

Statistical and Geometrical properties of the regularized kernel Kullback

Leibler divergence

Clémentine Chazal¹ Anna Korba¹ Francis Bach²

CREST/ENSAE, IP Paris¹ , INRIA, Paris²



Abstract

In this paper [2] we study the properties of the kernel Kullback-Leibler divergence (KKL), introduced in [1], with the aim of performing sampling by using the divergence as the objective of an optimisation problem. Our contributions are to propose a regularized version of the KKL, which is consistent for empirical measures and to derive a Wasserstein gradient of the KKL which has enabled to implement a sampling algorithm.

Introduction and motivations

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

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

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

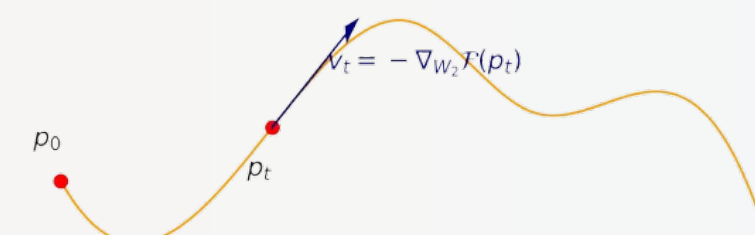
$$\mathcal{F}((I_d + \varepsilon h)_{\#} p) = \mathcal{F}(p) + \varepsilon \langle \nabla_{W_2} \mathcal{F}(p), h \rangle_p + o(\varepsilon),$$

holds, then $\nabla_{W_2} \mathcal{F}(p) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is the **Wasserstein gradient** of \mathcal{F} .

- Analogy between gradient flow and Wasserstein gradient flow

$$\begin{cases} \text{Gradient Flow} \\ x(0) = x_0, \\ x'(t) = -\nabla f(x(t)). \end{cases}$$

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



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 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 $x_1, \dots, x_n \sim p, y_1, \dots, y_m$ and note $\hat{p} = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$ and $\hat{q} = \frac{1}{m} \sum_{j=1}^m \delta_{y_j}$. 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 $S(x) = (\frac{1}{\sqrt{n}} k(x, x_i))_i, T(x) = ((\frac{1}{\sqrt{\alpha}} k(x, x_i))_i, (\frac{1}{\sqrt{1-\alpha}} k(x, y_j))_j)$ and A is a matrix depending on the eigenvalues and eigenvectors of K .

Theoretical properties of the regularized KKL

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

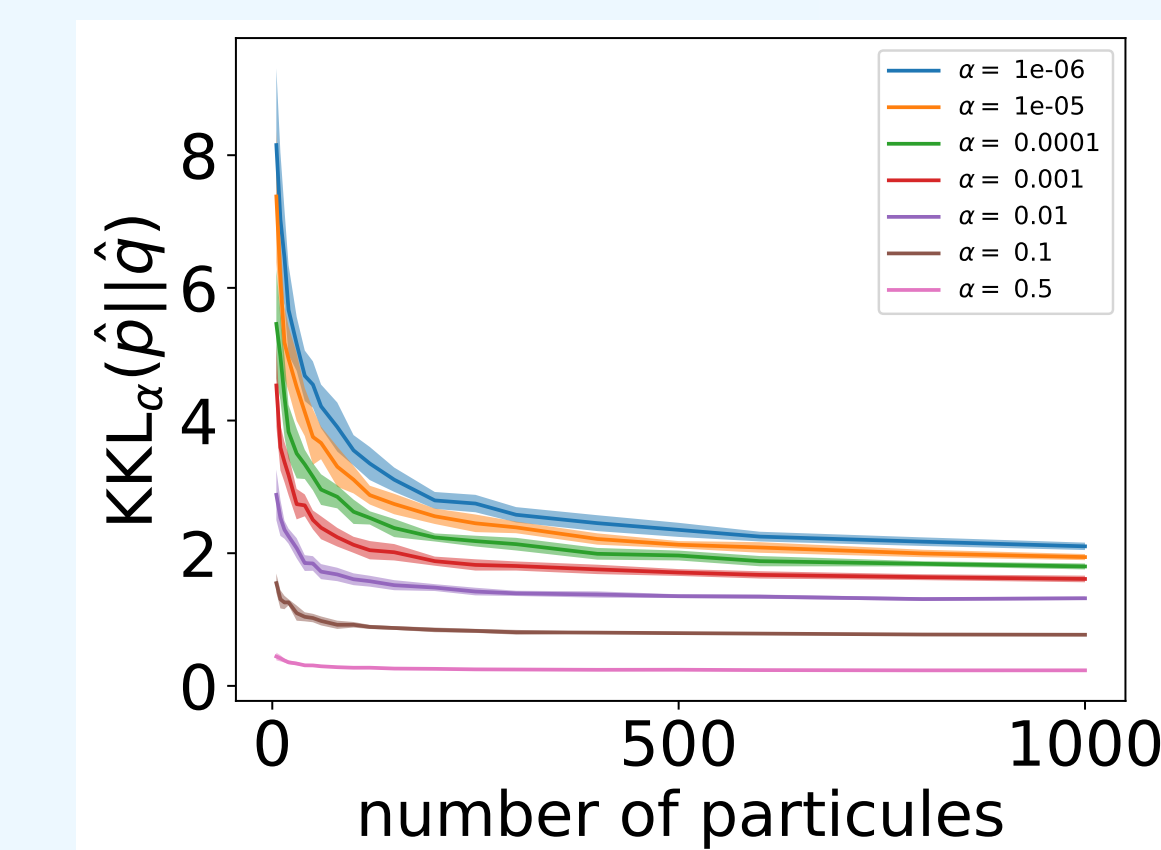
$$\text{KKL}_\alpha(p||q) \xrightarrow{\alpha \rightarrow 0} \text{KKL}(p||q).$$

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

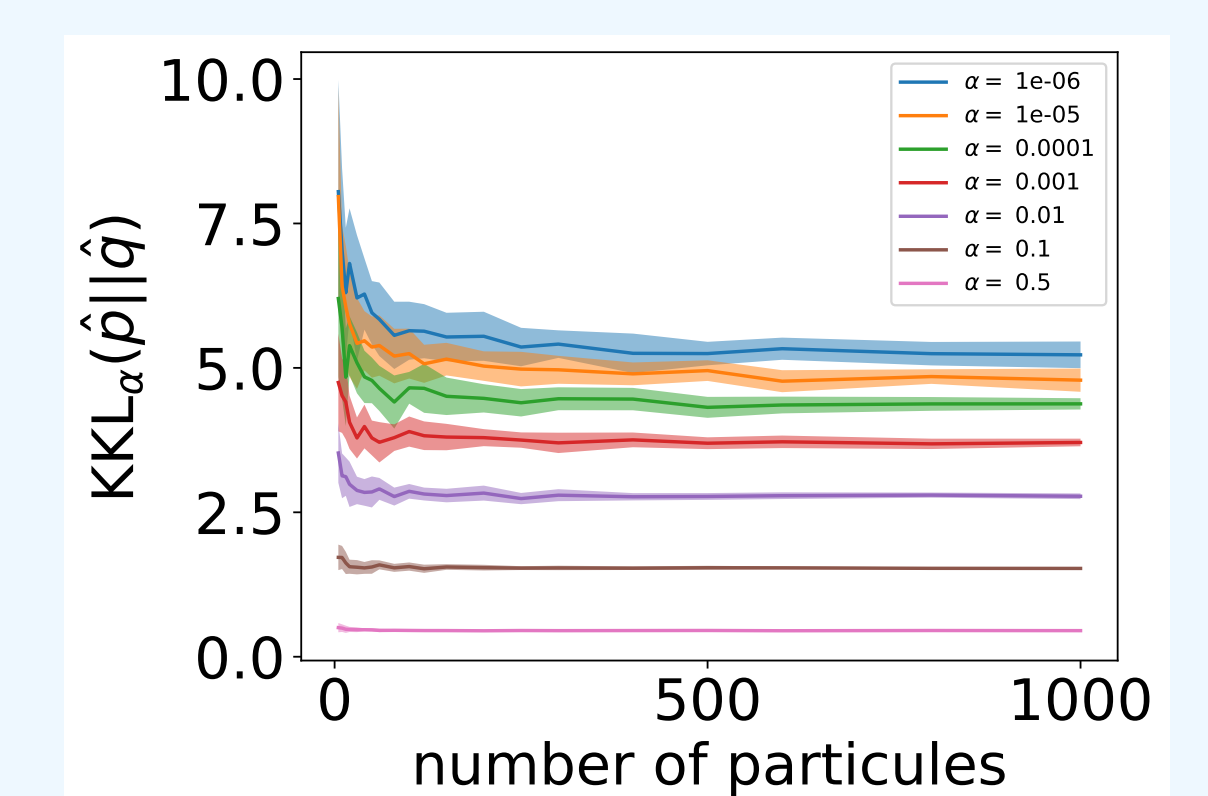
- Consistency of the regularized KKL for empirical measures:

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

The following experiments illustrate the previous theoretical results.



$d = 10$



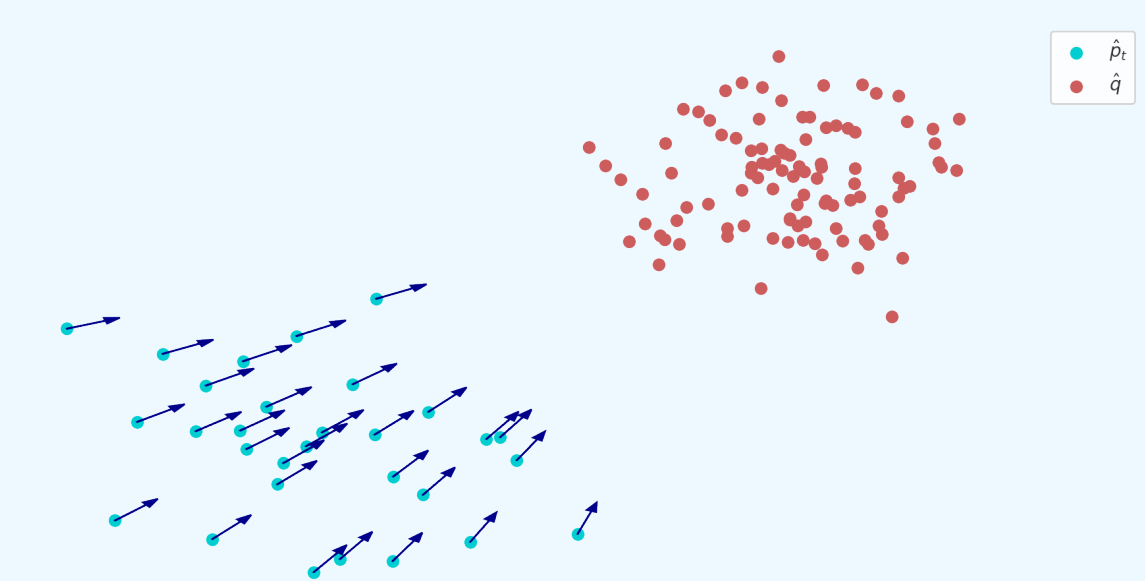
$d = 2$

Sampling experiments

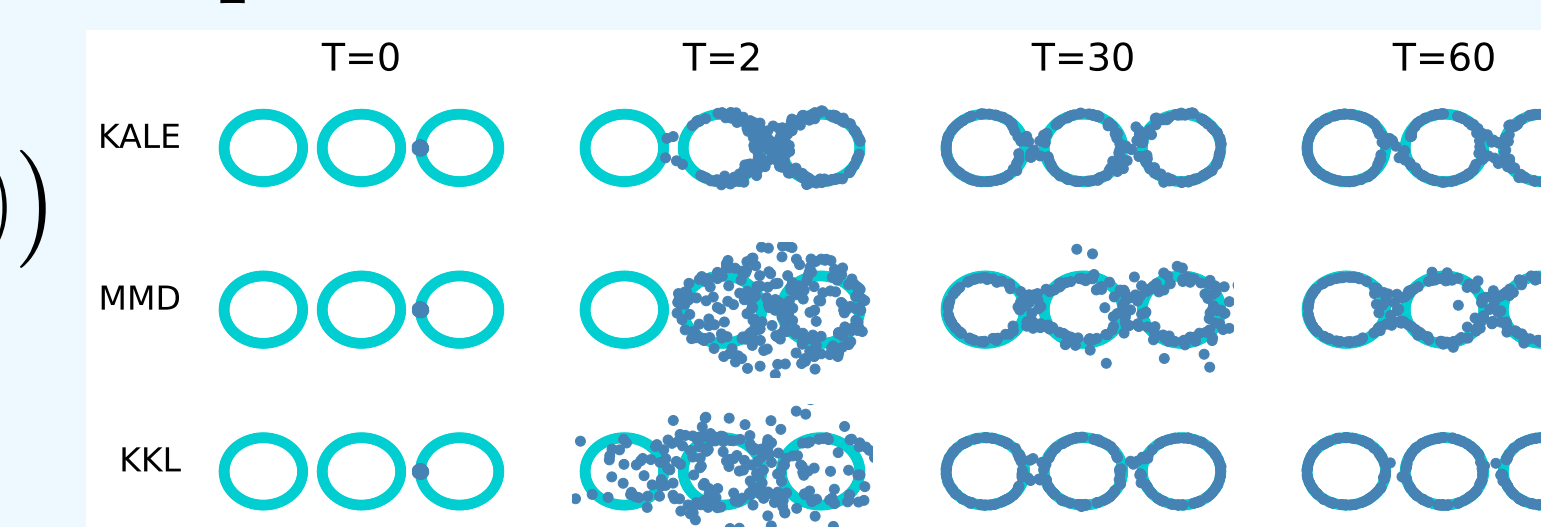
Now we fix \hat{q} , we optimize \hat{p} by a discretisation of the Wasserstein gradient flow of the regularized KKL.

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

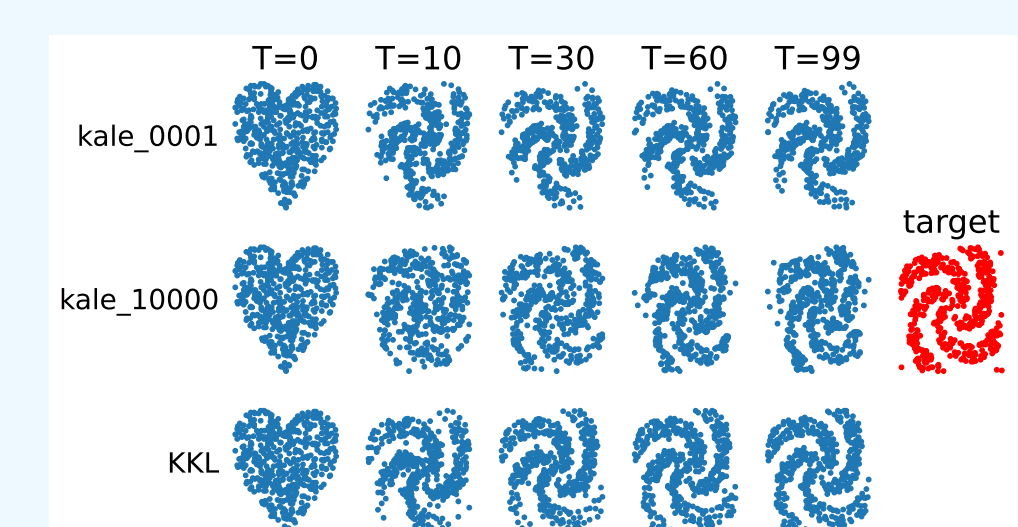
- $x_{t+1}^i = x_t^i - \gamma \nabla_{W_2} \mathcal{F}(\hat{p}_t)(x_t^i)$
- $\hat{p}_{t+1} = (I_d - \gamma \nabla_{W_2} \mathcal{F}(\hat{p}_t))_{\#} \hat{p}_t$



Experiments:



MMD, KALE and KKL flow for 3 rings target.



Shape transfer

Reference

- Francis Bach. Information theory with kernel methods. *IEEE Transactions on Information Theory*, 69(2):752–775, 2022.
- Clémentine Chazal, Anna Korba, and Francis Bach. Statistical and geometrical properties of regularized kernel kullback-leibler divergence. *NeurIPS*, 2024.