

# Minimax estimation on manifolds via optimal transport

Séminaire des doctorants proba-stats  
3 juin 2020

Vincent Divol

DataShape - Inria Saclay

LMO - Université Paris-Sud

[vincent.divol@inria.fr](mailto:vincent.divol@inria.fr)

# Density estimation: the basics

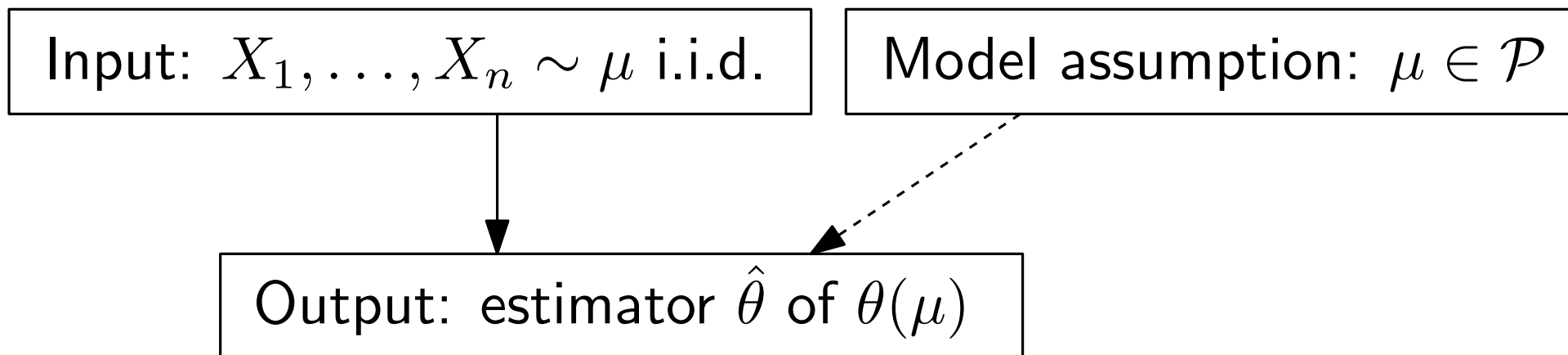
---

A (very) concise summary:

# Density estimation: the basics

---

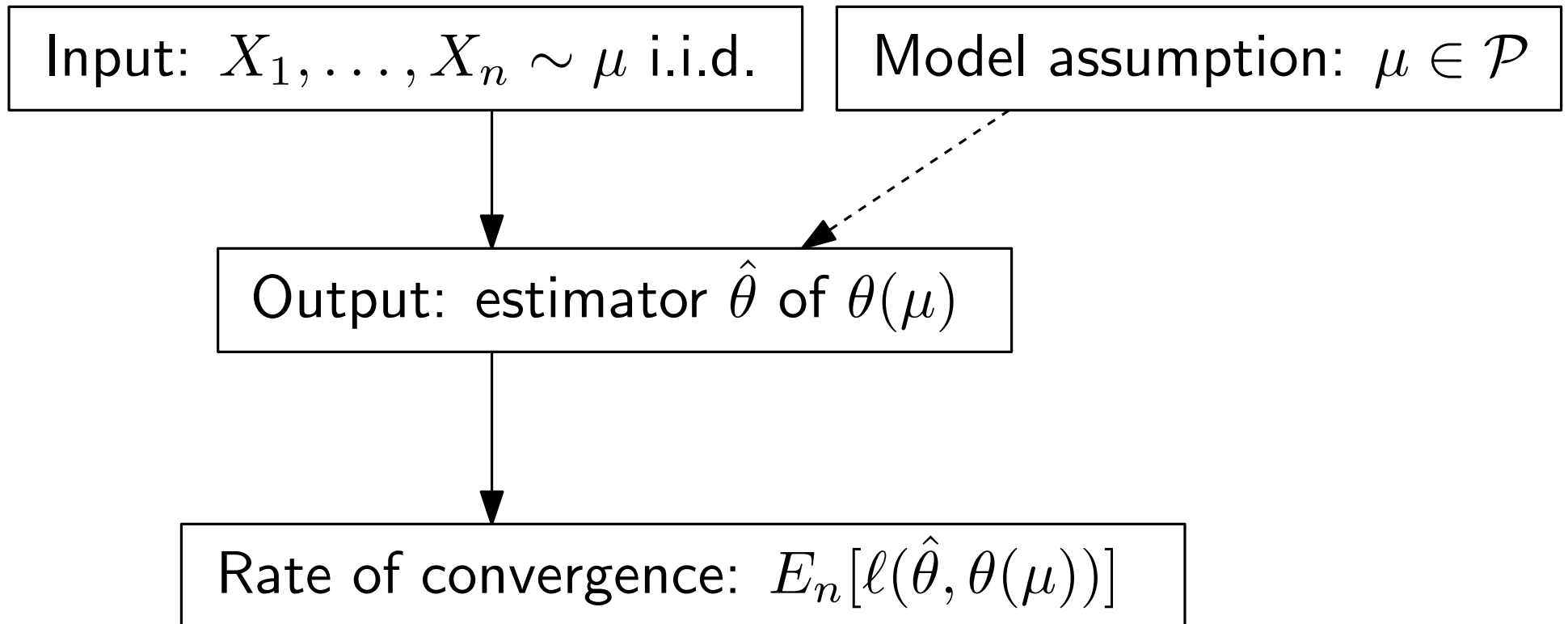
A (very) concise summary:



# Density estimation: the basics

---

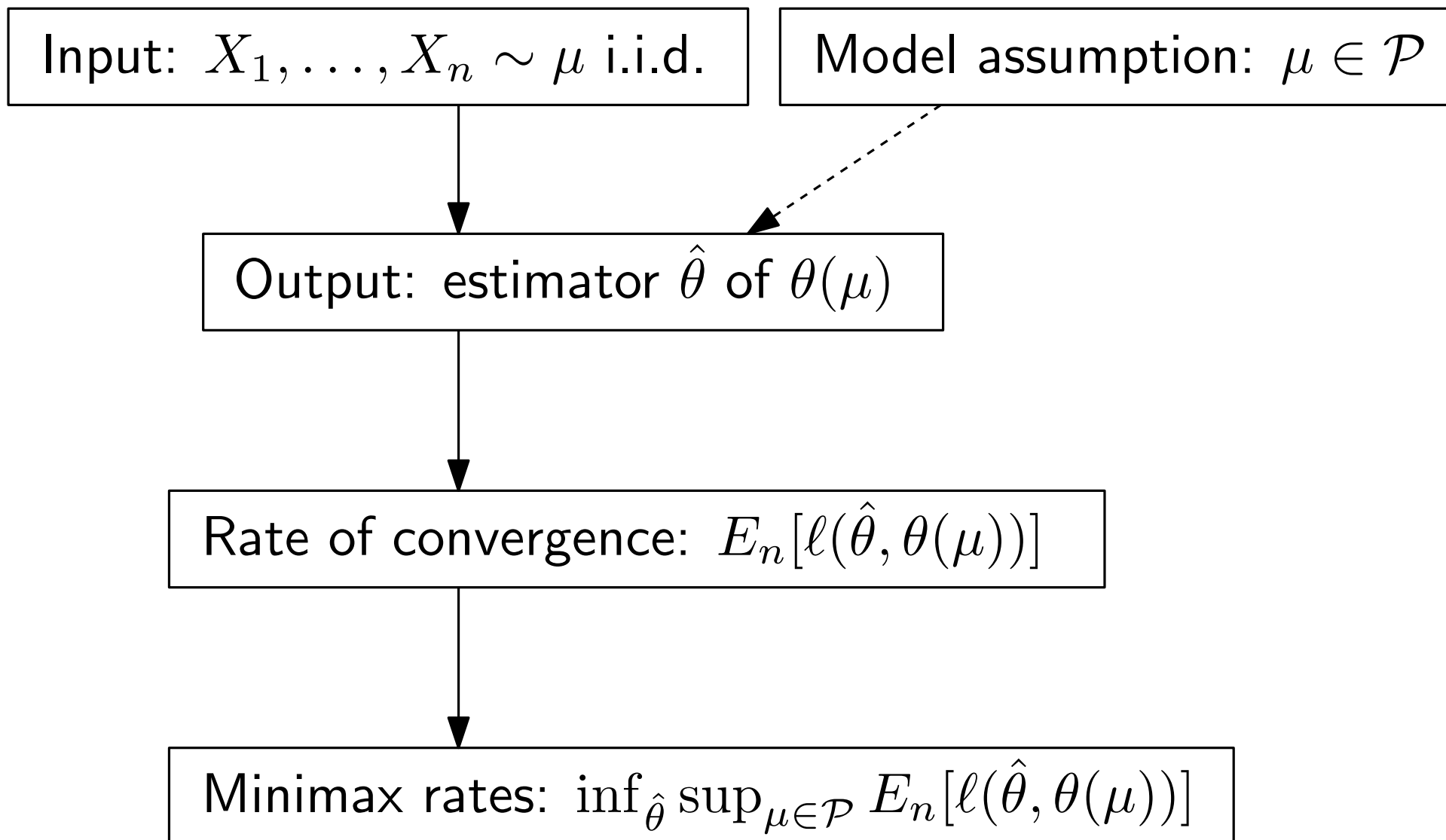
A (very) concise summary:



# Density estimation: the basics

---

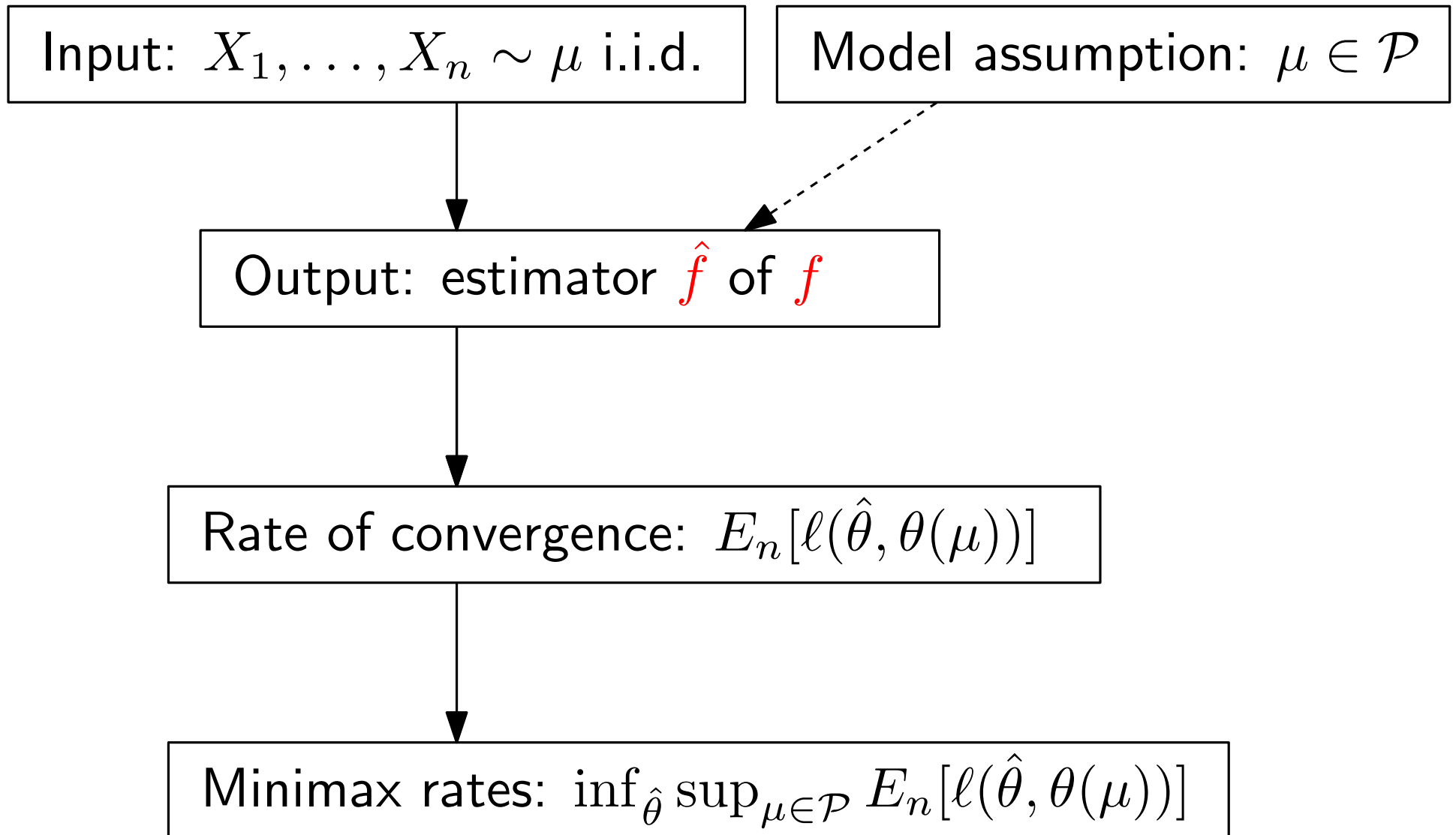
A (very) concise summary:



# Density estimation: the basics

---

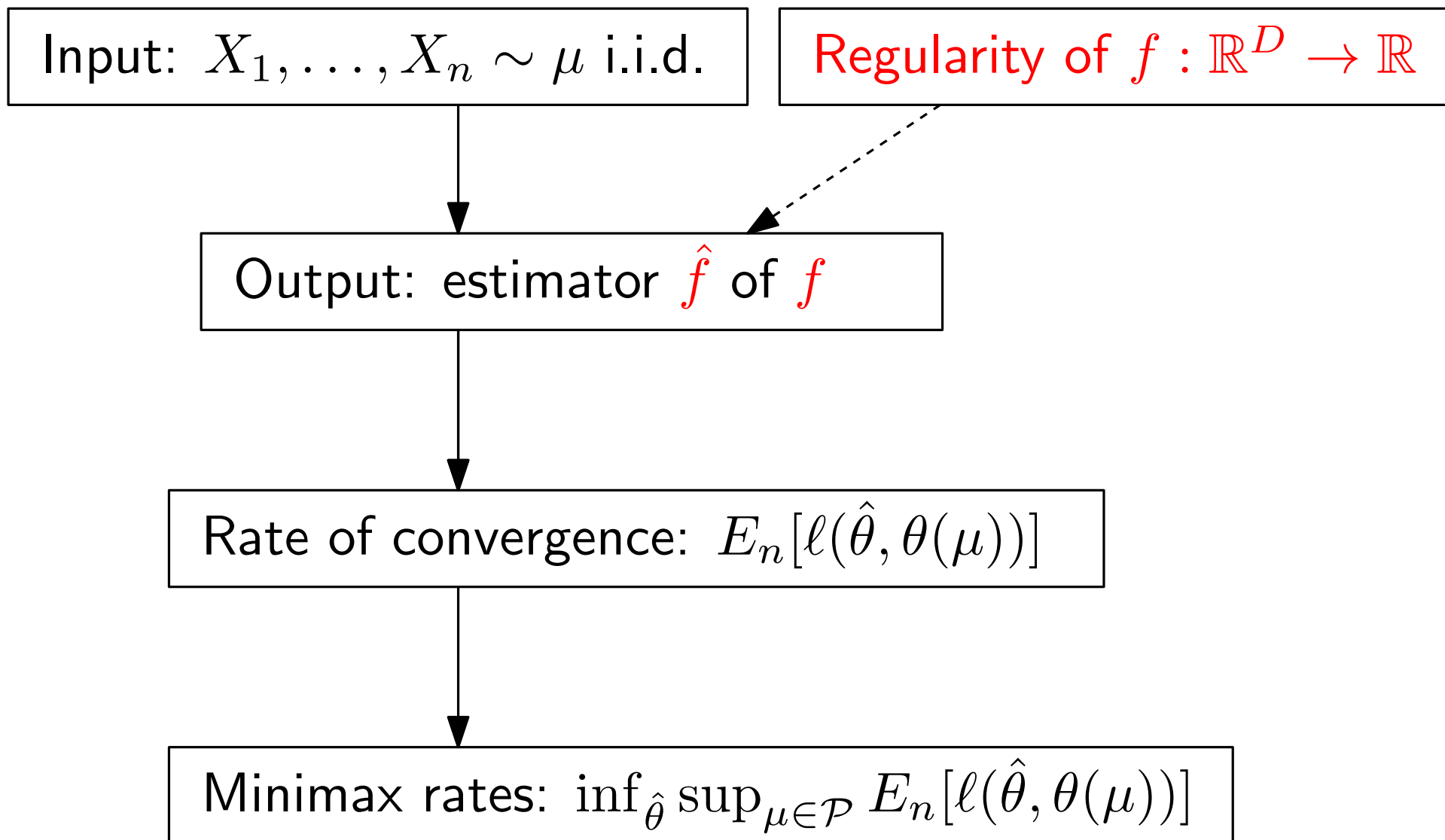
A (very) concise summary:



# Density estimation: the basics

---

A (very) concise summary:



# Density estimation: the basics

---

## Basic example

- $f \in H_2^s(\mathbb{R}^D) = \{f, \int |\mathcal{F}(f)(\xi)|^2 (1 + |\xi|^2)^s d\xi < \infty\}$
- Kernel density estimator:  $K$  with Fourier transform  $\kappa$

$$\hat{f}_h(x) = K_h * \mu_n(x) = \frac{1}{n} \sum_{i=1}^n h^{-D} K\left(\frac{x - X_i}{h}\right)$$



# Density estimation: the basics

---

## Basic example

- $f \in H_2^s(\mathbb{R}^D) = \{f, \int |\mathcal{F}(f)(\xi)|^2 (1 + |\xi|^2)^s d\xi < \infty\}$
- Kernel density estimator:  $K$  with Fourier transform  $\kappa$

$$\hat{f}_h(x) = K_h * \mu_n(x) = \frac{1}{n} \sum_{i=1}^n h^{-D} K\left(\frac{x - X_i}{h}\right)$$

- Expected value:  $f_h = E[\hat{f}_h] = K_h * f$

# Density estimation: the basics

---

## Basic example

- $f \in H_2^s(\mathbb{R}^D) = \{f, \int |\mathcal{F}(f)(\xi)|^2 (1 + |\xi|^2)^s d\xi < \infty\}$
- Kernel density estimator:  $K$  with Fourier transform  $\kappa$

$$\hat{f}_h(x) = K_h * \mu_n(x) = \frac{1}{n} \sum_{i=1}^n h^{-D} K\left(\frac{x - X_i}{h}\right)$$

- Expected value:  $f_h = E[\hat{f}_h] = K_h * f$
- $f(x) - \hat{f}_h(x) = f(x) - f_h(x) + f_h(x) - \hat{f}_h(x)$

# Density estimation: the basics

---

## Basic example

- $f \in H_2^s(\mathbb{R}^D) = \{f, \int |\mathcal{F}(f)(\xi)|^2 (1 + |\xi|^2)^s d\xi < \infty\}$
- Kernel density estimator:  $K$  with Fourier transform  $\kappa$

$$\hat{f}_h(x) = K_h * \mu_n(x) = \frac{1}{n} \sum_{i=1}^n h^{-D} K\left(\frac{x - X_i}{h}\right)$$

- Expected value:  $f_h = E[\hat{f}_h] = K_h * f$
- $f(x) - \hat{f}_h(x) = \underbrace{f(x) - f_h(x)}_{\text{biais}} + \underbrace{f_h(x) - \hat{f}_h(x)}_{\text{variance}}$

# Density estimation: the basics

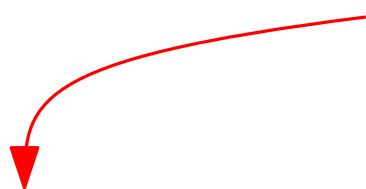
---

## Basic example

- $f \in H_2^s(\mathbb{R}^D) = \{f, \int |\mathcal{F}(f)(\xi)|^2 (1 + |\xi|^2)^s d\xi < \infty\}$
- Kernel density estimator:  $K$  with Fourier transform  $\kappa$

$$\hat{f}_h(x) = K_h * \mu_n(x) = \frac{1}{n} \sum_{i=1}^n h^{-D} K\left(\frac{x - X_i}{h}\right)$$

- Expected value:  $f_h = E[\hat{f}_h] = K_h * f$
- $f(x) - \hat{f}_h(x) = \underbrace{f(x) - f_h(x)}_{\text{biais}} + \underbrace{f_h(x) - \hat{f}_h(x)}_{\text{variance}}$



$$\begin{aligned} \|f - f_h\|_{L_2} &= \|\mathcal{F}(f)(\kappa(h\xi) - 1)\|_{L_2} \\ &\leq C_0 \|f\|_{H_2^s} h^s \end{aligned}$$

# Density estimation: the basics

---

## Basic example

- $f \in H_2^s(\mathbb{R}^D) = \{f, \int |\mathcal{F}(f)(\xi)|^2 (1 + |\xi|^2)^s d\xi < \infty\}$
- Kernel density estimator:  $K$  with Fourier transform  $\kappa$

$$\hat{f}_h(x) = K_h * \mu_n(x) = \frac{1}{n} \sum_{i=1}^n h^{-D} K\left(\frac{x - X_i}{h}\right)$$

- Expected value:  $f_h = E[\hat{f}_h] = K_h * f$

- $f(x) - \hat{f}_h(x) = \underbrace{f(x) - f_h(x)}_{\text{biais}} + \underbrace{f_h(x) - \hat{f}_h(x)}_{\text{variance}}$

$$\begin{aligned} \|f - f_h\|_{L_2} &= \|\mathcal{F}(f)(\kappa(h\xi) - 1)\|_{L_2} \\ &\leq C_0 \|f\|_{H_2^s} h^s \end{aligned}$$

$$E\|f_h - \hat{f}_h\|_{L_2} \leq \frac{C_1}{\sqrt{nh^D}}$$

# Density estimation: the basics

---

## Basic example

- $f \in H_2^s(\mathbb{R}^D) = \{f, \int |\mathcal{F}(f)(\xi)|^2 (1 + |\xi|^2)^s d\xi < \infty\}$
- Kernel density estimator:  $K$  with Fourier transform  $\kappa$

$$\hat{f}_h(x) = K_h * \mu_n(x) = \frac{1}{n} \sum_{i=1}^n h^{-D} K\left(\frac{x - X_i}{h}\right)$$

- Expected value:  $f_h = E[\hat{f}_h] = K_h * f$

- $f(x) - \hat{f}_h(x) = \underbrace{f(x) - f_h(x)}_{\text{biais}} + \underbrace{f_h(x) - \hat{f}_h(x)}_{\text{variance}}$

$$\begin{aligned} \|f - f_h\|_{L_2} &= \|\mathcal{F}(f)(\kappa(h\xi) - 1)\|_{L_2} \\ &\leq C_0 \|f\|_{H_2^s} h^s \end{aligned}$$

$$E\|f_h - \hat{f}_h\|_{L_2} \leq \frac{C_1}{\sqrt{nh^D}}$$

- Choose  $h \sim n^{-1/(2s+D)} \rightarrow \text{Risk} \sim n^{-s/(2s+D)}$ .

# Density estimation: the basics

---

## Basic example

- $f \in H_2^s(\mathbb{R}^D) = \{f, \int |\mathcal{F}(f)(\xi)|^2 (1 + |\xi|^2)^s d\xi < \infty\}$
- Kernel density estimator:  $K$  with Fourier transform  $\kappa$

$$\hat{f}_h(x) = K_h * \mu_n(x) = \frac{1}{n} \sum_{i=1}^n h^{-D} K\left(\frac{x - X_i}{h}\right)$$

- Expected value:  $f_h = E[\hat{f}_h] = K_h * f$

- $f(x) - \hat{f}_h(x) = \underbrace{f(x) - f_h(x)}_{\text{biais}} + \underbrace{f_h(x) - \hat{f}_h(x)}_{\text{variance}}$

$$\begin{aligned} \|f - f_h\|_{L_2} &= \|\mathcal{F}(f)(\kappa(h\xi) - 1)\|_{L_2} \\ &\leq C_0 \|f\|_{H_2^s} h^s \end{aligned}$$

$$E\|f_h - \hat{f}_h\|_{L_2} \leq \frac{C_1}{\sqrt{nh^D}}$$

- Choose  $h \sim n^{-1/(2s+D)} \rightarrow \text{Risk} \sim n^{-s/(2s+D)}$ .

# Curse of dimensionality (and what to do about it)

---

## Structural assumption on the signal

- Multi-index model, sparsity, small dimensional support, etc.

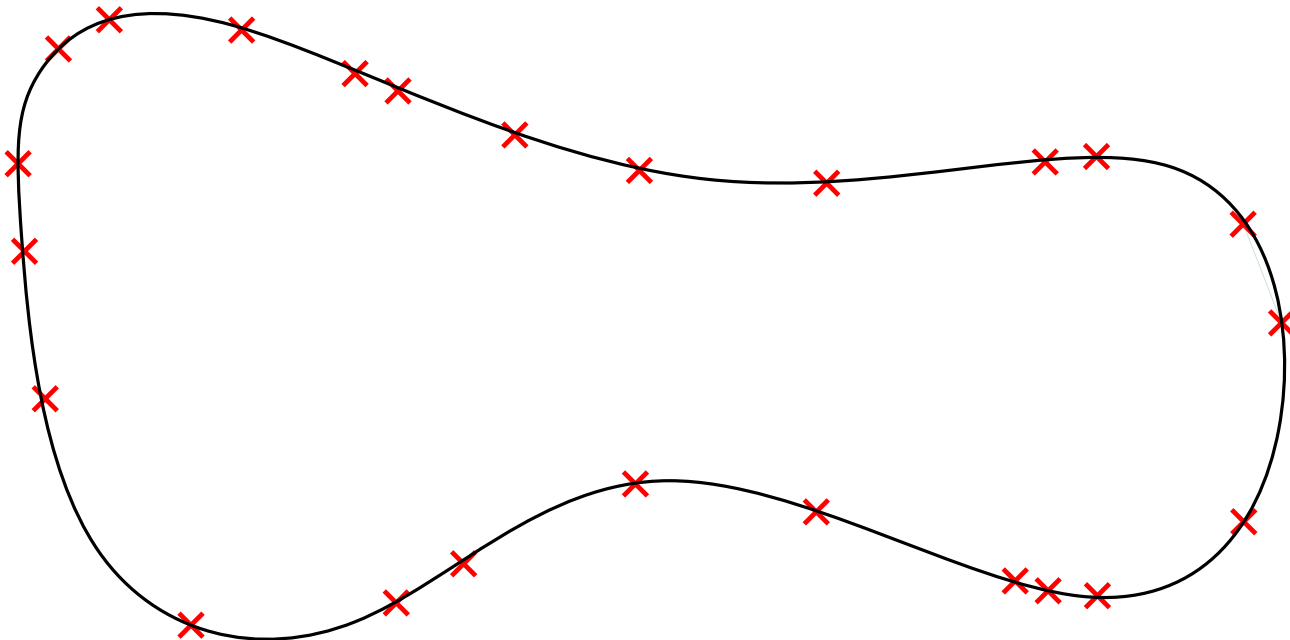


# Curse of dimensionality (and what to do about it)

---

## Structural assumption on the signal

- Multi-index model, sparsity, **small dimensional support**, etc.  
→ Manifold assumption:  $\mu$  on  $\mathbb{R}^D$  is supported on some **unknown** manifold  $M$  of dimension  $d \ll D$ .

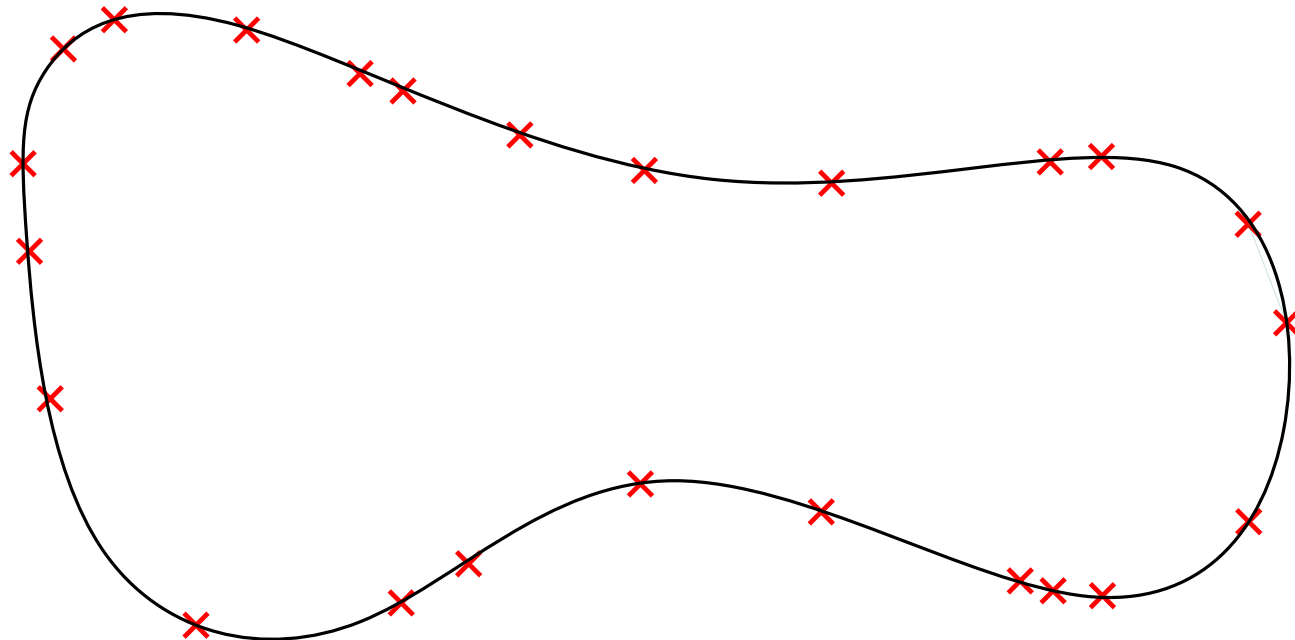


# Curse of dimensionality (and what to do about it)

---

## Structural assumption on the signal

- Multi-index model, sparsity, **small dimensional support**, etc.  
→ Manifold assumption:  $\mu$  on  $\mathbb{R}^D$  is supported on some **unknown** manifold  $M$  of dimension  $d \ll D$ .

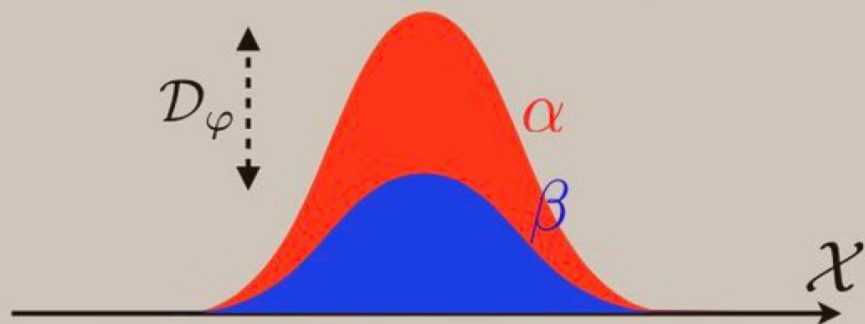


Loss?

# Distances between probability measures

Csiszár divergences:

$$\mathcal{D}_\varphi(\alpha|\beta) \stackrel{\text{def.}}{=} \int_{\mathcal{X}} \varphi\left(\frac{d\alpha}{d\beta}\right) d\beta$$



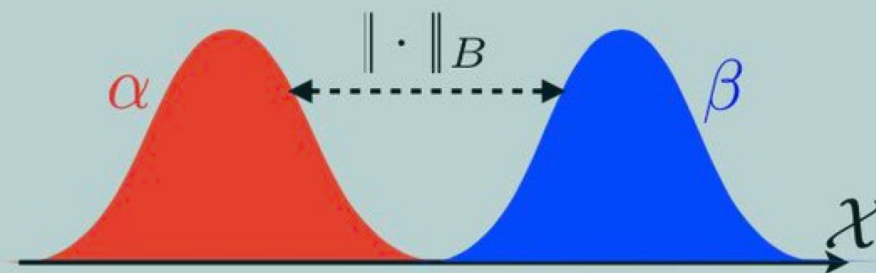
*Strong topology*

$\varphi$  convex,  $\varphi(1) = 0$ .

$$\varphi(r) = \begin{cases} r \log(r) & \rightarrow \text{KL} \\ |r - 1| & \rightarrow \text{TV} \\ |\sqrt{r} - 1|^2 & \rightarrow \text{Hellinger} \\ \dots & \end{cases}$$

Dual norms:

$$\|\alpha - \beta\|_B \stackrel{\text{def.}}{=} \max_{f \in B} \int_{\mathcal{X}} f(x)(d\alpha(x) - d\beta(x))$$



*Weak topology*

$$B = \{f ; \|f\|_B \leq 1\}$$

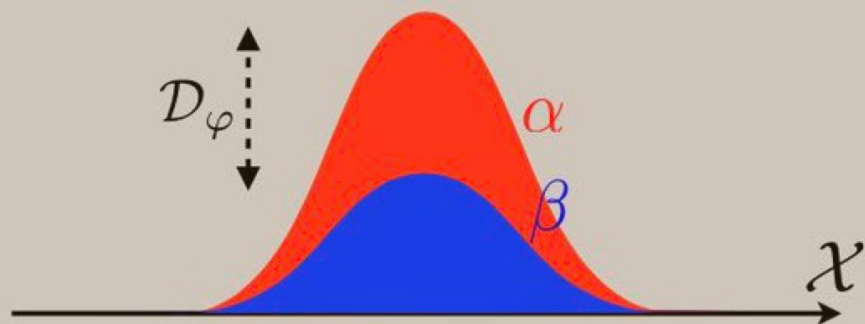
$$\|f\|_B = \begin{cases} \|\nabla f\|_\infty & \rightarrow W_1 \\ \|f\|_\infty + \|\nabla f\|_\infty & \rightarrow \text{flat} \\ \|\nabla^k f\|_2 & \rightarrow \text{MMD} \\ \dots & \end{cases}$$

[Source: Gabriel Peyré]

# Distances between probability measures

Csiszár divergences:

$$\mathcal{D}_\varphi(\alpha|\beta) \stackrel{\text{def.}}{=} \int_{\mathcal{X}} \varphi\left(\frac{d\alpha}{d\beta}\right) d\beta$$



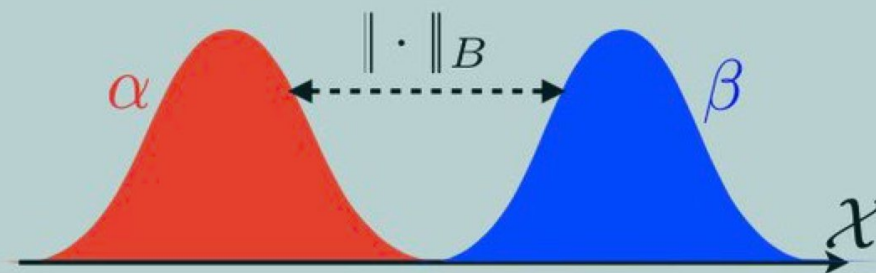
*Strong topology*

$\varphi$  convex,  $\varphi(1) = 0$ .

$$\varphi(r) = \begin{cases} r \log(r) & \rightarrow \text{KL} \\ |r - 1| & \rightarrow \text{TV} \\ |\sqrt{r} - 1|^2 & \rightarrow \text{Hellinger} \\ \dots & \end{cases}$$

Dual norms:

$$\|\alpha - \beta\|_B \stackrel{\text{def.}}{=} \max_{f \in B} \int_{\mathcal{X}} f(x)(d\alpha(x) - d\beta(x))$$



*Weak topology*

$$B = \{f ; \|f\|_B \leq 1\}$$

$$\|f\|_B = \begin{cases} \|\nabla f\|_\infty & \rightarrow W_1 \\ \|f\|_\infty + \|\nabla f\|_\infty & \rightarrow \text{flat} \\ \|\nabla^k f\|_2 & \rightarrow \text{MMD} \\ \dots & \end{cases}$$

[Source: Gabriel Peyré]

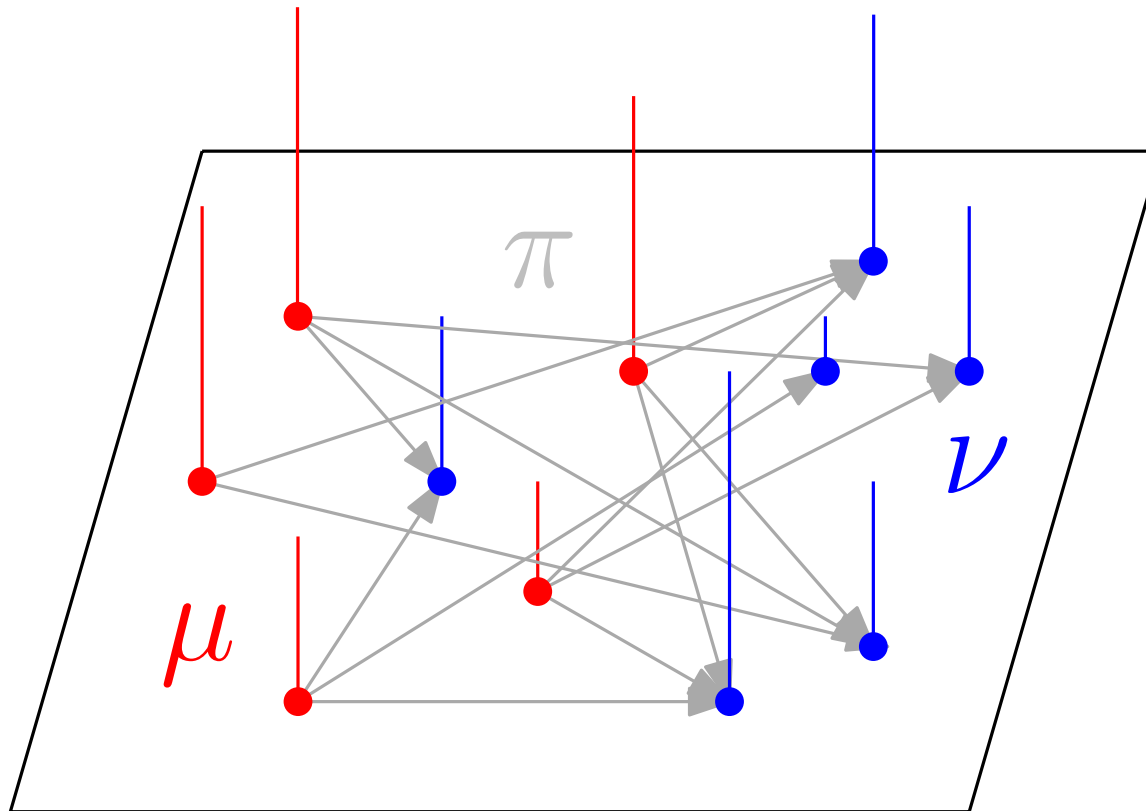
# Distances between probability measures

---

## Wasserstein distances:

- Let  $\mu, \nu$  be probability measures on  $\mathbb{R}^D$ .
- Let  $\Pi(\mu, \nu)$  be the set of couplings between  $\mu$  and  $\nu$ .

$$W_2(\mu, \nu) = \inf_{\pi \in \Pi(\mu, \nu)} \left( \int_{\mathbb{R}^D \times \mathbb{R}^D} |x - y|^2 d\pi(x, y) \right)^{1/2}$$



# Distances between probability measures

---

## Wasserstein distances:

- Let  $\mu, \nu$  be probability measures on  $\mathbb{R}^D$ .
- Let  $\Pi(\mu, \nu)$  be the set of couplings between  $\mu$  and  $\nu$ .

$$W_2(\mu, \nu) = \inf_{\pi \in \Pi(\mu, \nu)} \left( \int_{\mathbb{R}^D \times \mathbb{R}^D} |x - y|^2 d\pi(x, y) \right)^{1/2}$$

## Key observation:

If  $\mu, \nu$  have densities  $f, g$  on a manifold  $M$  with

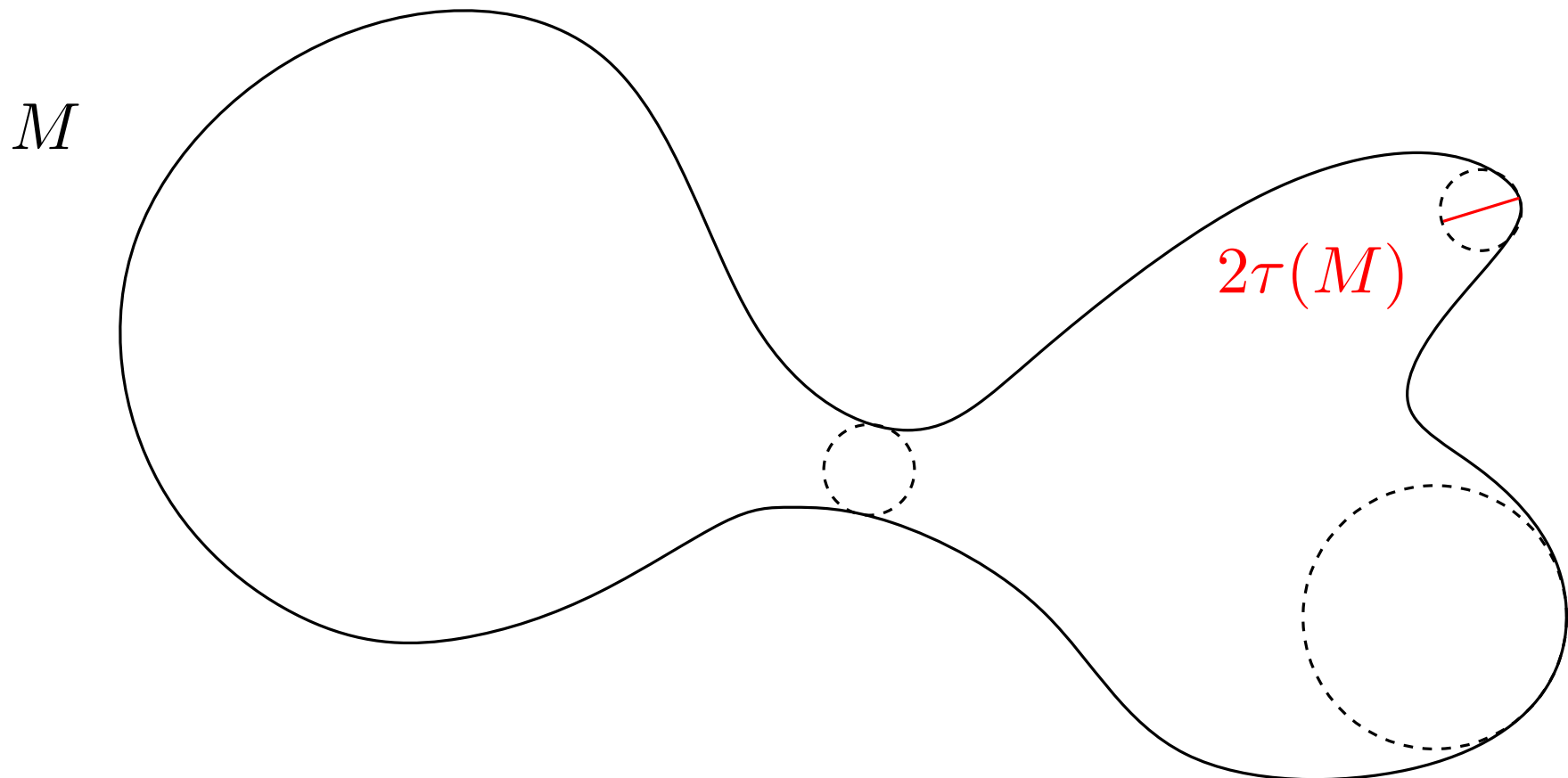
$$f, g \geq a, \text{ then } W_2(\mu, \nu) \leq a^{-1/2} \|f - g\|_{H_2^{-1}(M)}$$

- In  $1D$ : distance between c.d.f.
- For  $p = 1$ ,  $W_1(\mu, \nu) = \sup\{(\mu - \nu)(f), f \text{ 1-Lip}\}$ .

# Model assumptions

---

- $M \subset \mathbb{R}^D$  is a  $d$ -dimensional manifold with reach  $\tau(M) \geq \tau_{\min}$ .
- $\mu$  has a density  $f$  on  $M$  in  $H_2^s(M)$  with  $f \geq a > 0$ .



# Model assumptions

---

- $M \subset \mathbb{R}^D$  is a  $d$ -dimensional with reach  $\tau(M) \geq \tau_{\min}$ .
- $\mu$  has a density  $f$  on  $M$  in  $H_2^s(M)$  with  $f \geq a > 0$ .

## Goal:

Given  $X_1, \dots, X_n \sim \mu$ , produce an estimator  $\hat{\mu}$  of  $\mu$  with  $E_n W_2(\hat{\mu}, \mu)$  small.



# Model assumptions

---

- $M \subset \mathbb{R}^D$  is a  $d$ -dimensional with reach  $\tau(M) \geq \tau_{\min}$ .
- $\mu$  has a density  $f$  on  $M$  in  $H_2^s(M)$  with  $f \geq a > 0$ .

## Goal:

Given  $X_1, \dots, X_n \sim \mu$ , produce an estimator  $\hat{\mu}$  of  $\mu$  with  $E_n W_2(\hat{\mu}, \mu)$  small.

- Kernel estimators still work! Let  $\hat{\mu}_h$  have a density  $\hat{f}_h$  on  $M$ , given by

$$\hat{f}_h(x) = \frac{1}{n} \sum_{i=1}^n h^{-d} K\left(\frac{x - X_i}{h}\right) \quad \text{for } x \in M$$

# Model assumptions

---

Bias

$$\|f - f_h\|_{H_2^{-1}(M)} \leq Ch^{s+1}$$

Variance

$$E\|\hat{f}_h - f_h\|_{H_2^{-1}(M)} \leq \frac{C_1 h}{\sqrt{nh^d}}$$

- Choose  $h \sim n^{-1/(2s+d)}$

$$E_n W_2(\hat{\mu}_h, \mu) \lesssim n^{-(s+1)/(2s+d)}$$

# Model assumptions

---

Bias

$$\|f - f_h\|_{H_2^{-1}(M)} \leq Ch^{s+1}$$

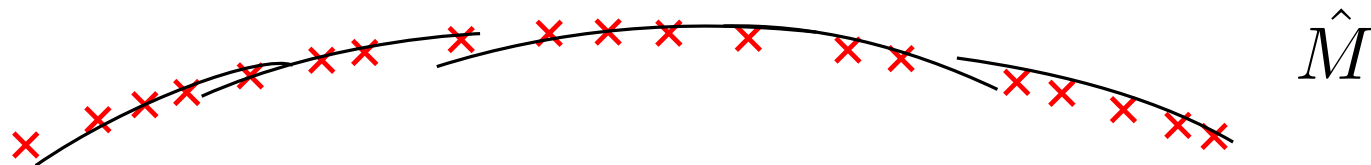
Variance

$$E\|\hat{f}_h - f_h\|_{H_2^{-1}(M)} \leq \frac{C_1 h}{\sqrt{nh^d}}$$

- Choose  $h \sim n^{-1/(2s+d)}$

$$E_n W_2(\hat{\mu}_h, \mu) \lesssim n^{-(s+1)/(2s+d)}$$

- $\hat{\mu}_h$  is supported on  $M$ ... which is unknown!  
→ Use a manifold estimator  $\hat{M}$  in a preprocessing step



# Take-home message

---

- To avoid the curse of dimensionality, a structural assumption has to be made on the signal to be estimate.
- Signal on  $d$ -manifold  $M \subset \mathbb{R}^D \rightarrow$  rates of estimation depending on  $d$  and not  $D$ .
  - $\rightarrow$  via kernel estimation on an estimated manifold