

Statistical and Geometrical properties of the regularized kernel Kullback Leibler



Clémentine Chazal¹

Anna Korba¹

Francis Bach²



CREST/ENSAE, IP Paris¹, INRIA, Paris²

Contributions

- Introduction of a regularized definition of the KKL form [1] which is defined for any probability distributions and study of its properties.
- Derivation of closed form expressions for the regularized KKL and its Wasserstein Gradient on empirical measures.
- Implementation of a sampling algorithm following a gradient descent scheme that obtains results on low-dimensional experiments.

Introduction and motivations

Problem: To approximate a target distribution q on \mathbb{R}^d , we solve the optimization problem

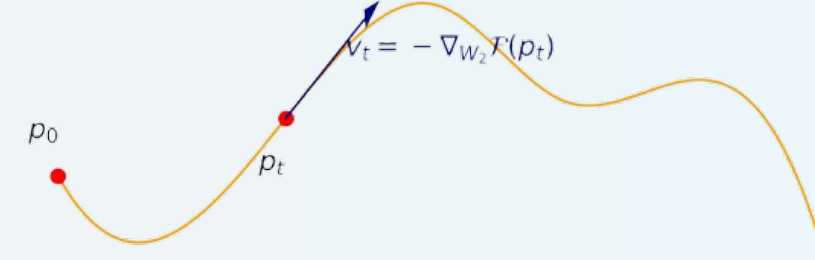
$$\min_{p \in \mathcal{P}(\mathbb{R}^d)} \mathcal{F}(p)$$

where $\mathcal{F}(p) = D(p||q)$ with D a **divergence** or a **distance**.

Wasserstein gradient flow:

- If for any function $h : \mathbb{R}^d \rightarrow \mathbb{R}^d$, $\varepsilon > 0$, the expansion
$$\mathcal{F}((I_d + \varepsilon h)_{\#} p) = \mathcal{F}(p) + \varepsilon \langle \nabla_{W_2} \mathcal{F}(p), h \rangle_p + o(\varepsilon),$$
holds, then $\nabla_{W_2} \mathcal{F}(p) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is the **Wasserstein gradient** of \mathcal{F} .

$$\begin{cases} \text{Wasserstein Gradient Flow} \\ p(0) = p_0, \\ \partial_t p(t) = -\nabla_{W_2} \mathcal{F}(p(t)). \end{cases}$$



The choice of D dictates the overall dynamics. In this project we selected the regularized Kernel Kullback Leibler Divergence.

Kernel Kullback Leibler divergence (KKL)

Kernel Kullback Leibler divergence (KKL): Given \mathcal{H} a RKHS with **reproducing kernel** k , for $p \ll q$, the KKL divergence is

$$\text{KKL}(p||q) := \text{Tr}[\Sigma_p(\log \Sigma_p - \log \Sigma_q)]$$

where

$$\Sigma_p = \int k(\cdot, x) k(\cdot, x)^* dp(x).$$

If k^2 is universal and $\forall x \in \mathbb{R}^d$, $k(x, x) = 1$ then

$$\text{KKL}(p||q) = 0 \Leftrightarrow p = q.$$

Regularized KKL: To handle cases where $p \not\ll q$, the regularized KKL is defined for $\alpha \in]0, 1[$ as

$$\text{KKL}_\alpha(p || q) := \text{KKL}(p || (1 - \alpha)q + \alpha p)$$

Closed form for regularized KKL on empirical distributions

Regularized KKL for empirical distributions: Let $x^1, \dots, x^n \sim p$, $y^1, \dots, y^m \sim q$ and note $\hat{p} = \frac{1}{n} \sum_{i=1}^n \delta_{x^i}$ and $\hat{q} = \frac{1}{m} \sum_{j=1}^m \delta_{y^j}$.

divergence

Regularized KKL admits a closed form expression

$$\text{KKL}_\alpha(\hat{p}||\hat{q}) = \text{Tr} \left(\frac{1}{n} K_{\hat{p}} \log \frac{1}{n} K_{\hat{p}} \right) - \text{Tr} (I_\alpha K \log(K)),$$

$$I_\alpha = \begin{pmatrix} \frac{1}{\alpha} I & 0 \\ 0 & 0 \end{pmatrix} \text{ and } K = \begin{pmatrix} \frac{\alpha}{n} K_{\hat{p}} & \sqrt{\frac{\alpha(1-\alpha)}{nm}} K_{\hat{p}, \hat{q}} \\ \sqrt{\frac{\alpha(1-\alpha)}{nm}} K_{\hat{q}, \hat{p}} & \frac{1-\alpha}{m} K_{\hat{q}} \end{pmatrix}$$

and $K_{\hat{p}} = (k(x^i, x^j))_{i,j=1}^n$, $K_{\hat{q}} = (k(y^i, y^j))_{i,j=1}^m$, $K_{\hat{p}, \hat{q}} = (k(x^i, y^j))_{i,j=1}^{n,m}$.

Wasserstein gradient for empirical measures:

$$\nabla_{W_2} \mathcal{F}(\hat{p})(x) = \nabla_x (S(x)^T g(K_{\hat{p}}) S(x) - T(x)^T g(K) T(x) - T(x)^T A T(x))$$

where $S(x) = (\frac{1}{\sqrt{n}} k(x, x^i))_i$, $T(x) = ((\frac{\alpha}{\sqrt{n}} k(x, x^i))_i, (\sqrt{\frac{1-\alpha}{m}} k(x, y^j))_j)$ and A is a matrix depending on the eigenvalues and eigenvectors of K .

Theoretical properties of the regularized KKL

- The regularized KKL is constant to the true KKL for $p \ll q$ when $\alpha \rightarrow 0$:

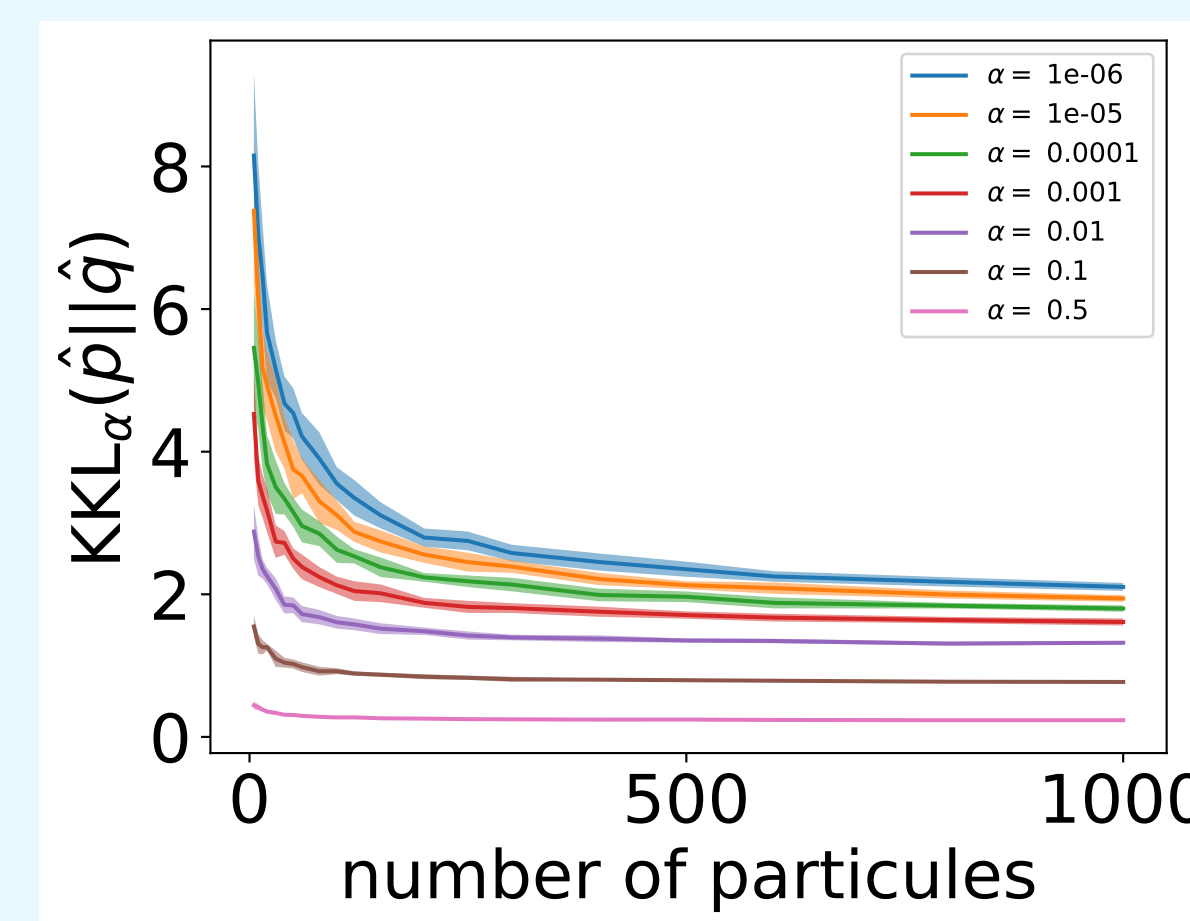
$$\text{KKL}_\alpha(p||q) \xrightarrow{\alpha \rightarrow 0} \text{KKL}(p||q).$$

- $\alpha \rightarrow \text{KKL}_\alpha(p||q)$ is decreasing.

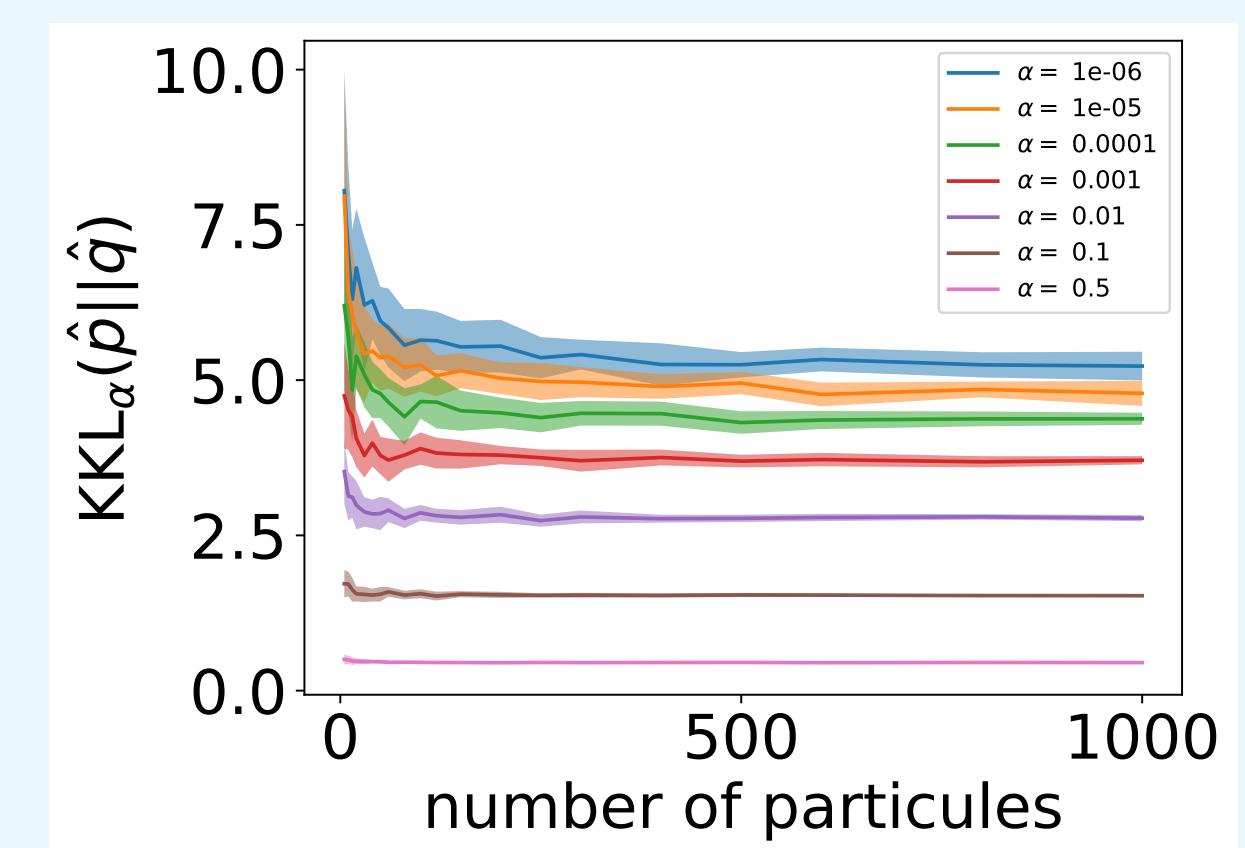
- Consistency of the regularized KKL for empirical measures:

$$\mathbb{E} |\text{KKL}_\alpha(\hat{p}||\hat{q}) - \text{KKL}_\alpha(p||q)| \leq C_{p,\alpha} \frac{\log n}{\sqrt{m \wedge n}} + C'_{p,\alpha} \frac{\log^2 n}{m \wedge n}.$$

The following experiments illustrate the previous theoretical results.



$d = 10$



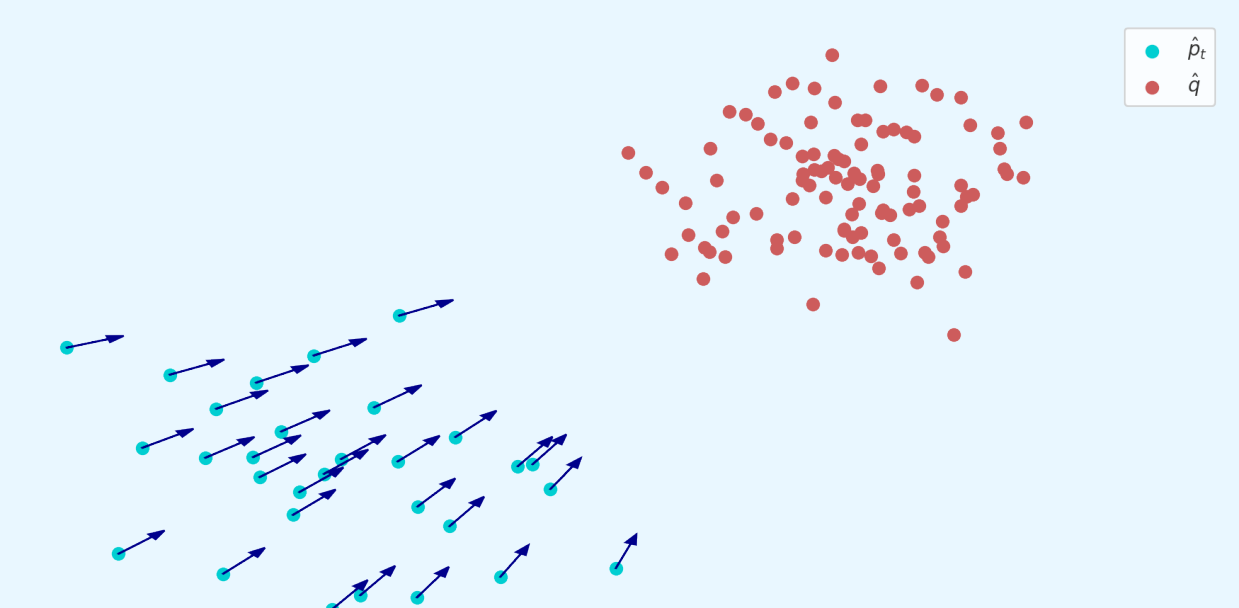
$d = 2$

Sampling experiments

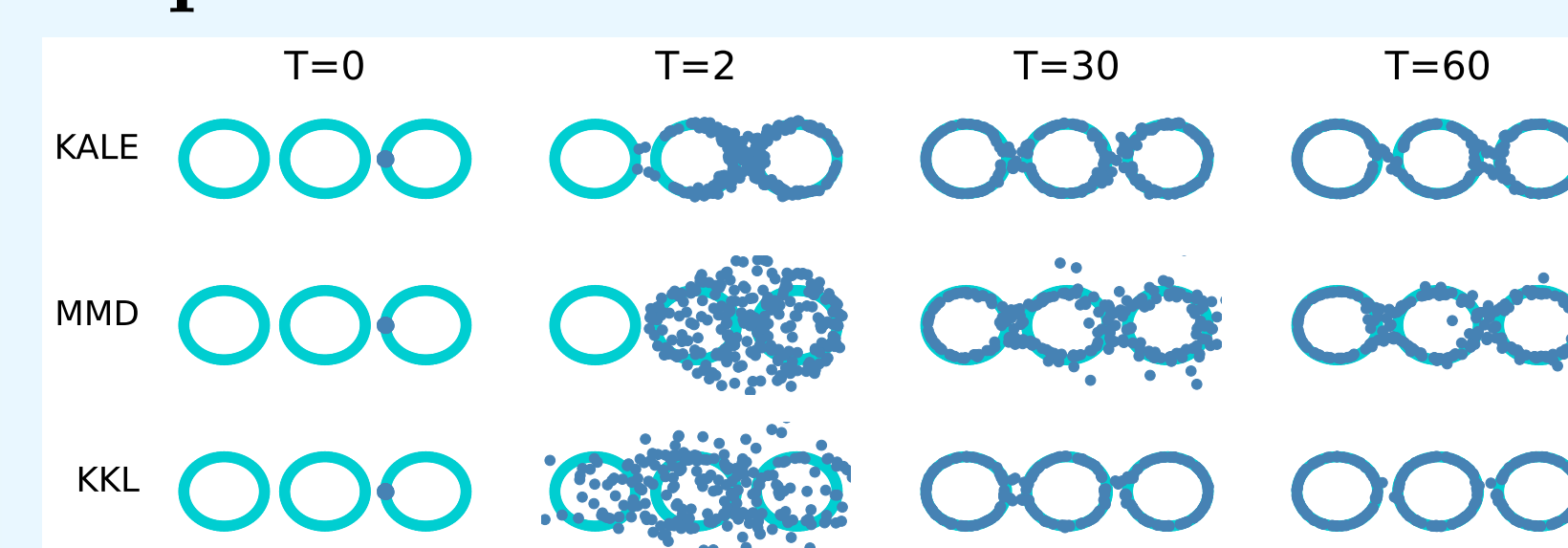
Now we fix \hat{q} , we optimize \hat{p} by a discretisation of the Wasserstein gradient flow of the regularized KKL.

Descent scheme: Let $\hat{p}_t = \frac{1}{n} \sum_{i=1}^n \delta_{x_t^i}$, $\gamma > 0$, $t = 1, \dots, T$.

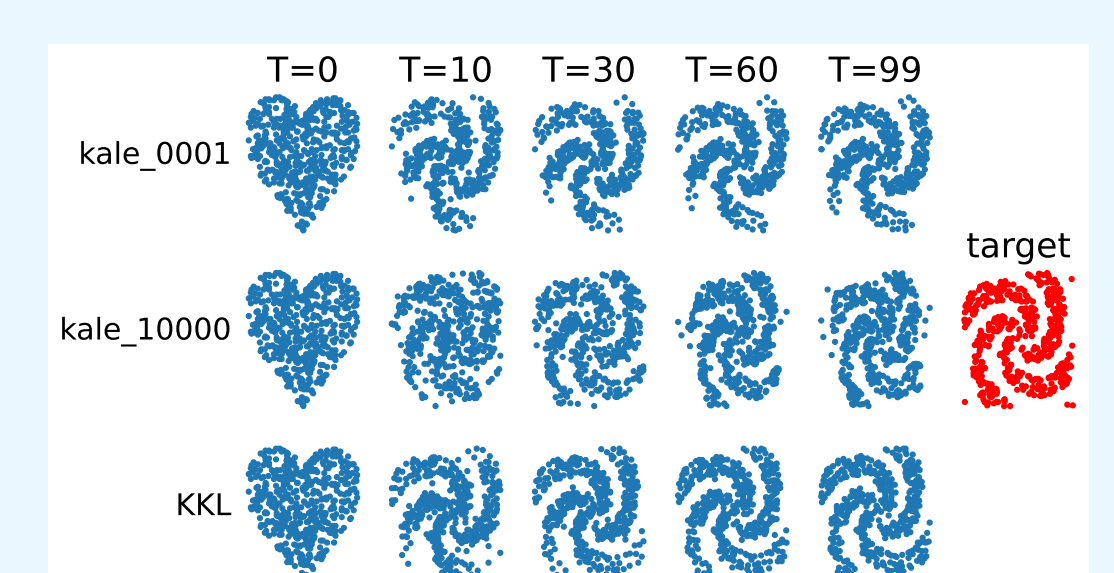
- $x_{t+1}^i = x_t^i - \gamma \nabla_{W_2} \mathcal{F}(\hat{p}_t)(x_t^i)$
- $\hat{p}_{t+1} = (I_d - \gamma \nabla_{W_2} \mathcal{F}(\hat{p}_t))_{\#} \hat{p}$



Experiments:



MMD, KALE and KKL flow for 3 rings target.



Shape transfer

Reference :

[1] Francis Bach. Information theory with kernel methods. *IEEE Transactions on Information Theory*, 69(2):752–775, 2022.