X)$

✦ Lifts the geometry of $X \to \mathscr{P}(X)$

# From Wasserstein to MMD

✦ **Classical optimal transport (in a nutshell)**



$\mu = \sum_{i=1}^{n} a_i \delta_{\mathbf{x}_i}$

$n \neq k$

$d(\mathbf{x}_i, \mathbf{y}_j)$

$\nu = \sum_{j=1}^{k} b_j \delta_{\mathbf{y}_j}$

**a**   **b**

**Find plan** $T \in \mathbb{R}_+^{n \times k}$

$$\min_T \sum_{ij} d(\mathbf{x}_i, \mathbf{y}_j) T_{ij}$$

**which constraints ?**

**Coupling**

$\Pi(a, b)$

$T 1_k = a$

$T^\top 1_n = b$

$T \in [0,1]^{n \times k}$

✦ **In machine learning**

  ✦ Domain adaptation



$\pi \leftarrow \mathrm{TO}(\mu_s, \mu_t)$

  ✦ Generative modeling

  ✦ Analysis of NN convergence

  ✦ ML on graphs, fairness

  ✦ And many other …

✦ **Wasserstein distance**

$\mu \in \mathscr{P}(X)$
$\nu \in \mathscr{P}(X)$

$$W_p(\mu, \nu) = \left( \min_T \int_{X \times X} d(x, y)^p \, \mathrm{d}T(x, y) \right)^{1/p}$$

  ✦ It is always **well-defined**

  ✦ It is a proper distance on $\mathscr{P}(X)$

  ✦ Lifts the geometry of $X \to \mathscr{P}(X)$

# From Wasserstein to MMD

✦ **Kernel theory (in a nutshell)**

$$\kappa : X \times X \to \mathbb{C} \text{ a PSD kernel}$$

# From Wasserstein to MMD

✦ **Kernel theory (in a nutshell)**

$\kappa : X \times X \to \mathbb{C}$ a PSD kernel

$\forall x, y \;\; \kappa(x, y) = \overline{\kappa(y, x)}$

$\forall (x_1, \cdots, x_n), \; K = [\kappa(x_i, x_j)]_{ij} \;$ is PSD

# From Wasserstein to MMD

✦ **Kernel theory (in a nutshell)**     $H_\kappa$ **the RKHS of** $\kappa$

$\kappa : X \times X \to \mathbb{C}$ a PSD kernel  $\iff$  A **Hilbert** space of **functions** from $X \to \mathbb{C}$

# From Wasserstein to MMD

✦ **Kernel theory (in a nutshell)**

$H_\kappa$ **the RKHS of** $\kappa$

$$\kappa : X \times X \to \mathbb{C} \text{ a PSD kernel} \iff \text{A } \textbf{Hilbert} \text{ space of } \textbf{functions} \text{ from } X \to \mathbb{C}$$

$\forall x, \kappa(\,\cdot\,, x) \in H_\kappa$

$f \in H_\kappa \implies \forall x, f(x) = \langle \kappa(\,\cdot\,, x), f \rangle_{H_\kappa}$

$\kappa(x, y) = \langle \kappa(\,\cdot\,, x), \kappa(\,\cdot\,, y) \rangle_{H_\kappa}$

# From Wasserstein to MMD

✦ **Kernel theory (in a nutshell)**

$H_\kappa$ **the RKHS of** $\kappa$

$\kappa : X \times X \to \mathbb{C}$ a PSD kernel $\iff$ A **Hilbert** space of **functions** from $X \to \mathbb{C}$

✦ **Translation invariant kernels**
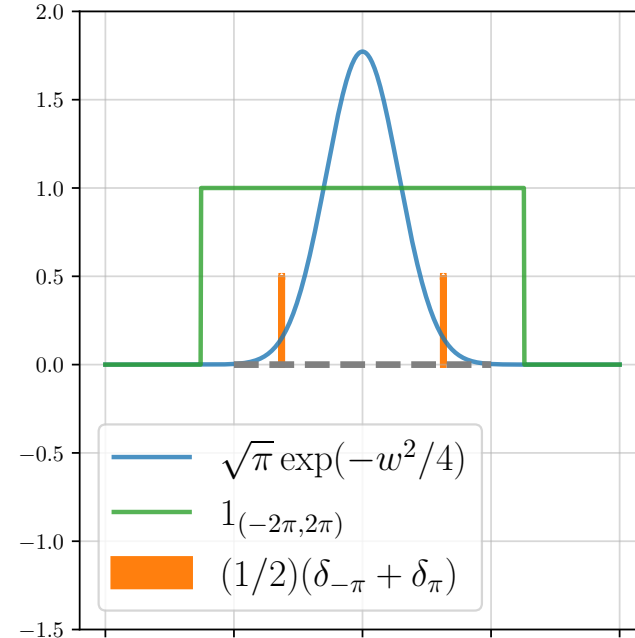
$X = \mathbb{R}^d \quad \kappa(x,y) = \kappa_0(x - y)$

# From Wasserstein to MMD

✦ **Kernel theory (in a nutshell)**

$H_\kappa$ **the RKHS of** $\kappa$

$$\kappa : X \times X \to \mathbb{C} \text{ a PSD kernel} \iff \text{A \textbf{Hilbert} space of \textbf{functions} from } X \to \mathbb{C}$$
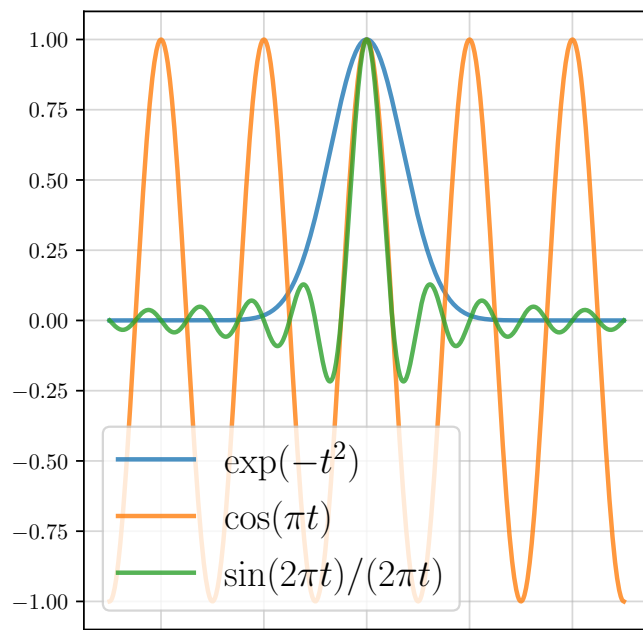
✦ **Translation invariant kernels**

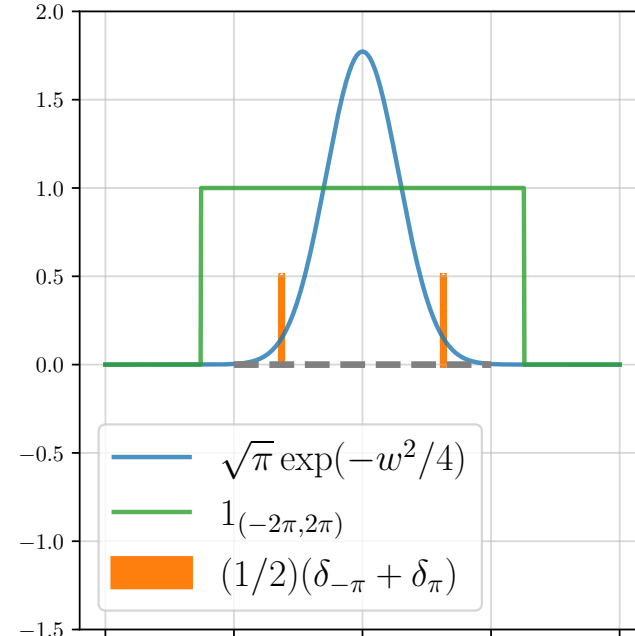$$X = \mathbb{R}^d \quad \kappa(x,y) = \kappa_0(x-y)$$

✦ Is a PSD kernel $\iff \forall \omega, \widehat{\kappa_0}(\omega) \geq 0$

(Bochner)



Functions — legend: $\exp(-t^2)$, $\cos(\pi t)$, $\sin(2\pi t)/(2\pi t)$

Fourier transforms — legend: $\sqrt{\pi}\exp(-w^2/4)$, $1_{(-2\pi, 2\pi)}$, $(1/2)(\delta_{-\pi} + \delta_\pi)$

# From Wasserstein to MMD

✦ **Kernel theory (in a nutshell)**

$H_\kappa$ **the RKHS of** $\kappa$

$$\kappa : X \times X \to \mathbb{C} \text{ a PSD kernel} \iff \text{A } \textbf{Hilbert} \text{ space of } \textbf{functions} \text{ from } X \to \mathbb{C}$$

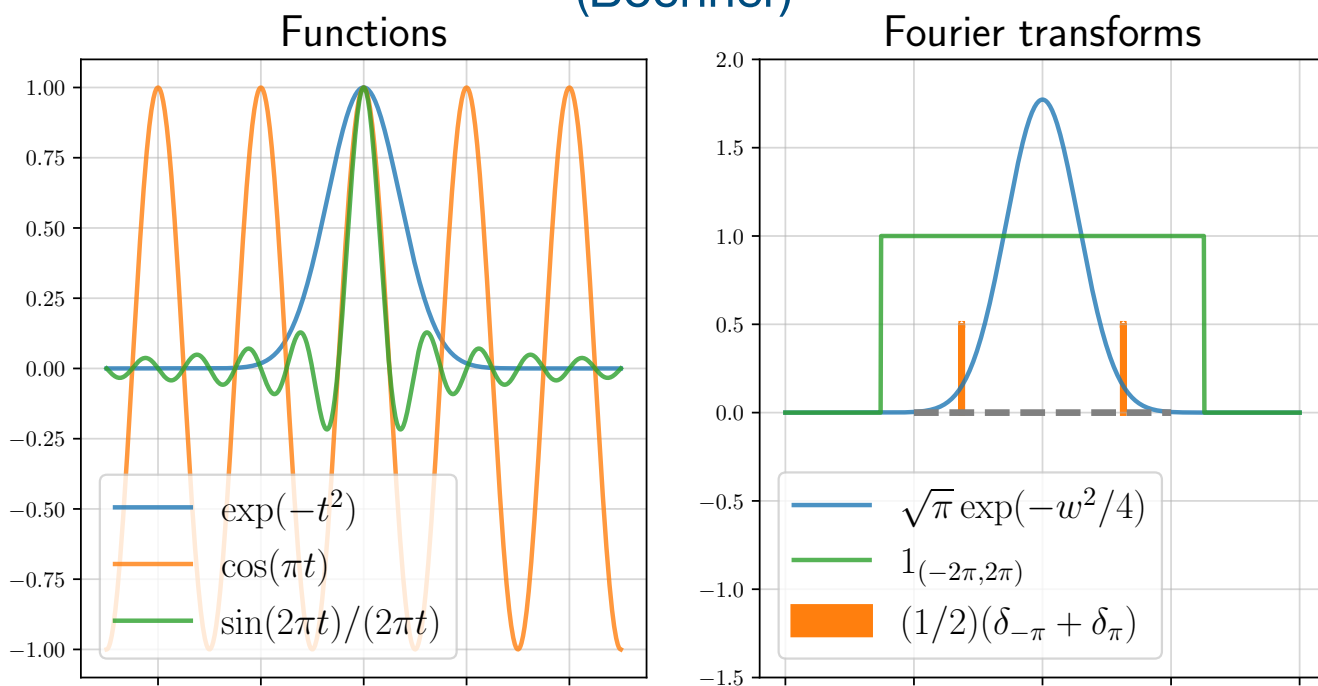✦ **Translation invariant kernels**

$$X = \mathbb{R}^d \quad \kappa(x, y) = \kappa_0(x - y)$$

✦ Is a PSD kernel $\iff \forall \omega, \widehat{\kappa_0}(\omega) \geq 0$

(Bochner)



Functions

Fourier transforms

- $\exp(-t^2)$
- $\cos(\pi t)$
- $\sin(2\pi t)/(2\pi t)$

- $\sqrt{\pi} \exp(-w^2/4)$
- $\mathbb{1}_{(-2\pi, 2\pi)}$
- $(1/2)(\delta_{-\pi} + \delta_\pi)$

(Rahimi, 2017)

✦ With the formula:

**RFF:**

$$\kappa(x, y) = \mathbb{E}_{\omega \sim \Lambda}[e^{-i\langle \omega, x-y \rangle}] \approx \langle \phi(x), \phi(y) \rangle_{\mathbb{R}^m}$$

# From Wasserstein to MMD

**✦ Kernel theory (in a nutshell)**

$H_\kappa$ **the RKHS of** $\kappa$

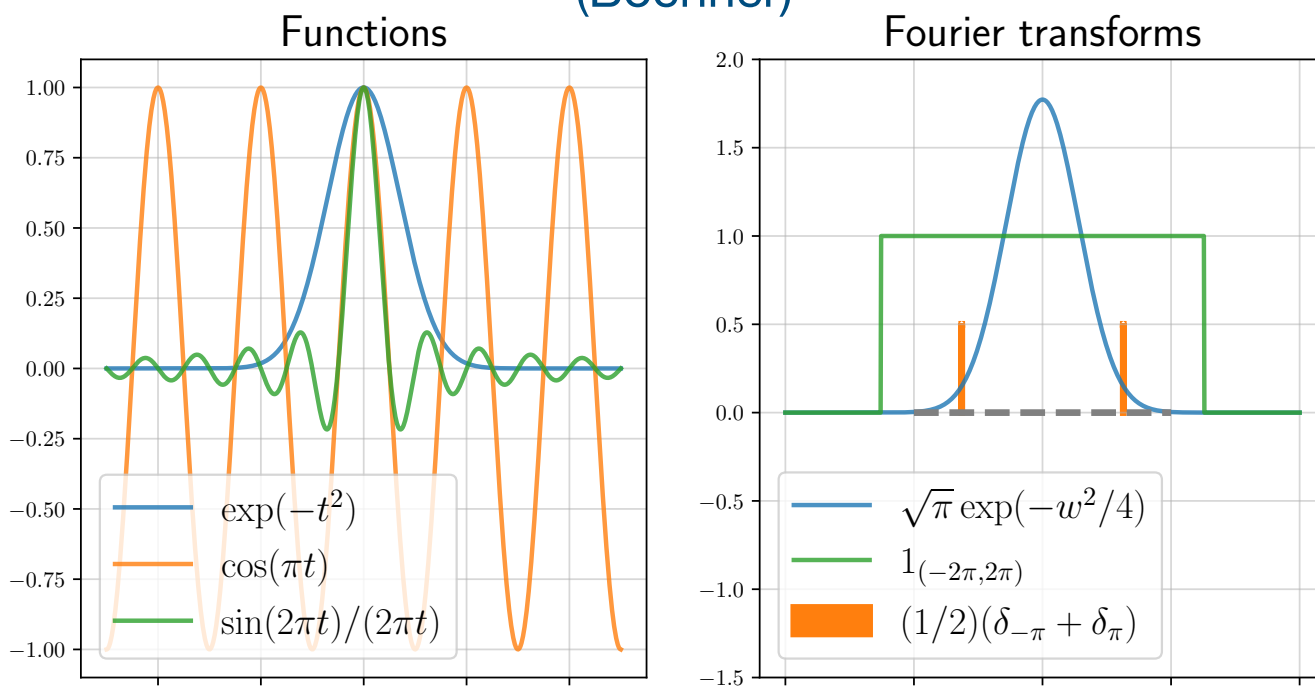$$\kappa : X \times X \to \mathbb{C} \text{ a PSD kernel} \iff \text{A \textbf{Hilbert} space of \textbf{functions} from } X \to \mathbb{C}$$

**✦ Translation invariant kernels**

$$X = \mathbb{R}^d \quad \kappa(x,y) = \kappa_0(x-y)$$

✦ Is a PSD kernel $\iff \forall \omega, \widehat{\kappa_0}(\omega) \geq 0$

(Bochner)

**✦ Maximum mean discrepancy**

$$\mu \in \mathscr{P}(X) \quad \nu \in \mathscr{P}(X)$$

$$\mathrm{MMD}_\kappa(\mu, \nu)$$
$$=$$
$$\| \int_X \kappa(\cdot, x)\mathrm{d}\mu(x) - \int_X \kappa(\cdot, y)\mathrm{d}\nu(y) \|_{H_\kappa}$$



Functions

Fourier transforms

- $\exp(-t^2)$
- $\cos(\pi t)$
- $\sin(2\pi t)/(2\pi t)$

- $\sqrt{\pi}\exp(-w^2/4)$
- $1_{(-2\pi, 2\pi)}$
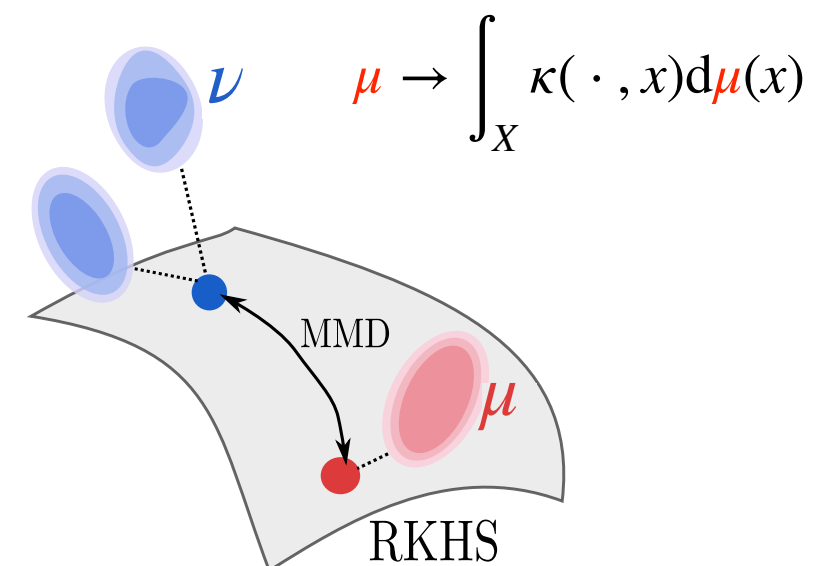- $(1/2)(\delta_{-\pi} + \delta_\pi)$

✦ With the formula:

**RFF:**

$$\kappa(x,y) = \mathbb{E}_{\omega \sim \Lambda}[e^{-i\langle \omega, x-y \rangle}] \approx \langle \phi(x), \phi(y) \rangle_{\mathbb{R}^m}$$

# From Wasserstein to MMD

**✦ Kernel theory (in a nutshell)**

$H_\kappa$ **the RKHS of** $\kappa$

$$\kappa : X \times X \to \mathbb{C} \text{ a PSD kernel} \iff \text{A \textbf{Hilbert} space of \textbf{functions} from } X \to \mathbb{C}$$

**✦ Translation invariant kernels**

$$X = \mathbb{R}^d \quad \kappa(x,y) = \kappa_0(x-y)$$

✦ Is a PSD kernel $\iff \forall \omega, \widehat{\kappa_0}(\omega) \geq 0$

(Bochner)



Functions

Fourier transforms

- $\exp(-t^2)$
- $\cos(\pi t)$
- $\sin(2\pi t)/(2\pi t)$

- $\sqrt{\pi}\exp(-w^2/4)$
- $\mathbb{1}_{(-2\pi, 2\pi)}$
- $(1/2)(\delta_{-\pi} + \delta_\pi)$

✦ With the formula:

$$\kappa(x,y) = \mathbb{E}_{\omega \sim \Lambda}[e^{-i\langle \omega, x-y \rangle}] \approx \langle \phi(x), \phi(y) \rangle_{\mathbb{R}^m}$$

**RFF:**

**✦ Maximum mean discrepancy**

$$\mu \in \mathscr{P}(X) \quad \nu \in \mathscr{P}(X)$$

$$\mathrm{MMD}_\kappa(\mu, \nu)$$
$$=$$
$$\| \int_X \kappa(\cdot, x)\mathrm{d}\mu(x) - \int_X \kappa(\cdot, y)\mathrm{d}\nu(y) \|_{H_\kappa}$$

✦ Distance in the embedding

$$\mu \to \int_X \kappa(\cdot, x)\mathrm{d}\mu(x)$$



$\nu$

$\mu$

MMD

RKHS

# From Wasserstein to MMD

**✦ Kernel theory (in a nutshell)**

$H_\kappa$ **the RKHS of** $\kappa$

$$\kappa : X \times X \to \mathbb{C} \text{ a PSD kernel} \iff \text{A \textbf{Hilbert} space of \textbf{functions} from } X \to \mathbb{C}$$

**✦ Translation invariant kernels**

$$X = \mathbb{R}^d \quad \kappa(x,y) = \kappa_0(x-y)$$

✦ Is a PSD kernel $\iff \forall \omega, \widehat{\kappa_0}(\omega) \geq 0$

(Bochner)



Functions

Fourier transforms

$\exp(-t^2)$
$\cos(\pi t)$
$\sin(2\pi t)/(2\pi t)$

$\sqrt{\pi}\exp(-w^2/4)$
$1_{(-2\pi, 2\pi)}$
$(1/2)(\delta_{-\pi} + \delta_{\pi})$

**✦ Maximum mean discrepancy**

$$\mu \in \mathscr{P}(X) \quad \nu \in \mathscr{P}(X)$$

$$\mathrm{MMD}_\kappa(\mu, \nu)$$
$$=$$
$$\| \int_X \kappa(\,\cdot\,, x)\mathrm{d}\mu(x) - \int_X \kappa(\,\cdot\,, y)\mathrm{d}\nu(y)\|_{H_\kappa}$$

$$\gamma \in \mathscr{M}(X)$$

$$\|\gamma\|_\kappa := \left( \int_{X \times X} \kappa(x,y)\ \mathrm{d}\gamma(x)\mathrm{d}\gamma(y) \right)^{1/2}$$

✦ Semi-norm on $\mathscr{M}(X)$

✦ Alternative formula:

$$\mathrm{MMD}_\kappa(\mu, \nu) = \|\mu - \nu\|_\kappa$$

✦ With the formula:

$$\kappa(x,y) = \mathbb{E}_{\omega \sim \Lambda}[e^{-i\langle\omega, x-y\rangle}] \approx \langle \phi(x), \phi(y) \rangle_{\mathbb{R}^m}$$

# Are they both equivalent ?

$$\forall \mu, \nu : C_1 \cdot W_p(\mu, \nu) \leq \text{MMD}(\mu, \nu) \leq C_2 \cdot W_p(\mu, \nu)$$

# Controlling MMDs by Wasserstein distances
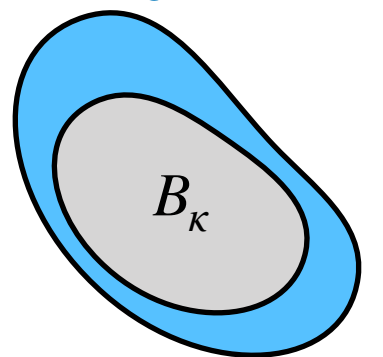
✦ **Can we find** $C > 0$ **such that**

$$(\,\star\,) \quad \forall \mu, \nu : \mathrm{MMD}_\kappa(\mu, \nu) \leq C \cdot \mathrm{W}_p(\mu, \nu)$$

# Controlling MMDs by Wasserstein distances

✦ **Can we find $C > 0$ such that**

$$(\star) \quad \forall \textcolor{red}{\mu}, \textcolor{blue}{\nu} : \mathrm{MMD}_\kappa(\textcolor{red}{\mu}, \textcolor{blue}{\nu}) \leq C \cdot \mathrm{W}_p(\textcolor{red}{\mu}, \textcolor{blue}{\nu})$$

**A characterization** $\qquad \exists C > 0 \ \text{ such that } (\star)$

$$\Longleftrightarrow \quad \forall x, y \ \sqrt{\kappa(x,x) + \kappa(y,y) - 2\kappa(x,y)} \leq C \cdot d(x,y)$$

$\mathrm{Lip}_C(X, \mathbb{R})$

$B_\kappa$

# Controlling MMDs by Wasserstein distances

✦ **Can we find $C > 0$ such that**

$$(\star) \quad \forall \mu, \nu : \mathrm{MMD}_\kappa(\mu, \nu) \leq C \cdot \mathrm{W}_p(\mu, \nu)$$

**A characterization** $\quad \exists C > 0$ such that $(\star)$

$$\Longleftrightarrow \quad \forall x, y \; \sqrt{\kappa(x,x) + \kappa(y,y) - 2\kappa(x,y)} \leq C \cdot d(x,y)$$

$\mathrm{Lip}_C(X, \mathbb{R})$

$B_\kappa$

✦ **Corollary for TI kernels** $\quad \kappa(x, y) = \kappa_0(x - y) \quad d(x, y) = \|x - y\|_2$

$(\star)$ **always holds** with $C = \kappa_0(0)\sqrt{\lambda_{\max}(-\nabla^2[\kappa_0](0))}$

✦ The MMD is a weaker notion of metric for smooth TI kernel

✦ This direction is easy !!

# Controlling Wasserstein distances by MMDs

✦ **Can we find** $C > 0$ **such that**

$$(\star\star) \quad \forall \mu, \nu : \mathrm{W}_p(\mu, \nu) \leq C \cdot \mathrm{MMD}_\kappa(\mu, \nu)$$

# Controlling Wasserstein distances by MMDs

✦ **Can we find $C > 0$ such that**

$$(\star \star) \quad \forall \mu, \nu : W_p(\mu, \nu) \leq C \cdot \mathrm{MMD}_\kappa(\mu, \nu)$$

> **Not without any assumption !!**

✦ e.g : if $\kappa$ bounded then $\forall \mu, \nu \in \mathscr{P}(X) \quad \mathrm{MMD}_\kappa(\mu, \nu) \leq \mathrm{cte}$

✦ but not Wasserstein distances !

# Controlling Wasserstein distances by MMDs

$\mathscr{P}(X)$

★ **Can we find** $C_\Sigma > 0$ **such that**

$$(\star\,\star\,\star)\ \forall\mu, \nu \in \Sigma : \mathrm{W}_p(\mu, \nu) \leq C_\Sigma \cdot \mathrm{MMD}_\kappa(\mu, \nu)$$

**model set**

$\Sigma$

$\mu$

$\nu$

# Controlling Wasserstein distances by MMDs

$\mathscr{P}(X)$

**✦ Can we find** $C_\Sigma > 0$ **such that**

**model set**

$(\star\star\star)\ \forall \mu, \nu \in \Sigma : W_p(\mu, \nu) \leq C_\Sigma \cdot \mathrm{MMD}_\kappa(\mu, \nu)$
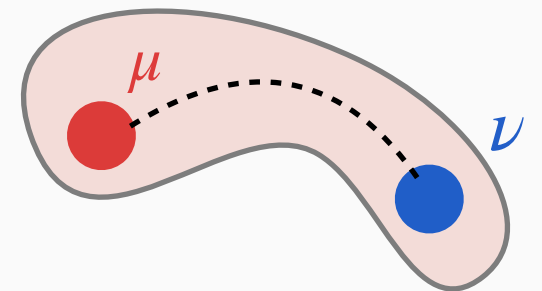


✦ If $\kappa$ bounded then necessarily $\Sigma$ must be bounded

$(\star\star\star) \implies \sup_{\mu,\nu \in \Sigma} \| \mathrm{mean}(\mu) - \mathrm{mean}(\nu) \|_2 < +\infty$

# Controlling Wasserstein distances by MMDs

$\mathscr{P}(X)$

✦ **Can we find** $C_\Sigma > 0$ **such that**

$$(\star\star\star)\ \forall \mu, \nu \in \Sigma : W_p(\mu, \nu) \leq C_\Sigma \cdot \text{MMD}_\kappa(\mu, \nu)$$

**model set**

$\Sigma$

$\mu$   $\nu$

✦ If $\kappa$ bounded then necessarily $\Sigma$ must be bounded

$$(\star\star\star) \implies \sup_{\mu,\nu \in \Sigma} \|\text{mean}(\mu) - \text{mean}(\nu)\|_2 < +\infty$$

✦ If $\Sigma$ contains $[\mu, \nu]$ with $\text{supp}(\mu) \cap \text{supp}(\nu) = \varnothing$

$\mu$   $\nu$

$$(\star\star\star) \text{ impossible with } p > 1$$

**We must find a « larger » definition**

# Controlling Wasserstein distances by MMDs

✦ **Definition: embeddability** $\Sigma \subset \mathscr{P}(X), \delta \in [0,1]$

$$(\Sigma, W_p) \text{ is } (\kappa, \delta) - \text{embeddable when}$$

$$\exists C > 0, \forall \mu, \nu \in \Sigma : W_p(\mu, \nu) \leq C \cdot \text{MMD}_\kappa^\delta(\mu, \nu)$$

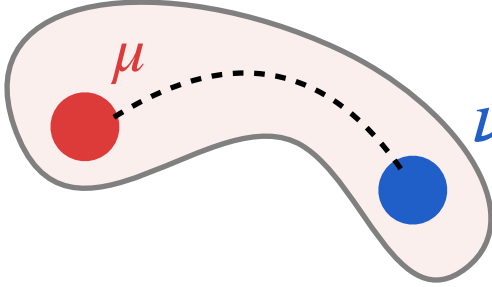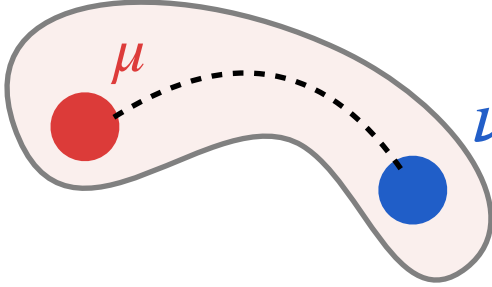# Controlling Wasserstein distances by MMDs

✦ **Definition: embeddability** $\Sigma \subset \mathscr{P}(X), \delta \in [0,1]$

$(\Sigma, W_p)$ is $(\kappa, \delta)-$ embeddable when

$$\exists C > 0, \forall \mu, \nu \in \Sigma : W_p(\mu, \nu) \leq C \cdot \text{MMD}_\kappa^\delta(\mu, \nu)$$

✦ **Some necessary conditions:**

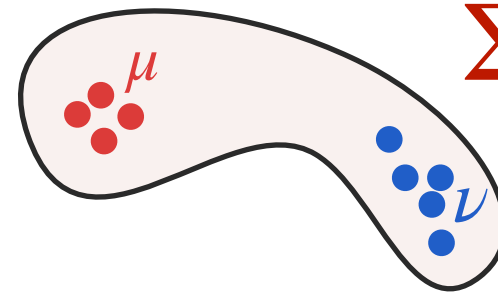✦ If  $(\kappa, \delta)-$ embeddable $\implies \delta \leq 1/p$

# Controlling Wasserstein distances by MMDs

✦ **Definition: embeddability** $\Sigma \subset \mathscr{P}(X), \delta \in [0,1]$

$$(\Sigma, W_p) \text{ is } (\kappa, \delta) - \text{embeddable when}$$

$$\exists C > 0, \forall \mu, \nu \in \Sigma : W_p(\mu, \nu) \leq C \cdot \text{MMD}_\kappa^\delta(\mu, \nu)$$

✦ **Some necessary conditions:**

✦ If  $(\kappa, \delta) - \text{embeddable} \implies \delta \leq 1/p$

✦ If $\Sigma = \{\mu \in \mathscr{P}(X) : \mu(B(0,R)) = 1\}$ and $\kappa$ bounded

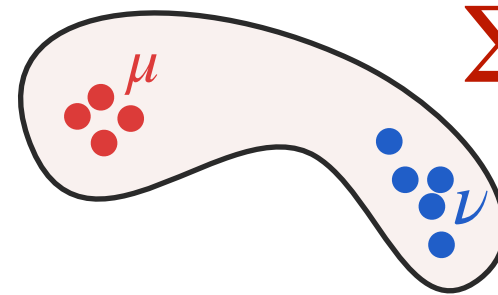$$(\kappa, \delta) - \text{embeddable} \implies \delta \leq 2/d$$

# Controlling Wasserstein distances by MMDs

✦ **Definition: embeddability** $\Sigma \subset \mathscr{P}(X), \delta \in [0,1]$

$(\Sigma, W_p)$ is $(\kappa, \delta)-$ embeddable when

$$\exists C > 0, \forall \mu, \nu \in \Sigma : W_p(\mu, \nu) \leq C \cdot \mathrm{MMD}_\kappa^\delta(\mu, \nu)$$

✦ **Some necessary conditions:**

$\Sigma$ discrete distributions **with K atoms**

$$\Sigma = \{ \sum_{i=1}^{K} a_i \delta_{\mathbf{x}_i} : \mathbf{x}_i \in B(0,R) \}$$

$$\kappa(x, y) = \kappa_0(x - y) \text{ with } \kappa_0 \text{ smooth}$$

# Controlling Wasserstein distances by MMDs

✦ **Definition: embeddability** $\Sigma \subset \mathcal{P}(X), \delta \in [0,1]$

$(\Sigma, W_p)$ is $(\kappa, \delta) - $ embeddable when

$$\exists C > 0, \forall \mu, \nu \in \Sigma : W_p(\mu, \nu) \leq C \cdot \mathrm{MMD}_\kappa^\delta(\mu, \nu)$$

✦ **Some necessary conditions:**



$\Sigma$ discrete distributions **with K atoms**

$$\Sigma = \{ \sum_{i=1}^K a_i \delta_{\mathbf{x}_i} : \mathbf{x}_i \in B(0,R) \}$$

$\kappa(x, y) = \kappa_0(x - y)$ with $\kappa_0$ smooth

$$(\kappa, \delta) - \text{embeddable} \implies \delta \leq 2/K$$

# Controlling Wasserstein distances by MMDs

✦ **Definition: embeddability** $\Sigma \subset \mathscr{P}(X), \delta \in [0,1]$

$(\Sigma, W_p)$ is $(\kappa, \delta) - $ embeddable when

$$\exists C > 0, \forall \mu, \nu \in \Sigma : W_p(\mu, \nu) \leq C \cdot \mathrm{MMD}_\kappa^\delta(\mu, \nu)$$

We **<u>cannot</u>** control Wass by MMD **uniformly over all discrete distrib.** (even in a compact) **for a smooth TI kernel**

# Controlling Wasserstein distances by MMDs

✦ **Some positive results**

   ✦ Distrib. with density in Sobolev space + bounded moments

$$\Sigma = \{\mu = f\mathrm{d}\mathbf{x} : \|f\|_s \leq B, \mathrm{M}_r(\mu) \leq M\}$$

with annotations: $\int |\partial^{|s|}f(\mathbf{x})|^2 \mathrm{d}\mathbf{x}$ and $\int \|\mathbf{x}\|_2^r \mathrm{d}\mu(\mathbf{x})$

# Controlling Wasserstein distances by MMDs

✦ **Some positive results**

   ✦ Distrib. with <span style="color:green">density in Sobolev space</span> + <span style="color:orange">bounded moments</span>

$$\int |\partial^{|s|} f(\mathbf{x})|^2 \, d\mathbf{x} \qquad \int \|\mathbf{x}\|_2^r \, d\mu(\mathbf{x})$$

$$\Sigma = \{\mu = f \, d\mathbf{x} : \|f\|_s \leq B, \mathbf{M}_r(\mu) \leq M\}$$

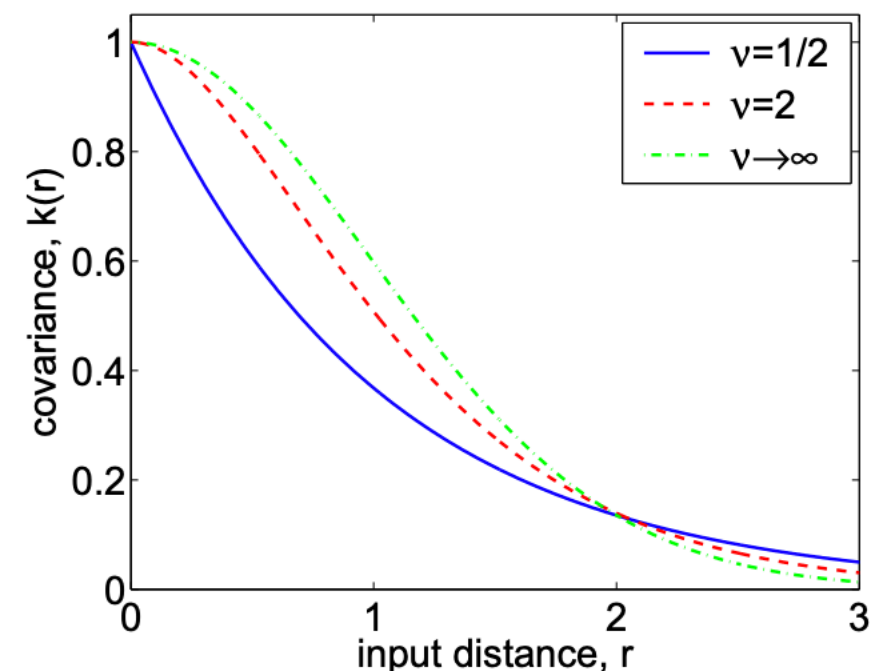   ✦ Any TI, PSD kernels $\kappa(x, y) = \kappa_0(x - y)$ <span style="color:blue">+ some regularity</span>

| ~~Gaussian~~ $\qquad 1/\widehat{\kappa_0}(\omega) = O_{\omega \to +\infty}(\|\omega\|^{2s})$

| **Matérn class**, splines,
polyharmonic curves

# Controlling Wasserstein distances by MMDs

✦ **Some positive results**

  ✦ Distrib. with density in Sobolev space + bounded moments

$$\int |\partial^{|s|} f(\mathbf{x})|^2 \, \mathrm{d}\mathbf{x} \qquad \int \|\mathbf{x}\|_2^r \mathrm{d}\mu(\mathbf{x})$$

$$\Sigma = \{\mu = f\mathrm{d}\mathbf{x} : \|f\|_s \leq B, \mathrm{M}_r(\mu) \leq M\}$$

  ✦ Any TI, PSD kernels $\kappa(x, y) = \kappa_0(x - y)$ + some regularity

    | Gaussian $\qquad 1/\widehat{\kappa_0}(\omega) = O_{\omega \to +\infty}(\|\omega\|^{2s})$

    | **Matérn class**, splines, polyharmonic curves

$(\Sigma, \mathrm{W}_p)$ is $(\kappa, \delta) -$ embeddable

with $\delta = \dfrac{r - p}{p(d + 2r)}$

# Controlling Wasserstein distances by MMDs

✦ **Some positive results**

✦ Distrib. with density in Sobolev space + bounded moments

$$\int |\partial^{|s|} f(\mathbf{x})|^2 \, d\mathbf{x} \qquad \int \|\mathbf{x}\|_2^r d\mu(\mathbf{x})$$

$$\Sigma = \{\mu = f d\mathbf{x} : \|f\|_s \leq B, \mathbf{M}_r(\mu) \leq M\}$$

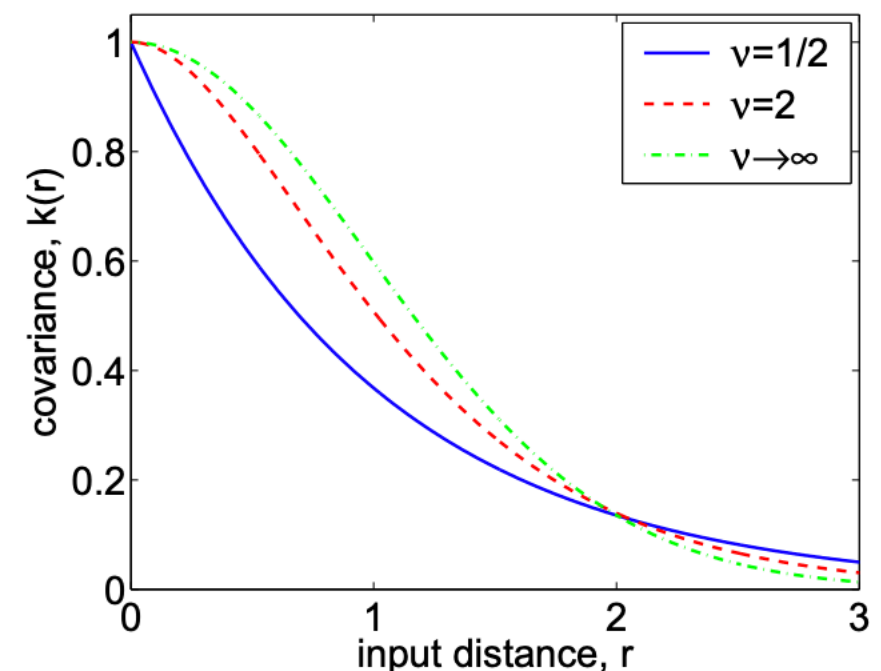✦ Any TI, PSD kernels $\kappa(x, y) = \kappa_0(x - y)$ + some regularity

$$1/\widehat{\kappa_0}(\omega) = O_{\omega \to +\infty}(\|\omega\|^{2s})$$

| Gaussian

| **Matérn class**, splines, polyharmonic curves

$(\Sigma, W_p)$ is $(\kappa, \delta) - $ embeddable

with $\delta \approx \dfrac{1}{2p}$ **($r$ big)**

# Controlling Wasserstein distances by MMDs

✦ **Some positive results**

    ✦ Distrib. with smooth densities + compact support

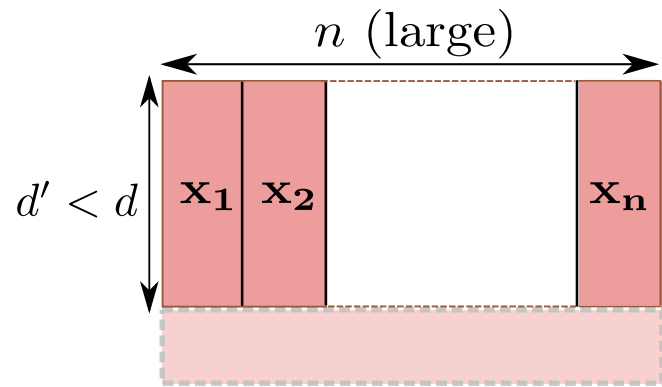    ✦ Any TI, PSD kernels $\kappa(x, y) = \kappa_0(x - y)$   + Matérn class

$$(\Sigma, W_p) \text{ is } (\kappa, \delta) - \text{embeddable with } \delta = 1/2p$$

# Controlling Wasserstein distances by MMDs

✦ **Some positive results**

   ✦ Distrib. with smooth densities + compact support

   ✦ Any TI, PSD kernels $\kappa(x, y) = \kappa_0(x - y)$ + Matérn class

$$(\Sigma, W_p) \text{ is } (\kappa, \delta) - \text{embeddable with } \delta = 1/2p$$
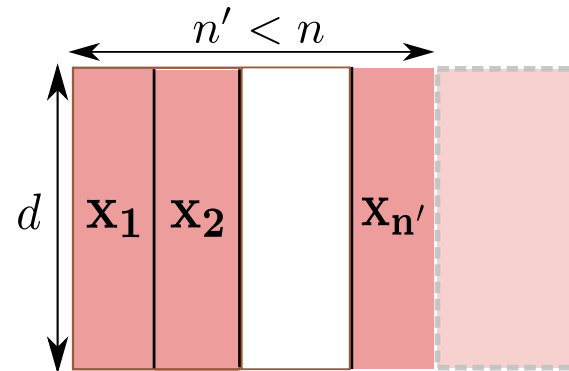
✦ **Other results**

   ✦ Larger class of distrib. / kernels if we allow an error $\eta > 0$

   ✦ For unbounded + conditionally PSD kernels (Chafaï, 2016)

   ✦ Other connections (Modeste, 2022), (Goldfeld, 2020)
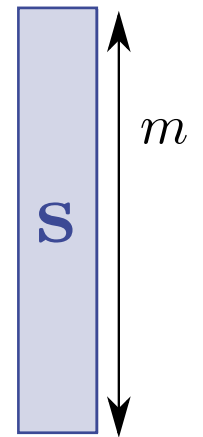
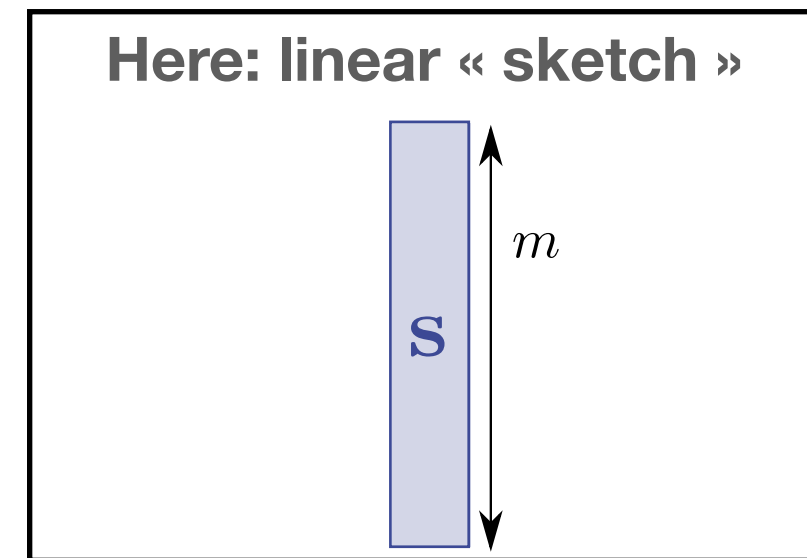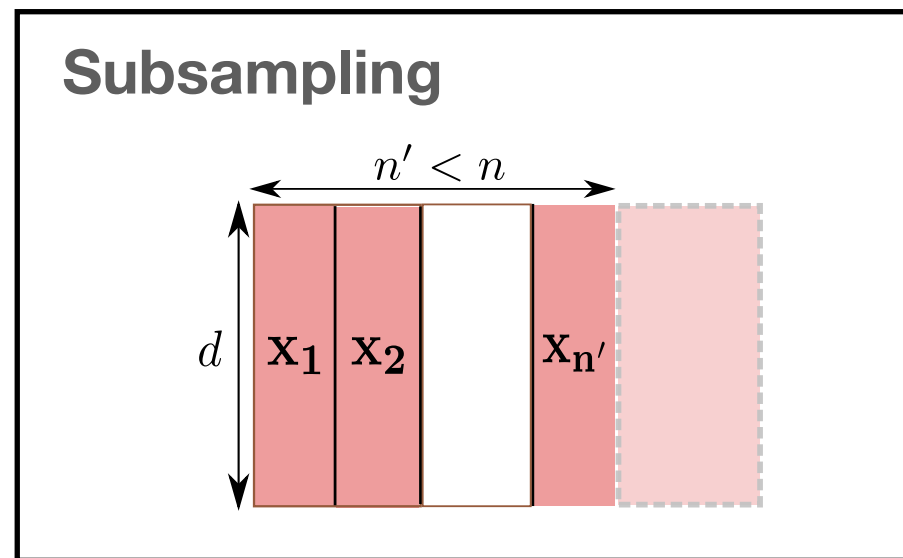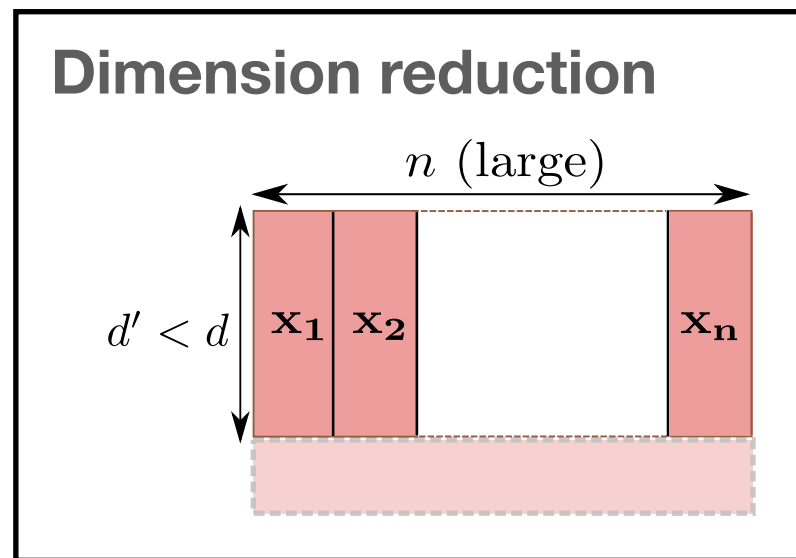# Motivations: compressive learning



Dimension reduction

$n$ (large)

$d' < d$

$\mathbf{x_1}$ $\mathbf{x_2}$ $\mathbf{x_n}$

Subsampling

$n' < n$

$d$

$\mathbf{x_1}$ $\mathbf{x_2}$ $\mathbf{x_{n'}}$
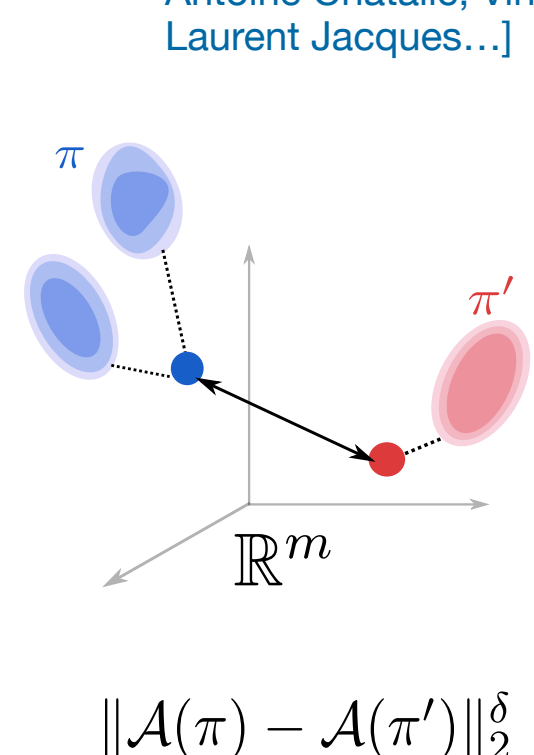
Here: linear « sketch »
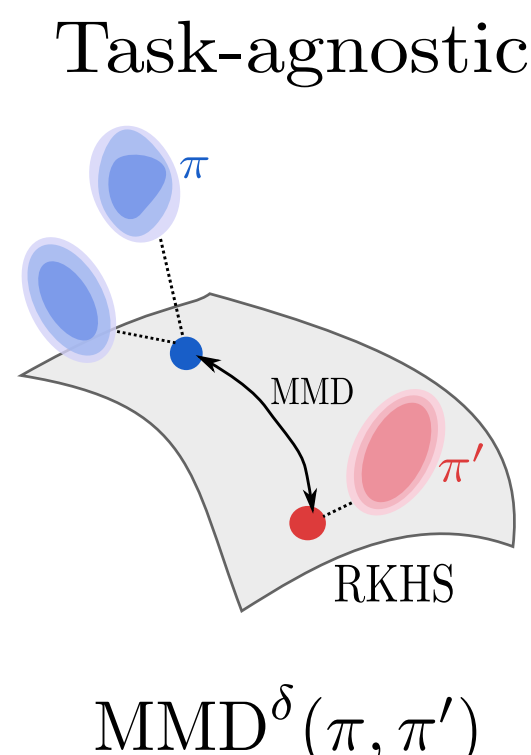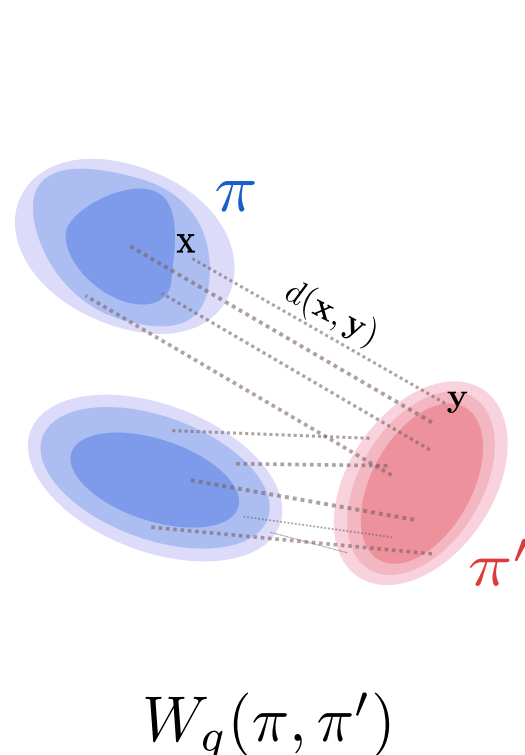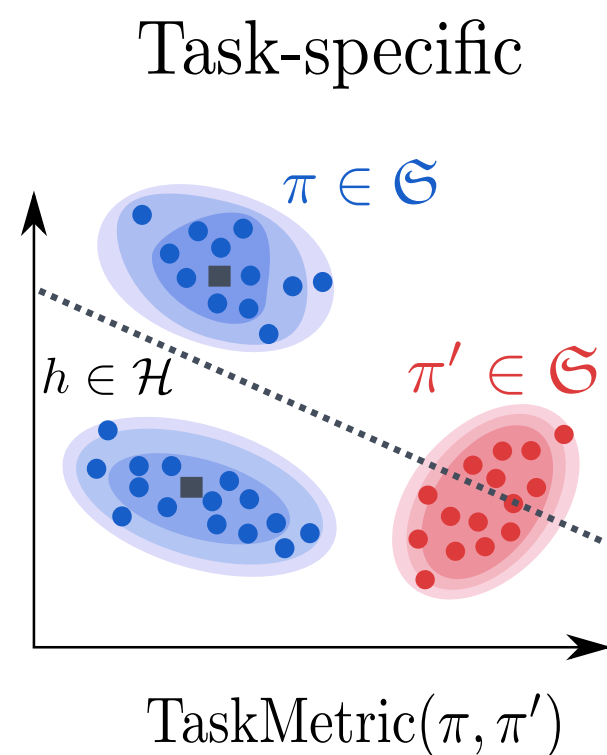
$\mathbf{S}$

$m$

[Rémi Gribonval, Gilles Blanchard, Nicolas Keriven, Yann Traonmilin, Antoine Chatalic, Vincent Schellekens, Laurent Jacques…]

# Motivations: compressive learning



**Dimension reduction**

$n$ (large)

$d' < d$

$\mathbf{x_1}$ $\mathbf{x_2}$ $\mathbf{x_n}$

**Subsampling**

$n' < n$

$d$

$\mathbf{x_1}$ $\mathbf{x_2}$ $\mathbf{x_{n'}}$

**Here: linear « sketch »**

$\mathbf{S}$

$m$

[Rémi Gribonval, Gilles Blanchard, Nicolas Keriven, Yann Traonmilin, Antoine Chatalic, Vincent Schellekens, Laurent Jacques...]

Task-specific

$\pi \in \mathfrak{S}$

$h \in \mathcal{H}$

$\pi' \in \mathfrak{S}$

$\mathrm{TaskMetric}(\pi, \pi')$

Task-agnostic

$\pi$

$\mathbf{x}$

$d(\mathbf{x}, \mathbf{y})$

$\mathbf{y}$

$\pi'$

$W_q(\pi, \pi')$

$\pi$

MMD

$\pi'$

RKHS

$\mathrm{MMD}^{\delta}(\pi, \pi')$

$\pi$

$\pi'$

$\mathbb{R}^m$

$\|\mathcal{A}(\pi) - \mathcal{A}(\pi')\|_2^{\delta}$

Wasserstein Learnability

Kernel Hölder LRIP

CSL guarantees

Hölder LRIP