

# Improved learning theory for kernel distribution regression with two-stage sampling

**François Bachoc**

Institut de Mathématiques de Toulouse  
Université Paul Sabatier

Joint work with **Louis Béthune** (Toulouse), **Alberto González-Sanz** (Columbia) and **Jean-Michel Loubes** (Toulouse)

Journée de lancement  
ANR BACKUP  
December 2023

- 1 Distribution regression, Hilbertian embedding and two-stage sampling
- 2 Near-unbiased condition and improved rates

# Distribution regression

We observe i.i.d. pairs

$$(\mu_i, Y_i), \quad i = 1, \dots, n.$$

- $Y_i \in \mathbb{R}$ .
- $\mu_i$  is a probability distribution on  $\Omega$ .
- $\Omega$  is compact in  $\mathbb{R}^d$ .

**Goal:** constructing a regression function

$$\hat{f}_n : \mathcal{P}(\Omega) \rightarrow \mathbb{R},$$

- where  $\mathcal{P}(\Omega)$  is the set of probability distributions on  $\Omega$ .

**Application fields** described in [Szabó et al., 2015, Szabó et al., 2016, Meunier et al., 2022, Bachoc et al., 2023a].

## Hilbertian embedding

$$\begin{aligned}x &: \mathcal{P}(\Omega) \rightarrow \mathcal{H} \\ \mu &\mapsto x_\mu,\end{aligned}$$

where  $\mathcal{H}$  is a Hilbert space.

$\Rightarrow$  In order to use kernels on Hilbert spaces (see later)!

# Hilbertian embedding 1: mean embedding

Consider a kernel  $k$  on  $\Omega$ .

## Very quick introduction to kernels and RKHS

- $k : \Omega \times \Omega \rightarrow \mathbb{R}$ .
- For any  $\ell \in \mathbb{N}$ ,  $t_1, \dots, t_\ell \in \Omega$ , the  $\ell \times \ell$  matrix  $[k(t_i, t_j)]$  is **symmetric non-negative definite**.
- There is a (unique) Hilbert space  $\mathcal{H}_k$  of functions from  $\Omega$  to  $\mathbb{R}$ ,
  - with inner product  $\langle \cdot, \cdot \rangle_{\mathcal{H}_k}$
  - with norm  $\| \cdot \|_{\mathcal{H}_k}$such that
  - $\mathcal{H}_k$  contains all functions  $k_t := k(t, \cdot)$  for  $t \in \Omega$ ,
  - for all  $g \in \mathcal{H}_k$ , for all  $t \in \Omega$ ,  $g(t) = \langle g, k_t \rangle_{\mathcal{H}_k}$  **reproducing property**. $\implies \mathcal{H}_k$  is the **reproducing kernel Hilbert space (RKHS)** of  $k$ .

Then **mean embedding**

$$x_\mu := \left( t \mapsto \int_{\Omega} k(t, x) d\mu(x) \right) = \int_{\Omega} k_x d\mu(x),$$

[Szabó et al., 2015, Szabó et al., 2016, Muandet et al., 2017].

# Hilbertian embedding 2: sliced Wasserstein

The **sliced Wasserstein distance** [Kolouri et al., 2018, Manole et al., 2022, Meunier et al., 2022]

$$SW(\mu, \nu)^2 := \int_{\mathcal{S}^{d-1}} \int_0^1 (F_{\mu_\theta}^{-1}(t) - F_{\nu_\theta}^{-1}(t))^2 dt d\Lambda(\theta),$$

- with two distributions  $\mu, \nu \in \mathcal{P}(\Omega)$ ,
- where  $\mathcal{S}^{d-1}$  is the unit sphere ,
- where  $\Lambda$  is the uniform distribution on  $\mathcal{S}^{d-1}$ ,
- where  $\mu_\theta$  is the **univariate** distribution of  $\langle \theta, X \rangle$  for  $X \sim \mu$ ,
- where  $F_{\mu_\theta}^{-1}$  is the quantile function of  $\mu_\theta$ .

Hilbert distance of a **Hilbertian embedding**

- $\mathcal{H} = \mathcal{S}^{d-1} \times [0, 1]$ ,
- $x_\mu(\theta, t) = F_{\mu_\theta}^{-1}(t)$ .

## Hilbertian embedding 3: Sinkhorn distance and dual potential

Dual formulation of entropic-regularized (Sinkhorn) optimal transport  
[Genevay, 2019]

$$\sup_{h \in L^1(\mu), g \in L^1(\mathcal{U})} \int_{\Omega} h(x) d\mu(x) + \int_{\Omega} g(y) d\mathcal{U}(y) - \epsilon \int_{\Omega \times \Omega} e^{\frac{1}{\epsilon}(h(x) + g(y) - \frac{1}{2}\|x-y\|^2)} d\mu(x) d\mathcal{U}(y).$$

- $\epsilon > 0$  regularization parameter.
- Fixed  $\mathcal{U} \in \mathcal{P}(\Omega)$  called reference measure.
- For any  $\mu \in \mathcal{P}(\Omega)$ .

### Hilbertian embedding

There is a unique optimal  $(h^*, g^*)$  such that  $g^*$  is centered w. r. t.  $\mathcal{U}$ .  
Also  $g^* \in L^2(\mathcal{U})$ .

[Bachoc et al., 2023a]:

- $x_{\mu} := g^*$ .
- $\mathcal{H} := L^2(\mathcal{U})$ .

# Kernel ridge regression on Hilbert space

- **Hilbertian covariates:** for  $i = 1, \dots, n$ , let

$$x_i := x_{\mu_i}.$$

- **Squared exponential kernel on  $\mathcal{H}$ :** for  $u, v \in \mathcal{H}$ ,

$$K(u, v) := e^{-\|u-v\|_{\mathcal{H}}^2}.$$

$\implies$  Yields the RKHS  $\mathcal{H}_K$  of functions from  $\mathcal{H}$  to  $\mathbb{R}$ .

- **Ridge regression**

$$\hat{f}_n = \operatorname{argmax}_{f \in \mathcal{H}_K} R_n(f)$$

with

$$R_n(f) := \frac{1}{n} \sum_{i=1}^n (f(x_i) - Y_i)^2 + \lambda \|f\|_{\mathcal{H}_K}^2,$$

- where  $\lambda > 0$  is a **regularization parameter**.

# Two-stage sampling

Studied in [Szabó et al., 2015, Szabó et al., 2016, Meunier et al., 2022].

- For  $i = 1, \dots, n$ ,  $\mu_i$  is **unobserved**.
- We observe i. i. d.  $(X_{i,j})_{j=1, \dots, N}$  with  $X_{i,j} \sim \mu_i$ .
- We let

$$\mu_i^N = \frac{1}{N} \sum_{j=1}^N \delta_{X_{i,j}}$$

and

$$x_{N,i} = x_{\mu_i^N}.$$

## Ridge regression with approximate covariates

$$\hat{f}_{n,N} = \operatorname{argmax}_{f \in \mathcal{H}_K} R_{n,N}(f)$$

with

$$R_{n,N}(f) := \frac{1}{n} \sum_{i=1}^n (f(x_{N,i}) - Y_i)^2 + \lambda \|f\|_{\mathcal{H}_K}^2.$$

1 Distribution regression, Hilbertian embedding and two-stage sampling

2 Near-unbiased condition and improved rates

# Existing error bounds on $\hat{f}_n - \hat{f}_{n,N}$

- [Szabó et al., 2015, Szabó et al., 2016, Meunier et al., 2022] address their respective distribution regression settings.
- But their results are naturally made general.

## Existing bounds

For all  $s \geq 1$ , conditionally to  $(x_i, Y_i)_{i=1}^n$ ,

$$\mathbb{E} \left[ \left\| \hat{f}_n - \hat{f}_{n,N} \right\|_{\mathcal{H}_K}^s \right]^{1/s} \leq \frac{\text{constant} \left( \left\| \hat{f}_n \right\|_{\mathcal{H}_K} + Y_{\max,n} \right)}{\sqrt{N\lambda}}$$

- with  $Y_{\max,n} = \max_{i=1,\dots,n} |Y_i|$ .

# Existing proofs are improvable?

- Proofs based on explicit expressions of  $\hat{f}_n$  and  $\hat{f}_{n,N}$ .
- Somewhere:

$$\begin{aligned} & \left\| \frac{1}{n} \sum_{i=1}^n \left( \hat{f}_n(x_i) K_{x_i} - \hat{f}_n(x_{N,i}) K_{x_{N,i}} \right) \right\|_{\mathcal{H}_K} \\ & \leq \frac{1}{n} \sum_{i=1}^n \left\| \hat{f}_n(x_i) K_{x_i} - \hat{f}_n(x_{N,i}) K_{x_{N,i}} \right\|_{\mathcal{H}_K}. \end{aligned}$$

- But  $\hat{f}_n(x_i) K_{x_i} - \hat{f}_n(x_{N,i}) K_{x_{N,i}}$  are **independent** conditionally on  $(x_i, Y_i)_{i=1}^n$ .
- Do they have **approximately zero mean**?

# Near-unbiased condition

In [\[Bachoc et al., 2023b\]](#).

## Near-unbiased condition

- For  $i = 1, \dots, n$ , there are random  $a_{N,i}$  and  $b_{N,i}$  such that

$$x_{N,i} - x_i = a_{N,i} + b_{N,i}.$$

- $\|a_{N,i}\|_{\mathcal{H}}$  has order  $\frac{1}{\sqrt{N}}$ .
- $\mathbb{E}(a_{N,i}|\mu_i) = 0 \in \mathcal{H}$ .
- $\|b_{N,i}\|_{\mathcal{H}}$  has order  $\frac{1}{N}$ .

## For the 3 examples of Hilbertian embedding

- **Mean embedding:**  $b_{N,i} = 0$  (exactly unbiased).
- **Sinkhorn:** indeed near unbiased, relying on [\[González-Sanz et al., 2022\]](#).
- **Sliced Wasserstein:** indeed near unbiased under conditions.

In [\[Bachoc et al., 2023b\]](#).

## Theorem

Up to constant

$$\sqrt{\mathbb{E}_n \left[ \|\hat{f}_n - \hat{f}_{n,N}\|_{\mathcal{H}_K}^2 \right]} \leq \frac{Y_{\max,n} + \|\hat{f}_n\|_{\mathcal{H}_K}}{\lambda N} + \frac{Y_{\max,n} + \|\hat{f}_n\|_{\mathcal{H}_K}}{\lambda \sqrt{n} \sqrt{N}} \\ + \left( 1 + \frac{\sqrt{N}}{\sqrt{n}} \right)^{-1} \left( \frac{Y_{\max,n} + \|\hat{f}_n\|_{\mathcal{H}_K}}{\lambda n} + \frac{Y_{\max,n} + \|\hat{f}_n\|_{\mathcal{H}_K}}{\lambda^2 n \sqrt{N}} \right)$$

- with  $Y_{\max,n} = \max_{i=1,\dots,n} |Y_i|$ ,
- where  $\mathbb{E}_n$  denotes the conditional expectation given  $(\mu_i, Y_i)_{i=1}^n$ .

The  $\sqrt{n}$  we gain comes from average of independent centered variables.

# Proof ingredient 1

$$\hat{f}_n = \operatorname{argmin}_{f \in \mathcal{H}_K} \frac{1}{n} \sum_{i=1}^n (f(x_i) - Y_i)^2 + \lambda \|f\|_{\mathcal{H}_K}^2$$

and

$$\hat{f}_{n,N} = \operatorname{argmin}_{f \in \mathcal{H}_K} \frac{1}{n} \sum_{i=1}^n (f(x_{N,i}) - Y_i)^2 + \lambda \|f\|_{\mathcal{H}_K}^2.$$

Then, **exploiting convexity**,

$$\begin{aligned} \lambda \|\hat{f}_n - \hat{f}_{n,N}\|_{\mathcal{H}_K}^2 &\leq \frac{1}{n} \sum_{i=1}^n \left\{ \left[ \hat{f}_n(x_{N,i}) - \hat{f}_{n,N}(x_{N,i}) \right] \hat{f}_n(x_{N,i}) \right. \\ &\quad \left. - \left[ \hat{f}_n(x_i) - \hat{f}_{n,N}(x_i) \right] \hat{f}_n(x_i) \right\} \\ &\quad + \frac{1}{n} \sum_{i=1}^n Y_i \left\{ \left[ \hat{f}_n(x_i) - \hat{f}_{n,N}(x_i) \right] - \left[ \hat{f}_n(x_{N,i}) - \hat{f}_{n,N}(x_{N,i}) \right] \right\}. \end{aligned}$$

## Proof ingredient 2

We are led to bound (in  $\mathbb{R}$ !) terms such as

$$\mathbb{E} \left( \frac{1}{n} \sum_{i=1}^n Y_i \left\{ \left[ \hat{f}_n(x_i) - \hat{f}_{n,N}(x_i) \right] - \left[ \hat{f}_n(x_{N,i}) - \hat{f}_{n,N}(x_{N,i}) \right] \right\} \right).$$

By coupling arguments, we approximate by

$$\mathbb{E} \left( \frac{1}{n} \sum_{i=1}^n Y_i \left\{ \left[ \hat{f}_n(x_i) - \tilde{f}_{n,N}(x_i) \right] - \left[ \hat{f}_n(x_{N,i}) - \tilde{f}_{n,N}(x_{N,i}) \right] \right\} \right),$$

- with  $\tilde{f}_{n,N}$  constructed from new independent  $(\tilde{x}_{N,i})_{i=1}^n$ .

# Application to sufficient $N$ for minimax rate (1/2)

- [Caponnetto and De Vito, 2007] provide minimax rates as  $n \rightarrow \infty$  with **one-stage sampling** (for  $\hat{f}_n$ ).
- **Target:** conditional expectation function

$$f^* = \mathbb{E}(Y_i | x_i = \cdot) \quad \text{assumed to be in } \mathcal{H}_K.$$

- We let  $\mathcal{L}$  be the distribution of  $x_i$ .

## Problem class on $\mathcal{H}$

Hardness of  $(\mathcal{L}, K, f^*)$  measured by

- $b > 1$  **effective dimension of  $\mathcal{H}_K$  w. r. t. distribution  $\mathcal{L}$ ,**
- $c \in (1, 2]$  **complexity of  $f^*$ .**

## Minimax rate

$$\sqrt{\int_{\mathcal{H}} \left( f^*(x) - \hat{f}_n(x) \right)^2 d\mathcal{L}(x)} = \mathcal{O}_{\mathbb{P}} \left( n^{-\frac{bc}{2(bc+1)}} \right).$$

- With  $\lambda = n^{-\frac{b}{bc+1}}$ .

# Application to sufficient $N$ for minimax rate (2/2)

In [Bachoc et al., 2023b], from our bounds:

## Sufficient $N$ for minimax

$$\sqrt{\int_{\mathcal{H}} \left( f^*(x) - \hat{f}_{n,N}(x) \right)^2 d\mathcal{L}(x)} = \mathcal{O}_{\mathbb{P}} \left( n^{-\frac{bc}{2(bc+1)}} \right).$$

- With  $\lambda = n^{-\frac{b}{bc+1}}$ .
- With  $N = n^a$ ,

$$\begin{cases} a = \max\left(\frac{b+\frac{bc}{2}}{bc+1}, \frac{2b-1}{bc+1}, \frac{4b-bc-2}{bc+1}\right) (\leq 1) & \text{if } b\left(1 - \frac{c}{2}\right) \leq \frac{3}{4} \\ a = \max\left(\frac{b+\frac{bc}{2}}{bc+1}, \frac{2b-\frac{1}{2}}{bc+1}\right) (> 1) & \text{if } b\left(1 - \frac{c}{2}\right) > \frac{3}{4} \end{cases}.$$

In [Szabó et al., 2015, Szabó et al., 2016], same result for mean embedding with  $N = n^{\frac{b(c+1)}{bc+1}}$ ,

- $\frac{b(c+1)}{bc+1} > a$ .

# Conclusion

- Hilbertian embedding for (symmetric non-negative definite) kernels.
- Two-stage sampling as an additional source of error.
- Main contribution: tighter control of this error.
- The paper [[Bachoc et al., 2023b](#)]: arXiv:2308.14335.
- Paper [[Bachoc et al., 2023a](#)] on Sinkhorn kernel.
- Public Python codes (links in papers).

**Thank you for your attention!**

# Bibliography I

-  Bachoc, F., Béthune, L., González-Sanz, A., and Loubes, J.-M. (2023a).  
Gaussian processes on distributions based on regularized optimal transport.  
*In International Conference on Artificial Intelligence and Statistics.*
-  Bachoc, F., Béthune, L., González-Sanz, A., and Loubes, J.-M. (2023b).  
Improved learning theory for kernel distribution regression with two-stage sampling.  
*arXiv:2308.14335.*
-  Caponnetto, A. and De Vito, E. (2007).  
Optimal rates for the regularized least-squares algorithm.  
*Foundations of Computational Mathematics*, 7:331–368.
-  Genevay, A. (2019).  
*Entropy-regularized optimal transport for machine learning.*  
PhD thesis, Paris Sciences et Lettres (ComUE).

-  González-Sanz, A., Loubes, J.-M., and Niles-Weed, J. (2022). Weak limits of entropy regularized optimal transport; potentials, plans and divergences. *arXiv preprint arXiv:2207.07427*.
-  Kolouri, S., Rohde, G. K., and Hoffmann, H. (2018). Sliced Wasserstein distance for learning Gaussian mixture models. In *IEEE Conference on Computer Vision and Pattern Recognition*.
-  Manole, T., Balakrishnan, S., and Wasserman, L. (2022). Minimax confidence intervals for the sliced Wasserstein distance. *Electronic Journal of Statistics*, 16(1):2252–2345.
-  Meunier, D., Pontil, M., and Ciliberto, C. (2022). Distribution regression with sliced wasserstein kernels. In *International Conference on Machine Learning*.

-  Muandet, K., Fukumizu, K., Sriperumbudur, B., and Schölkopf, B. (2017).  
Kernel mean embedding of distributions: A review and beyond.  
*Foundations and Trends® in Machine Learning*, 10(1-2):1–141.
-  Szabó, Z., Gretton, A., Póczos, B., and Sriperumbudur, B. (2015).  
Two-stage sampled learning theory on distributions.  
*In International Conference on Artificial Intelligence and Statistics*.
-  Szabó, Z., Sriperumbudur, B., Póczos, B., and Gretton, A. (2016).  
Learning theory for distribution regression.  
*The Journal of Machine Learning Research*, 17(1):5272–5311.