Sampling with Kernelized Wasserstein Gradient Flows

Anna Korba ENSAE/CREST

Journées ANR MAGA

Joint work with Adil Salim (Simons Institute), Giulia Luise (UCL), Michael Arbel (INRIA Grenoble), Arthur Gretton (UCL), Pierre-Cyril Aubin-Frankowski (INRIA Paris), Szymon Majewski (Polytechnique), Pierre Ablin (CNRS), Lantian Xu (CMU), Dejan Slepčev (CMU).

Outline

Problem/Motivation

Background

- Part I SVGD algorithm
- Some non-asymptotic results
- Part II : Sampling as optimization of the KSD/MMD
- Preliminaries on Kernel Stein Discrepancy

KSD Descent

Experiments

Theoretical properties of the KSD flow?

Sampling as optimization over distributions

Assume that $\pi \in \mathcal{P}_2(\mathbb{R}^d) = \{ \mu \in \mathcal{P}(\mathbb{R}^d), \int ||x||^2 d\mu(x) < \infty \}.$ The sampling task can be recast as an optimization problem:

$$\pi = \operatorname*{argmin}_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} D(\mu | \pi) := \mathcal{F}(\mu),$$

where *D* is a **dissimilarity functional**.

Starting from an initial distribution $\mu_0 \in \mathcal{P}_2(\mathbb{R}^d)$, one can then consider the **Wasserstein gradient flow** of \mathcal{F} over $\mathcal{P}_2(\mathbb{R}^d)$ to transport μ_0 to π .

Ex 1: Bayesian inference

Problem : Sample from a target distribution π over \mathbb{R}^d , whose density w.r.t. Lebesgue is known up to an intractable constant *Z* :

$$\pi(x) = \frac{\tilde{\pi}(x)}{Z}$$

Ex 1: Bayesian inference

Problem : Sample from a target distribution π over \mathbb{R}^d , whose density w.r.t. Lebesgue is known up to an intractable constant *Z* :

$$\pi(x) = \frac{\tilde{\pi}(x)}{Z}$$

Motivation : Bayesian statistics.

• Let
$$\mathcal{D} = (w_i, y_i)_{i=1,...,N}$$
 observed data.

Assume an underlying model parametrized by θ
 (e.g. p(y|w, θ) gaussian)

 \implies Likelihood: $p(\mathcal{D}|\theta) = \prod_{i=1}^{N} p(y_i|\theta, w_i).$

• Assume also $\theta \sim p$ (prior distribution).

Bayes' rule :
$$\pi(\theta) := p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{Z}$$
, $Z = \int_{\mathbb{R}^d} p(\mathcal{D}|\theta)p(\theta)d\theta$.

Can be written as an optimization problem on \mathcal{P} , e.g.

$$\min_{\mu \in \mathcal{P}} \mathsf{KL}(\mu | \pi)$$
4/6

Ex 2 : Regression with infinite width NN



Assume $\exists \pi \in \mathcal{P}$, $\mathbb{E}[y|X = x] = \mathbb{E}_{Z \sim \pi}[\phi_Z(x)]$, then the above problem corresponds to

 $\min_{\nu \in \mathcal{P}} \mathsf{MMD}^2(\nu,\pi)$

Outline

Problem/Motivation

Background

Part I - SVGD algorithm

Some non-asymptotic results

Part II : Sampling as optimization of the KSD/MMD

Preliminaries on Kernel Stein Discrepancy

KSD Descent

Experiments

Theoretical properties of the KSD flow?

Wasserstein gradient flows (WGF) [Ambrosio et al., 2008]

The differential of $\mu \mapsto \mathcal{F}(\mu)$ evaluated at $\mu \in \mathcal{P}$ is the unique function $\frac{\partial \mathcal{F}(\mu)}{\partial \mu} : \mathbb{R}^d \to \mathbb{R}$ s. t. for any $\mu, \mu' \in \mathcal{P}, \mu' - \mu \in \mathcal{P}$:

 $\lim_{\epsilon \to 0} \frac{1}{\epsilon} (\mathcal{F}(\mu + \epsilon(\mu' - \mu)) - \mathcal{F}(\mu)) = \int_{\mathbb{R}^d} \frac{\partial \mathcal{F}(\mu)}{\partial \mu} (x) (d\mu' - d\mu) (x).$

Wasserstein gradient flows (WGF) [Ambrosio et al., 2008]

The differential of $\mu \mapsto \mathcal{F}(\mu)$ evaluated at $\mu \in \mathcal{P}$ is the unique function $\frac{\partial \mathcal{F}(\mu)}{\partial \mu} : \mathbb{R}^d \to \mathbb{R}$ s. t. for any $\mu, \mu' \in \mathcal{P}, \mu' - \mu \in \mathcal{P}$:

$$\lim_{\epsilon \to 0} \frac{1}{\epsilon} (\mathcal{F}(\mu + \epsilon(\mu' - \mu)) - \mathcal{F}(\mu)) = \int_{\mathbb{R}^d} \frac{\partial \mathcal{F}(\mu)}{\partial \mu} (x) (d\mu' - d\mu) (x).$$

Then $\mu : [0, \infty] \to \mathcal{P}, t \mapsto \mu_t$ satisfies a Wasserstein gradient flow of \mathcal{F} if distributionnally:

$$\frac{\partial \mu_t}{\partial t} = \operatorname{div} \left(\mu_t \nabla_{W_2} \mathcal{F}(\mu_t) \right),$$

where $\nabla_{W_2} \mathcal{F}(\mu) := \nabla \frac{\partial \mathcal{F}(\mu)}{\partial \mu} \in L^2(\mu)$ denotes the Wasserstein gradient of \mathcal{F} .

Some time discretizations

1. Forward method :

$$\mu_{n+1} = exp_{\mu_n}(-\gamma \nabla_{W_2} \mathcal{F}(\mu_n)) = (I - \gamma \nabla_{W_2} \mathcal{F}(\mu_n))_{\#} \mu_n$$

where $exp_{\mu} : L^{2}(\mu) \to \mathcal{P}, \phi \mapsto (I + \phi)_{\#}\mu$, and which corresponds in \mathbb{R}^{d} to:

$$X_{n+1} = X_n - \gamma \nabla_{W_2} \mathcal{F}(\mu_n)(X_n) \sim \mu_{n+1}$$

2. Backward method :

$$\mu_{n+1} = JKO_{\gamma \mathcal{F}}(\mu_n)$$

where $JKO_{\gamma \mathcal{F}}(\nu) = \operatorname*{argmin}_{\mu \in \mathcal{P}} \mathcal{F}(\mu) + \frac{1}{2\gamma} W_2^2(\nu, \mu).$

3. Forward-Backward method, splitting :

$$\nu_{n+1} = (I - \gamma \nabla_{W_2} \mathcal{F}_1(\mu_n))_{\#} \mu_n$$

$$\mu_{n+1} = JKO_{\gamma \mathcal{F}_2}(\nu_{n+1})$$

 $\text{if }\mathcal{F}(\mu)=\mathcal{F}_1(\mu)+\mathcal{F}_2(\mu).$

Choice of the loss function

Many possibilities for the choice of D among Wasserstein distances, *f*-divergences, Integral Probability Metrics... e.g.:

D is the KL (Kullback-Leibler divergence):

$$\mathsf{KL}(\mu|\pi) = \begin{cases} \int_{\mathbb{R}^d} \log\left(\frac{\mu}{\pi}(x)\right) d\mu(x) & \text{if } \mu \ll \pi \\ +\infty & \text{otherwise.} \end{cases}$$

D is the MMD (Maximum Mean Discrepancy):

$$\begin{split} \mathsf{MMD}^2(\mu,\pi) &= \iint_{\mathbb{R}^d} k(x,y) d\mu(x) d\mu(y) \\ &+ \iint_{\mathbb{R}^d} k(x,y) d\pi(x) d\pi(y) - 2 \iint_{\mathbb{R}^d} k(x,y) d\mu(x) d\pi(y). \end{split}$$

where $k : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$ is a p.s.d. kernel.

D is the KSD (Kernel Stein Discrepancy):

$$\begin{split} \mathsf{KSD}^2(\mu|\pi) &= \iint_{\mathbb{R}^d} k_\pi(x, y) d\mu(x) d\mu(y), \text{ where} \\ k_\pi(x, y) &= \nabla \log \pi(x)^T \nabla \log \pi(y) k(x, y) + \nabla \log \pi(x)^T \nabla_y k(x, y) \\ &+ \nabla_x k(x, y)^T \nabla \log(y) + \nabla \cdot_x \nabla_y k(x, y). \end{split}$$

Background on kernels and RKHS [Steinwart and Christmann, 2008]

► Let
$$k : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$$
 a positive, semi-definite kernel, e.g. $k(x, x') = \exp\left(-\frac{\|x-x'\|^2}{h}\right), \exp\left(-\frac{\|x-x'\|}{h}\right), (c + \|x - x'\|)^{\beta}$ with $\beta \in]0, 1[...]$

H_k its corresponding RKHS (Reproducing Kernel Hilbert Space):

$$\mathcal{H}_{k} = \overline{\left\{\sum_{i=1}^{m} \alpha_{i} k(\cdot, \mathbf{x}_{i}); \ \mathbf{m} \in \mathbb{N}; \ \alpha_{1}, \ldots, \alpha_{\mathbf{m}} \in \mathbb{R}; \ \mathbf{x}_{1}, \ldots, \mathbf{x}_{\mathbf{m}} \in \mathbb{R}^{d}\right\}}$$

H_k is a Hilbert space with inner product (.,.)_{H_k} and norm ||.||_{H_k}. It satisfies the reproducing property:

$$\forall \quad f \in \mathcal{H}, \ x \in \mathbb{R}^d, \quad f(x) = \langle f, k(x, .) \rangle_{\mathcal{H}}$$

Consequence : for any *f* ∈ *H_k*, by the reproducing property and Cauchy-Schwartz inequality,

$$\left|\int_{\mathbb{R}^d} f(x) d\pi(x) - \int_{\mathbb{R}^d} f(x) d\mu(x) \right| \leq \|f\|_{\mathcal{H}_k} \operatorname{\mathsf{MMD}}(\mu,\pi).$$

• We denote by \mathcal{H}_k^d the Cartesian product RKHS

MMD and KSD Descent

For discrete measures $\mu_N = \frac{1}{N} \sum_{i=1}^N \delta_{X^i}$, we can define $F(X^1, \ldots, X^N) := \mathcal{F}(\mu_N)$ when it is well defined (e.g., for the MMD or KSD).

In that case, Forward method for WGF = Gradient descent for F on the position of the particles.

▶ If *D* is the MMD, the gradient of *F* is readily obtained as

$$\nabla_{x^i} F(X^1,\ldots,X^N) = \frac{1}{N} \sum_{j=1}^N \nabla_2 k(X^i,X^j) - \int \nabla_2 k(X^i,x) d\pi(x).$$

 \implies requires to know the density π , or at least samples from it ! In contrast, if *D* is the KSD,

$$\nabla_{x^i} F(X^1,\ldots,X^N) = \frac{1}{N} \sum_{j=1}^N \nabla_2 k_{\pi}(X^i,X^j).$$

Algorithm: at each time $n \ge 0$, for any i = 1, ..., N:

$$X_{n+1}^i = X_n - \gamma \nabla_{x^i} F(X_n^1, \dots, X_n^N).$$

The target distribution π is solution of :

$$\min_{\mu \in \mathcal{P}(\mathbb{R}^d)} \mathsf{KL}(\mu | \pi) \tag{1}$$

The target distribution π is solution of :

$$\min_{\mu\in\mathcal{P}(\mathbb{R}^d)}\mathsf{KL}(\mu|\pi)$$

1. Variants of Langevin Monte Carlo (LMC)

[Dalalyan, 2017], [Durmus and Moulines, 2016], [Durmus et al., 2019],

$$X_{k+1} = X_k + \gamma \nabla \log \pi(X_k) + \sqrt{2\gamma} \epsilon_k, \quad \epsilon_k \sim \mathcal{N}(0, I_d)$$

• generates a Markov chain whose law converges to π

- corresponds to a time-discretization of the gradient flow of the KL
- rates of convergence deteriorates quickly in high dimensions

The target distribution π is solution of :

$$\min_{\mu\in\mathcal{P}(\mathbb{R}^d)}\mathsf{KL}(\mu|\pi)$$

1. Variants of Langevin Monte Carlo (LMC)

[Dalalyan, 2017], [Durmus and Moulines, 2016], [Durmus et al., 2019],

$$X_{k+1} = X_k + \gamma \nabla \log \pi(X_k) + \sqrt{2\gamma} \epsilon_k, \quad \epsilon_k \sim \mathcal{N}(0, I_d)$$

- generates a Markov chain whose law converges to π
- corresponds to a time-discretization of the gradient flow of the KL
- rates of convergence deteriorates quickly in high dimensions

2. Variational Inference (VI):

[Alquier and Ridgway, 2017], [Zhang et al., 2018]

- restrict the search space in (1) to a parametric family
- tractable in the large scale setting
- only returns an approximation of π

The target distribution π is solution of :

$$\min_{\mu\in\mathcal{P}(\mathbb{R}^d)}\mathsf{KL}(\mu|\pi)$$

1. Variants of Langevin Monte Carlo (LMC)

[Dalalyan, 2017], [Durmus and Moulines, 2016], [Durmus et al., 2019],

$$X_{k+1} = X_k + \gamma \nabla \log \pi(X_k) + \sqrt{2\gamma} \epsilon_k, \quad \epsilon_k \sim \mathcal{N}(0, I_d)$$

- generates a Markov chain whose law converges to π
- corresponds to a time-discretization of the gradient flow of the KL
- rates of convergence deteriorates quickly in high dimensions

2. Variational Inference (VI):

[Alquier and Ridgway, 2017], [Zhang et al., 2018]

- restrict the search space in (1) to a parametric family
- tractable in the large scale setting
- only returns an approximation of π

 \Longrightarrow Other algorithms can be obtained by discretizing the W_2 gradient flow of the KL...

Forward method for the KL

Problem: $\nabla_{W_2} \operatorname{KL}(\mu_n | \pi) = \nabla \log(\frac{\mu_n}{\pi})$ where μ_n is unknown.

While $\nabla \log \pi$ is known, $\nabla \log \mu_n$ has to be estimated from *N* particles X_n^1, \ldots, X_n^N , e.g. with¹:

1. Kernel Density Estimation (KDE):

$$\mu_n(.)\approx \frac{1}{N}\sum_{i=1}^N k(X_n^i-.)$$

Then,

$$-\nabla_{W_2} \operatorname{KL}(\mu_n | \pi)(.) \approx -\left(\nabla V(.) + \frac{\sum_{i=1}^N \nabla k(.-X_n^i)}{\sum_{i=1}^N k(.-X_n^i)}\right)$$

<u>Remark</u>: it is not the W_2 gradient of some functional (see the next slide)

¹assume a symmetric, translation invariant kernel

2. Blob Method [Carrillo et al., 2019]: Instead of

$$\mathcal{U}(\mu) = \int \log(\mu(\mathbf{x})) d\mu(\mathbf{x}),$$

consider

$$\mathcal{U}_k(\mu) = \int \log(k \star \mu(x)) d\mu(x)$$
, where $k \star \mu(x) = \int k(x-y) d\mu(y)$.

Then,

$$\frac{\partial \mathcal{U}_{k}(\mu)}{\partial \mu}(.) = k \star \left(\frac{\mu}{k \star \mu}\right) + \log(k \star \mu)$$
$$\implies \nabla_{W_{2}}\mathcal{U}_{k}(\mu) = = \nabla k \star \left(\frac{\mu}{k \star \mu}\right) + \underbrace{\nabla \log(k \star \mu)}_{\frac{\nabla k \star \mu}{k \star \mu}}$$

$$\implies \nabla_{W_2} \operatorname{KL}(\mu_n | \pi)(.) \approx - (\nabla V(.) + \sum_{i=1}^{N} \frac{\nabla k(. - X_n^i)}{\sum_{m=1}^{N} k(X_n^i - X_n^m)} + \frac{\sum_{i=1}^{N} \nabla k(. - X_n^i)}{\sum_{i=1}^{N} k(. - X_n^i)} \right)$$

Stein Variational Gradient Descent [Liu and Wang, 2016]

$$-\nabla_{W_2} \operatorname{KL}(\mu_n | \pi)(.) \approx -\frac{1}{N} \left(\sum_{i=1}^N k(.-X_n^j) \nabla V(X_n^i) + \nabla_{X_n^i} k(.-X_n^i) \right)$$

3. Stein Variational Gradient Descent (SVGD)

[Liu, 2017], [Duncan et al., 2019]

- "non parametric" VI, only depends on the choice of some kernel k
- corresponds to a time-discretization of the gradient flow of the KL under a metric depending on k

$$W_k^2(\mu_0,\mu_1) = \inf_{\mu,\nu} \left\{ \int_0^1 \| v_t(x) \|_{\mathcal{H}_k^d}^2 dt(x) : \frac{\partial \mu_t}{\partial t} = \nabla \cdot (\mu_t v_t) \right\}.$$

https://chi-feng.github.io/mcmc-demo/app.html?
algorithm=SVGD&target=banana

Outline

Problem/Motivation

Background

Part I - SVGD algorithm

Some non-asymptotic results

Part II : Sampling as optimization of the KSD/MMD

Preliminaries on Kernel Stein Discrepancy

KSD Descent

Experiments

Theoretical properties of the KSD flow?

SVGD in the ML literature

Empirical performance demonstrated in various tasks:

- Bayesian inference [Liu and Wang, 2016, Feng et al., 2017, Liu and Zhu, 2018, Detommaso et al., 2018]
- learning deep probabilistic models [Wang and Liu, 2016, Pu et al., 2017]
- reinforcement learning [Liu et al., 2017]
- Theoretical guarantees :
 - asymptotic theory: (in continuous time, infinite number of particles) converges asymptotically to π [Lu et al., 2019] when V grows at most polynomially
 - non asymptotic theory: no rates of convergence, only partial results [Korba et al., 2020]

SVGD trick and the kernel integral operator

We assume $\int_{\mathbb{R}^d \times \mathbb{R}^d} k(x, x) d\mu(x) < \infty$ for any $\mu \in \mathcal{P}$. $\Longrightarrow \mathcal{H} \subset L^2(\mu)$. For instance assume $\|k(x, .)\|_{\mathcal{H}_k}^2 = k(x, x) \leq B^2$, then for $f \in \mathcal{H}_k$

$$\begin{split} \|f\|_{L^2(\mu)}^2 &= \int \|f(x)\|^2 d\mu(x) = \int \langle f, k(x, .) \rangle_{\mathcal{H}_k}^2 d\mu(x) \\ &\leq \|f\|_{\mathcal{H}_k}^2 \int k(x, x) d\mu(x) \leq B^2 \|f\|_{\mathcal{H}_k}^2 \end{split}$$

Then, the injection from $\iota : \mathcal{H} \to L^2(\mu)$ admits an adjoint $\iota^* = S_{\mu}$, where $S_{\mu} : L^2(\mu) \to \mathcal{H}$ is defined by:

$$\mathcal{S}_{\mu}f(\cdot)=\int k(x,.)f(x)d\mu(x),\quad f\in L^{2}(\mu).$$

SVGD trick and the kernel integral operator

We assume $\int_{\mathbb{R}^d \times \mathbb{R}^d} k(x, x) d\mu(x) < \infty$ for any $\mu \in \mathcal{P}$. $\Longrightarrow \mathcal{H} \subset L^2(\mu)$. For instance assume $\|k(x, .)\|_{\mathcal{H}_k}^2 = k(x, x) \leq B^2$, then for $f \in \mathcal{H}_k$

$$\begin{split} \|f\|_{L^2(\mu)}^2 &= \int \|f(x)\|^2 d\mu(x) = \int \langle f, k(x, .) \rangle_{\mathcal{H}_k}^2 d\mu(x) \\ &\leq \|f\|_{\mathcal{H}_k}^2 \int k(x, x) d\mu(x) \leq B^2 \|f\|_{\mathcal{H}_k}^2 \end{split}$$

Then, the injection from $\iota : \mathcal{H} \to L^2(\mu)$ admits an adjoint $\iota^* = S_{\mu}$, where $S_{\mu} : L^2(\mu) \to \mathcal{H}$ is defined by:

$$\mathcal{S}_{\mu}f(\cdot)=\int k(x,\cdot)f(x)d\mu(x),\quad f\in L^{2}(\mu).$$

We have for any $f,g\in L_2(\mu) imes \mathcal{H}$:

$$\langle f, \iota g \rangle_{L^2(\mu)} = \langle \iota^* f, g \rangle_{\mathcal{H}} = \langle S_\mu f, g \rangle_{\mathcal{H}}.$$

We will denote $P_{\mu} = \iota \circ S_{\mu}$.

SVGD algorithm

SVGD trick: applying this operator to the W_2 gradient of $KL(\cdot|\pi)$ leads to

$$egin{aligned} \mathcal{P}_{\mu}
abla \log\left(rac{\mu}{\pi}
ight)(\cdot) &= \int
abla \log\left(rac{\mu}{\pi}
ight)(x)k(x,.)d\mu(x) \ &= -\int [
abla \log \pi(x)k(x,\cdot) +
abla_x k(x,\cdot)]d\mu(x), \end{aligned}$$

under appropriate boundary conditions on *k* and π , e.g. $\lim_{\|x\|\to\infty} k(x,\cdot)\pi(x) \to 0.$

SVGD algorithm

SVGD trick: applying this operator to the W_2 gradient of $KL(\cdot|\pi)$ leads to

$$egin{aligned} & \mathcal{P}_{\mu}
abla \log\left(rac{\mu}{\pi}
ight)(\cdot) = \int
abla \log\left(rac{\mu}{\pi}
ight)(x)k(x,.)d\mu(x) \ &= -\int [
abla \log \pi(x)k(x,\cdot) +
abla_x k(x,\cdot)]d\mu(x), \end{aligned}$$

under appropriate boundary conditions on *k* and π , e.g. $\lim_{\|x\|\to\infty} k(x,\cdot)\pi(x) \to 0.$

Algorithm : Starting from *N* i.i.d. samples $(X_0^i)_{i=1,...,N} \sim \mu_0$, SVGD algorithm updates the *N* particles as follows :

$$X_{n+1}^{i} = X_{n}^{i} - \gamma \underbrace{\left[\frac{1}{N} \sum_{j=1}^{N} k(X_{n}^{i}, X_{n}^{j}) \nabla_{X_{n}^{j}} \log \pi(X_{n}^{j}) + \nabla_{X_{n}^{j}} k(X_{n}^{j}, X_{n}^{i})\right]}_{P_{\hat{\mu}_{n}} \nabla \log\left(\frac{\hat{\mu}_{n}}{\pi}\right)(X_{n}^{i}), \quad \text{with } \hat{\mu}_{n} = \frac{1}{N} \sum_{j=1}^{N} \delta_{X_{n}^{j}}}$$

SVGD algorithm

SVGD trick: applying this operator to the W_2 gradient of $KL(\cdot|\pi)$ leads to

$$egin{aligned} & \mathcal{P}_{\mu}
abla \log\left(rac{\mu}{\pi}
ight)(\cdot) = \int
abla \log\left(rac{\mu}{\pi}
ight)(x)k(x,.)d\mu(x) \ &= -\int [
abla \log \pi(x)k(x,\cdot) +
abla_x k(x,\cdot)]d\mu(x), \end{aligned}$$

under appropriate boundary conditions on *k* and π , e.g. $\lim_{\|x\|\to\infty} k(x,\cdot)\pi(x) \to 0.$

Algorithm : Starting from *N* i.i.d. samples $(X_0^i)_{i=1,...,N} \sim \mu_0$, SVGD algorithm updates the *N* particles as follows :

$$X_{n+1}^{i} = X_{n}^{i} - \gamma \underbrace{\left[\frac{1}{N}\sum_{j=1}^{N}k(X_{n}^{i}, X_{n}^{j})\nabla_{X_{n}^{j}}\log\pi(X_{n}^{j}) + \nabla_{X_{n}^{j}}k(X_{n}^{j}, X_{n}^{i})\right]}_{P_{\hat{\mu}_{n}}\nabla\log\left(\frac{\hat{\mu}_{n}}{\pi}\right)(X_{n}^{i}), \quad \text{with } \hat{\mu}_{n} = \frac{1}{N}\sum_{j=1}^{N}\delta_{X_{n}^{j}}}$$

Continuous-time dynamics of SVGD

SVGD gradient flow [Liu, 2017], [Lu et al., 2019]:

$$\frac{\partial \mu_t}{\partial t} + \operatorname{div}(\mu_t V_t) = \mathbf{0}, \qquad V_t := -\mathbf{P}_{\mu_t} \nabla \log\left(\frac{\mu_t}{\pi}\right)$$

Continuous-time dynamics of SVGD

SVGD gradient flow [Liu, 2017], [Lu et al., 2019]:

$$rac{\partial \mu_t}{\partial t} + \textit{div}(\mu_t V_t) = 0, \qquad V_t := - \textit{P}_{\mu_t}
abla \log\left(rac{\mu_t}{\pi}
ight)$$

How fast the KL decreases along SVGD dynamics?

$$\begin{split} \frac{d\operatorname{\mathsf{KL}}(\mu_t|\pi)}{dt} &= \left\langle V_t, \nabla \log\left(\frac{\mu_t}{\pi}\right) \right\rangle_{L^2(\mu_t)} \\ &= -\left\langle \iota S_{\mu_t} \nabla \log\left(\frac{\mu_t}{\pi}\right), \nabla \log\left(\frac{\mu_t}{\pi}\right) \right\rangle_{L^2(\mu_t)} \\ &= -\underbrace{\left\| S_{\mu_t} \nabla \log\left(\frac{\mu_t}{\pi}\right) \right\|_{\mathcal{H}}^2}_{\operatorname{\mathsf{KSD}}^2(\mu_t|\pi)} \operatorname{since} \iota^* = S_{\mu_t} \\ &\leq 0. \end{split}$$

Continuous-time dynamics of SVGD

SVGD gradient flow [Liu, 2017], [Lu et al., 2019]:

$$rac{\partial \mu_t}{\partial t} + \textit{div}(\mu_t V_t) = 0, \qquad V_t := - \textit{P}_{\mu_t}
abla \log\left(rac{\mu_t}{\pi}
ight)$$

How fast the KL decreases along SVGD dynamics?

$$\begin{aligned} \frac{d\operatorname{\mathsf{KL}}(\mu_t|\pi)}{dt} &= \left\langle V_t, \nabla \log\left(\frac{\mu_t}{\pi}\right) \right\rangle_{L^2(\mu_t)} \\ &= -\left\langle \iota S_{\mu_t} \nabla \log\left(\frac{\mu_t}{\pi}\right), \nabla \log\left(\frac{\mu_t}{\pi}\right) \right\rangle_{L^2(\mu_t)} \\ &= -\underbrace{\left\| S_{\mu_t} \nabla \log\left(\frac{\mu_t}{\pi}\right) \right\|_{\mathcal{H}}^2}_{\operatorname{KSD}^2(\mu_t|\pi)} \text{ since } \iota^* = S_{\mu_t} \\ &\leq 0 \end{aligned}$$

On the r.h.s. we have the squared Kernel Stein discrepancy (KSD) [Chwialkowski et al., 2016] or Stein Fisher information at μ_t .

Stein Fisher information

Stationary condition : $\text{KSD}^2(\mu_t | \pi) = \left\| S_{\mu_t} \nabla \log \left(\frac{\mu_t}{\pi} \right) \right\|_{\mathcal{H}}^2 = 0.$

Implies weak convergence of μ_t to π if [Gorham and Mackey, 2017]:

• π is distantly dissipative (e.g. gaussian mixtures):

$$\lim \inf_{r \to \infty} \kappa(r) > 0,$$

$$\kappa(r) = \inf\{-2 \frac{\langle \nabla \log \pi(x) - \nabla \log \pi(y), x - y \rangle}{\|x - y\|_2^2}; \|x - y\|_2^2 = r\}$$

k is translation invariant with a non-vanishing Fourier transform and the sequence is uniformly tight; or k is the IMQ kernel defined by k(x, y) = (c² + ||x − y||₂²)^β for c > 0 and β ∈ [−1,0] (slow decay rate).

Stein Fisher information

Stationary condition : $\text{KSD}^2(\mu_t | \pi) = \left\| S_{\mu_t} \nabla \log \left(\frac{\mu_t}{\pi} \right) \right\|_{\mathcal{H}}^2 = 0.$

Implies weak convergence of μ_t to π if [Gorham and Mackey, 2017]:

• π is distantly dissipative (e.g. gaussian mixtures):

$$\lim \inf_{r \to \infty} \kappa(r) > 0,$$

$$\kappa(r) = \inf\{-2 \frac{\langle \nabla \log \pi(x) - \nabla \log \pi(y), x - y \rangle}{\|x - y\|_2^2}; \|x - y\|_2^2 = r\}$$

k is translation invariant with a non-vanishing Fourier transform and the sequence is uniformly tight; or k is the IMQ kernel defined by k(x, y) = (c² + ||x − y||₂²)^β for c > 0 and β ∈ [−1,0] (slow decay rate).

Proposition:[Korba et al., 2020] if *k* is bounded, $\pi \propto \exp(-V)$ with H_V bounded above and if $\exists C > 0$, $\int ||x||^2 d\mu_t(x) < C$ for all t > 0, then $\text{KSD}^2(\mu_t|\pi) \to 0$.

Convergence of continuous-time dynamics

The convergence of the Stein Fisher information to 0 can be slow. When do we have fast convergence of SVGD dynamics?

 π satisfies the Stein log-Sobolev inequality [Duncan et al., 2019] with constant $\lambda > 0$ if for any μ :

$$\mathsf{KL}(\mu|\pi) \leq rac{1}{2\lambda} \operatorname{KSD}^2(\mu|\pi).$$

Convergence of continuous-time dynamics

The convergence of the Stein Fisher information to 0 can be slow. When do we have fast convergence of SVGD dynamics?

 π satisfies the Stein log-Sobolev inequality [Duncan et al., 2019] with constant $\lambda > 0$ if for any μ :

$$\mathsf{KL}(\mu|\pi) \leq rac{1}{2\lambda} \operatorname{KSD}^2(\mu|\pi).$$

If it holds,

$$rac{d\,\mathsf{KL}(\mu_t|\pi)}{dt} = -\,\mathsf{KSD}^2(\mu_t|\pi) \leq -2\lambda\,\mathsf{KL}(\mu_t|\pi)$$

and by integrating :

$$\mathsf{KL}(\mu_t|\pi) \leq e^{-2\lambda t} \mathsf{KL}(\mu_0|\pi).$$

"Classic" log-Sobolev inequality upper bounds the KL by the Fisher divergence :

$$\mathsf{KL}(\mu|\pi) \leq rac{1}{2\lambda} \left\|
abla \log\left(rac{\mu}{\pi}
ight)
ight\|_{L^2(\mu)}^2$$

satisfied as soon as π is λ -log concave, but it's more general.

"Classic" log-Sobolev inequality upper bounds the KL by the Fisher divergence :

$$\mathsf{KL}(\mu|\pi) \leq rac{1}{2\lambda} \left\|
abla \log\left(rac{\mu}{\pi}
ight)
ight\|_{L^2(\mu)}^2$$

satisfied as soon as π is λ -log concave, but it's more general.

When is Stein log-Sobolev satisfied? not as well known and understood [Duncan et al., 2019], but :

- it fails to hold if k is too regular with respect to π
- some working examples in dimension 1
- whether it holds in higher dimension is more challenging and subject to further research...
Outline

Problem/Motivation

Background

Part I - SVGD algorithm

Some non-asymptotic results

Part II : Sampling as optimization of the KSD/MMD Preliminaries on Kernel Stein Discrepancy KSD Descent

Theoretical properties of the KSD flow?

Gradient descent for $V : \mathbb{R}^d \to \mathbb{R}$ a $C^2(\mathbb{R}^d)$ s.t. $||H_V(x)|| \le M$ for any x.

$$\mathbf{x}_{n+1} = \mathbf{x}_n - \gamma \nabla \mathbf{V}(\mathbf{x}_n).$$

Gradient descent for $V : \mathbb{R}^d \to \mathbb{R}$ a $C^2(\mathbb{R}^d)$ s.t. $||H_V(x)|| \le M$ for any x.

$$x_{n+1} = x_n - \gamma \nabla V(x_n).$$

Denote $x(t) = x_n - t\nabla V(x_n)$ and $\varphi(t) = V(x(t))$. Using Taylor expansion :

$$\varphi(\gamma) = \varphi(\mathbf{0}) + \gamma \varphi'(\mathbf{0}) + \int_{\mathbf{0}}^{\gamma} (\gamma - t) \varphi''(t) dt.$$

Gradient descent for $V : \mathbb{R}^d \to \mathbb{R}$ a $C^2(\mathbb{R}^d)$ s.t. $||H_V(x)|| \le M$ for any x.

$$x_{n+1} = x_n - \gamma \nabla V(x_n).$$

Denote $x(t) = x_n - t \nabla V(x_n)$ and $\varphi(t) = V(x(t))$. Using Taylor expansion :

$$\varphi(\gamma) = \varphi(0) + \gamma \varphi'(0) + \int_0^{\gamma} (\gamma - t) \varphi''(t) dt.$$

Since $(\ddot{x}(t) = 0)$: $\varphi'(0) = \langle \nabla V(x(0)), \dot{x}(0) \rangle = \langle \nabla V(x(0)), -\nabla V(x_n) \rangle = -\|\nabla V(x_n)\|^2$, $\varphi''(t) = \langle \dot{x}(t), H_V(x(t))\dot{x}(t) \rangle \le M \|\dot{x}(t)\|^2 = M \|\nabla V(x_n)\|^2$,

Gradient descent for $V : \mathbb{R}^d \to \mathbb{R}$ a $C^2(\mathbb{R}^d)$ s.t. $||H_V(x)|| \le M$ for any x.

$$x_{n+1} = x_n - \gamma \nabla V(x_n).$$

Denote $x(t) = x_n - t\nabla V(x_n)$ and $\varphi(t) = V(x(t))$. Using Taylor expansion :

$$\varphi(\gamma) = \varphi(0) + \gamma \varphi'(0) + \int_0^{\gamma} (\gamma - t) \varphi''(t) dt.$$

Since $(\ddot{x}(t) = 0)$: $\varphi'(0) = \langle \nabla V(x(0)), \dot{x}(0) \rangle = \langle \nabla V(x(0)), -\nabla V(x_n) \rangle = - \|\nabla V(x_n)\|^2$, $\varphi''(t) = \langle \dot{x}(t), H_V(x(t))\dot{x}(t) \rangle \le M \|\dot{x}(t)\|^2 = M \|\nabla V(x_n)\|^2$, we have

$$V(x_{n+1}) \leq V(x_n) - \gamma \|\nabla V(x_n)\|^2 + M \int_0^\gamma (\gamma - t) \|\nabla V(x_n)\|^2 dt$$
$$V(x_{n+1}) - V(x_n) \leq -\gamma \left(1 - \frac{M\gamma}{2}\right) \|\nabla V(x_n)\|^2.$$

A descent lemma for SVGD

Assume that $\pi \propto \exp(-V)$ where $||H_V(x)|| \leq M$. The Hessian of the KL at μ is an operator on $L^2(\mu)$:

$$\langle f, \textit{Hess}_{\mathsf{KL}(.|\pi)}(\mu)f \rangle_{L^{2}(\mu)} = \mathbb{E}_{X \sim \mu}\left[\langle f(X), H_{V}(X)f(X) \rangle + \|Jf(X)\|_{HS}^{2}
ight]$$

and yet, this operator **is not bounded** due to the Jacobian term.

A descent lemma for SVGD

Assume that $\pi \propto \exp(-V)$ where $||H_V(x)|| \leq M$. The Hessian of the KL at μ is an operator on $L^2(\mu)$:

$$\langle f, \textit{Hess}_{\mathsf{KL}(.|\pi)}(\mu)f \rangle_{L^2(\mu)} = \mathbb{E}_{X \sim \mu} \left[\langle f(X), H_V(X)f(X) \rangle + \|Jf(X)\|_{\mathcal{HS}}^2
ight]$$

and yet, this operator **is not bounded** due to the Jacobian term.

In the case of SVGD, one restricts the descent directions f to \mathcal{H} . Under several assumptions (boundedness of k and ∇k , of Hessian of V and moments on the trajectory) we could show for γ small enough:

$$\mathsf{KL}(\mu_{n+1}|\pi) - \mathsf{KL}(\mu_n|\pi) \leq -c_{\gamma} \underbrace{\left\| S_{\mu_n} \nabla \log \left(\frac{\mu_n}{\pi} \right) \right\|_{\mathcal{H}}^2}_{\mathsf{KSD}^2(\mu_n|\pi)}$$

Rates in terms of the Stein Fisher Information

Consequence of the descent lemma: for γ small enough,

$$\min_{k=1,\dots,n} \mathsf{KSD}^2(\mu_n|\pi) \leq \frac{1}{n} \sum_{k=1}^n \mathsf{KSD}^2(\mu_k|\pi) \leq \frac{\mathsf{KL}(\mu_0|\pi)}{c_\gamma n}.$$

This result does not rely on:

- convexity of V
- nor on Stein log Sobolev inequality
- ▶ but only on smoothness of *V*.

unlike most convergence results on LMC which rely on Log Sobolev inequality or convexity of V.

Rates in terms of the KL objective?

To obtain rates, one may combine a descent lemma (1) of the form

$$\mathsf{KL}(\mu_{n+1}|\pi) - \mathsf{KL}(\mu_n|\pi) \leq -c_\gamma \left\| \mathcal{S}_{\mu_n}
abla \log\left(rac{\mu_n}{\pi}
ight)
ight\|_{\mathcal{H}}^2$$

and the Stein log-Sobolev inequality (2) with constant λ :

$$\mathsf{KL}(\mu_{n+1}|\pi) - \mathsf{KL}(\mu_n|\pi) \underbrace{\leq}_{(1)} - c_{\gamma} \left\| S_{\mu_n} \nabla \log \left(\frac{\mu_n}{\pi} \right) \right\|_{\mathcal{H}}^2 \underbrace{\leq}_{(2)} - c_{\gamma} 2\lambda \, \mathsf{KL}(\mu_n|\pi).$$

Iterating this inequality yields $KL(\mu_n|\pi) \leq (1 - 2c_{\gamma}\lambda)^n KL(\mu_0|\pi)$.

"Classic" approach in optimization [Karimi et al., 2016] or in the analysis of LMC.

Not possible to combine both....

Given that both the kernel and its derivative are bounded, the equation

$$\int \sum_{i=1}^{d} [(\partial_i V(x))^2 k(x,x) - \partial_i V(x)(\partial_i^1 k(x,x) + \partial_i^2 k(x,x)) + \partial_i^1 \partial_i^2 k(x,x)] d\pi(x) < \infty$$
(2)

reduces to a property on V which, as far as we can tell, always holds on \mathbb{R}^d ...

Not possible to combine both....

Given that both the kernel and its derivative are bounded, the equation

$$\int \sum_{i=1}^{d} [(\partial_i V(x))^2 k(x,x) - \partial_i V(x)(\partial_i^1 k(x,x) + \partial_i^2 k(x,x)) + \partial_i^1 \partial_i^2 k(x,x)] d\pi(x) < \infty$$
(2)

reduces to a property on V which, as far as we can tell, always holds on \mathbb{R}^d ...

and this implies that Stein LSI does not hold [Duncan et al., 2019].

Remark : Equation (2) does not hold for :

- k polynomial of order \geq 3, and
- π with exploding β moments with β ≥ 3 (ex: a student distribution, which belongs to P the set of distributions with bounded second moment).

Experiments



Figure: The particle implementation of the SVGD algorithm illustrates the convergence of $\text{KSD}^2(\mu_n|\pi)$ to 0.

We already have a bound on μ_n versus π . What about $\hat{\mu}_n$? Recall that the practical SVGD implementation is :

$$X_{n+1}^{i} = X_{n}^{i} - \gamma P_{\hat{\mu}_{n}} \nabla \log \left(\frac{\hat{\mu}_{n}}{\pi}\right) (X_{n}^{i}), \qquad \hat{\mu}_{n} = \frac{1}{N} \sum_{j=1}^{N} \delta_{X_{n}^{j}}.$$

where $\hat{\mu}_n$ denotes the empirical distribution of the interacting particles.

We already have a bound on μ_n versus π . What about $\hat{\mu}_n$? Recall that the practical SVGD implementation is :

$$X_{n+1}^{i} = X_{n}^{i} - \gamma P_{\hat{\mu}_{n}} \nabla \log \left(\frac{\hat{\mu}_{n}}{\pi}\right) (X_{n}^{i}), \qquad \hat{\mu}_{n} = \frac{1}{N} \sum_{i=1}^{N} \delta_{X_{n}^{i}}.$$

where $\hat{\mu}_n$ denotes the empirical distribution of the interacting particles.

Propagation of chaos result Let $n \ge 0$ and T > 0. Under boundedness and Lipschitzness assumptions for all $k, \nabla k, V$; for any $0 \le n \le \frac{T}{\gamma}$ we have :

$$\mathbb{E}[W_2^2(\bar{\mu}_n,\hat{\mu}_n)] \leq \frac{1}{2} \left(\frac{1}{\sqrt{N}} \sqrt{var(\mu_0)} e^{LT}\right) (e^{2LT} - 1)$$

where *L* is a constant depending on *k* and π and $\bar{\mu}_n = \frac{1}{N} \sum_{j=1}^N \delta_{\bar{X}_n^j}$ with $\bar{X}_n^j \sim \mu_n$ i.i.d.

Contributions and openings

- First rates of convergence for SVGD, using techniques from optimal transport and optimization (discrete time infinite number of particles)
- Propagation of chaos bound (finite number of particles regime)

Open questions

Rates in KL?

- Propagation of chaos : weaker assumptions? uniform in time (UIT)?
- Is it possible to obtain a unified convergence bound (decreasing as n, N → ∞)? (requires UIT)

$$D(\widehat{\mu}_n, \pi) \leq A_n + B_N$$

how good is SVGD quantisation?

Other kernels?

SVGD dynamics also appear in black-box variational inference and Gans [Chu et al., 2020], where the kernel is *the neural tangent kernel* and **depends on the current distribution** ($k \implies k_{\mu_n}$)

Outline

Problem/Motivation

Background

Part I - SVGD algorithm

Some non-asymptotic results

Part II : Sampling as optimization of the KSD/MMD Preliminaries on Kernel Stein Discrepancy KSD Descent Experiments

Theoretical properties of the KSD flow?

Outline

- **Problem/Motivation**
- Background
- Part I SVGD algorithm
- Some non-asymptotic results
- Part II : Sampling as optimization of the KSD/MMD
- Preliminaries on Kernel Stein Discrepancy
- KSD Descent
- Experiments
- Theoretical properties of the KSD flow?

Kernel Stein Discrepancy [Chwialkowski et al., 2016, Liu et al., 2016]

For $\mu, \pi \in \mathcal{P}_2(\mathbb{R}^d)$, the KSD of μ relative to π is defined as $\mathsf{KSD}^2(\mu|\pi) = \iint k_{\pi}(x, y) d\mu(x) d\mu(y),$

where $k_{\pi} : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$ is the **Stein kernel**, defined through

- the score function $s(x) = \nabla \log \pi(x)$,
- ▶ a p.s.d. kernel $k : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}, k \in C^2(\mathbb{R}^d)^2$

For
$$x, y \in \mathbb{R}^d$$
,
 $k_{\pi}(x, y) = s(x)^T s(y) k(x, y) + s(x)^T \nabla_2 k(x, y)$
 $+ \nabla_1 k(x, y)^T s(y) + \nabla \cdot \nabla_2 k(x, y)$
 $= \sum_{i=1}^d \frac{\partial \log \pi(x)}{\partial x_i} \cdot \frac{\partial \log \pi(y)}{\partial y_i} \cdot k(x, y) + \frac{\partial \log \pi(x)}{\partial x_i} \cdot \frac{\partial k(x, y)}{\partial y_i}$
 $+ \frac{\partial \log \pi(y)}{\partial y_i} \cdot \frac{\partial k(x, y)}{\partial x_i} + \frac{\partial^2 k(x, y)}{\partial x_i \partial y_i} \in \mathbb{R}.$
²e.g. : $k(x, y) = \exp(-||x - y||^2/h)$

We have seen that the KSD² is also as a kernelized Fisher divergence $(\|\nabla \log(\frac{\mu}{\pi})\|_{L^{2}(\mu)}^{2})$:

$$\mathsf{KSD}^2(\mu|\pi) = \left\| \mathcal{S}_{\mu,k} \nabla \log\left(\frac{\mu}{\pi}\right) \right\|_{\mathcal{H}_k}^2, \ \mathcal{S}_{\mu,k} : f \mapsto \int f(x) k(x,.) d\mu(x).$$

$$\begin{split} \left\| S_{\mu,k} \nabla \log\left(\frac{\mu}{\pi}\right) \right\|_{\mathcal{H}_{k}}^{2} &= \langle S_{\mu,k} \nabla \log\left(\frac{\mu}{\pi}\right), S_{\mu,k} \nabla \log\left(\frac{\mu}{\pi}\right) \rangle_{\mathcal{H}_{k}} \\ &= \int \int \nabla \log\left(\frac{\mu}{\pi}(x)\right) \nabla \log\left(\frac{\mu}{\pi}(y)\right) k(x,y) d\mu(x) d\mu(y) \end{split}$$

+ I.P.P 3 times³ ($\nabla \log \mu(x) d\mu(x) = \nabla \mu(x) dx$) recovers the formula of the previous slide.

³Assuming appropriate boundary conditions, e.g. $\lim_{\|x\|\to\infty} k(x,.)\mu(x) = 0.$

Stein identity and link with MMD

Under mild assumptions on *k* and π , the Stein kernel k_{π} is p.s.d. and satisfies a **Stein identity** [Oates et al., 2017]

$$\int_{\mathbb{R}^d} k_{\pi}(x,.) d\pi(x) = 0.$$

Consequently, **KSD is an MMD** with kernel k_{π} , since:

$$\begin{split} \mathsf{MMD}^2(\mu|\pi) &= \int k_\pi(x,y) d\mu(x) d\mu(y) + \int k_\pi(x,y) d\pi(x) d\pi(y) \\ &- 2 \int k_\pi(x,y) d\mu(x) d\pi(y) \\ &= \int k_\pi(x,y) d\mu(x) d\mu(y) \\ &= \mathsf{KSD}^2(\mu|\pi) \end{split}$$

Outline

- **Problem/Motivation**
- Background
- Part I SVGD algorithm
- Some non-asymptotic results
- Part II : Sampling as optimization of the KSD/MMD
- Preliminaries on Kernel Stein Discrepancy
- **KSD** Descent
- Experiments
- Theoretical properties of the KSD flow?

KSD Descent - algorithms

We propose two ways to implement KSD Descent:

Algorithm 1 KSD Descent GD

Input: initial particles $(x_0^i)_{i=1}^N \sim \mu_0$, number of iterations M, step-size γ for n = 1 to M do $[x_{n+1}^i]_{i=1}^N = [x_n^i]_{i=1}^N - \frac{2\gamma}{N^2} \sum_{j=1}^N [\nabla_2 k_\pi(x_n^j, x_n^i)]_{i=1}^N$, end for Return: $[x_M^i]_{i=1}^N$.

Algorithm 2 KSD Descent L-BFGS

Input: initial particles $(x_0^i)_{i=1}^N \sim \mu_0$, tolerance tol

Return: $[x_*^i]_{i=1}^N = L$ -BFGS $(L, \nabla L, [x_0^i]_{i=1}^N, \text{tol})$.

L-BFGS [Liu and Nocedal, 1989] is a quasi Newton algorithm that is faster and more robust than Gradient Descent, and **does not** require the choice of step-size!

L-BFGS

L-BFGS (Limited memory Broyden–Fletcher–Goldfarb–Shanno algorithm) is a quasi-Newton method:

$$x_{n+1} = x_n - \gamma_n B_n^{-1} \nabla L(x_n) := x_n + \gamma_n d_n$$
(3)

where B_n^{-1} is a p.s.d. matrix approximating the inverse Hessian at x_n . Step1. (requires ∇L) It computes a cheap version of d_n based on

BFGS recursion:

$$B_{n+1}^{-1} = \left(I - \frac{\Delta x_n y_n^T}{y_n^T \Delta x_n}\right) B_n^{-1} \left(I - \frac{y_n \Delta x_n^T}{y_n^T \Delta x_n}\right) + \frac{\Delta x_n \Delta x_n^T}{y_n^T \Delta x_n}$$

where
$$\Delta x_n = x_{n+1} - x_n$$

 $y_n = \nabla L(x_{n+1}) - \nabla L(x_n)$

Step2. (requires *L* and ∇L) A line-search is performed to find the best step-size in (3) :

$$L(x_n + \gamma_n d_n) \le L(x_n) + c_1 \gamma_n \nabla L(x_n)^T d_n$$
$$\nabla L(x_n + \gamma_n d_n)^T d_n \ge c_2 \nabla L(x_n)^T d_n$$

Outline

- **Problem/Motivation**
- Background
- Part I SVGD algorithm
- Some non-asymptotic results
- Part II : Sampling as optimization of the KSD/MMD
- Preliminaries on Kernel Stein Discrepancy
- KSD Descent

Experiments

Theoretical properties of the KSD flow?

Toy experiments - 2D standard gaussian



The green points represent the initial positions of the particles. The light grey curves correspond to their trajectories.

SVGD vs KSD Descent - importance of the step-size



Convergence speed of KSD and SVGD on a Gaussian problem in 1D, with 30 particles.

2D mixture of (isolated) Gaussians - failure cases



The green crosses indicate the initial particle positions the blue ones are the final positions The light red arrows correspond to the score directions.

More initializations



Green crosses : initial particle positions Blue crosses : final positions In the paper, we explain how particles can get stuck in planes of symmetry of the target π .

- We show that if a stationary measure µ∞ is full support, then F(µ∞) = 0.
- but we also show that if supp(µ₀) ⊂ M, where M is a plane of symmetry of π, then for any time *t* it remains true for µ_t: supp(µ_t) ⊂ M.

Isolated Gaussian mixture - annealing

Add an inverse temperature variable $\beta : \pi^{\beta}(x) \propto \exp(-\beta V(x))$, with $0 < \beta \le 1$ (i.e. multiply the score by β .)



This is a hard problem, even for Langevin diffusions, where tempering strategies also have been proposed.

Beyond Log-concavity: Provable Guarantees for Sampling Multi-modal Distributions using Simulated Tempering Langevin Monte Carlo. Rong Ge, Holden Lee, Andrej Risteski. 2017.

Real world experiments (10 particles)



Bayesian logistic regression.

Accuracy of the KSD descent and SVGD for 13 datasets ($d \approx 50$). Both methods yield similar results. KSD is better by 2% on one dataset.

Hint: convex likelihood.

Bayesian ICA.

Each dot is the Amari distance between an estimated matrix and the true unmixing matrix ($d \le 8$). **KSD is not better than random.** Hint: highly non-convex likelihood.

So.. when does it work?



Comparison of KSD Descent and Stein points on a "banana" distribution. Green points are the initial points for KSD Descent. Both methods work successfully here, **even though it is not a log-concave distribution.**

We posit that KSD Descent succeeds because there is no saddle point in the potential.

Outline

- **Problem/Motivation**
- Background
- Part I SVGD algorithm
- Some non-asymptotic results
- Part II : Sampling as optimization of the KSD/MMD
- Preliminaries on Kernel Stein Discrepancy
- KSD Descent
- Experiments
- Theoretical properties of the KSD flow?

First strategy : functional inequality? $\mathcal{F}(\mu|\pi) = \iint k_{\pi}(x, y) d\mu(x) d\mu(y).$

We have

$$rac{\partial \mathcal{F}(\mu)}{\partial \mu} = \int k_{\pi}(x,.) d\mu(x) = \mathbb{E}_{x \sim \mu}[k_{\pi}(x,.)]$$

and under appropriate growth assumptions on k_{π} :

$$\nabla_{W_2}\mathcal{F}(\mu) = \mathbb{E}_{\boldsymbol{x} \sim \mu}[\nabla_2 \boldsymbol{k}_{\pi}(\boldsymbol{x}, \cdot)],$$

Hence

$$\begin{aligned} \frac{d\mathcal{F}(\mu_t)}{dt} &= \langle \nabla_{W_2} \mathcal{F}(\mu_t), -\nabla_{W_2} \mathcal{F}(\mu_t) \rangle_{L^2(\mu_t)} \\ &= -\mathbb{E}_{\mathbf{y} \sim \mu_t} \left[\|\mathbb{E}_{\mathbf{x} \sim \mu_t} [\nabla_2 \mathbf{k}_{\pi}(\mathbf{x}, \mathbf{y})] \|^2 \right] \leq \mathbf{0}. \end{aligned}$$

 \Rightarrow Difficult to identify a functional inequality to relate $d\mathcal{F}(\mu_t)/dt$ to $\mathcal{F}(\mu_t)$, and establish convergence in continuous time (similar to [Arbel et al., 2019]).

Second strategy : geodesic convexity of the KSD?

Let $\psi \in C_c^{\infty}(\mathbb{R}^d)$ and the path $\rho_t = (I + t\nabla \psi)_{\#}\mu$ for $t \in [0, 1]$. Define the quadratic form $\operatorname{Hess}_{\mu} \mathcal{F}(\psi, \psi) := \frac{d^2}{dt^2}\Big|_{t=0} \mathcal{F}(\rho_t)$, which is related to the W_2 **Hessian of** \mathcal{F} at μ .

For $\psi \in C^{\infty}_{c}(\mathbb{R}^{d})$, we have

$$\begin{aligned} \mathsf{Hess}_{\mu}\,\mathcal{F}(\psi,\psi) &= \mathbb{E}_{x,y\sim\mu}\left[\nabla\psi(x)^{\mathsf{T}}\nabla_{\mathsf{1}}\nabla_{\mathsf{2}}k_{\pi}(x,y)\nabla\psi(y)\right] \\ &+ \mathbb{E}_{x,y\sim\mu}\left[\nabla\psi(x)^{\mathsf{T}}H_{\mathsf{1}}k_{\pi}(x,y)\nabla\psi(x)\right].\end{aligned}$$

The first term is always positive but not the second one.

 \implies the KSD is not convex w.r.t. W_2 geodesics.
Third strategy : curvature near equilibrium?

What happens near equilibrium π ? the second term vanishes due to the Stein property of k_{π} and :

$$\operatorname{Hess}_{\pi}\mathcal{F}(\psi,\psi) = \|\mathcal{S}_{\pi,k_{\pi}}\mathcal{L}_{\pi}\psi\|_{\mathcal{H}_{k_{\pi}}}^{2} \geq 0$$

where

$$\mathcal{L}_{\pi}: f \mapsto -\Delta f - \langle \nabla \log \pi, \nabla f \rangle_{\mathbb{R}^d}$$

 $\mathcal{S}_{\mu,k_{\pi}}: f \mapsto \int k_{\pi}(x,.)f(x)d\mu(x) \in \mathcal{H}_{k_{\pi}} = \overline{\{k_{\pi}(x,.), x \in \mathbb{R}^d\}}$

Question: can we bound from below the Hessian at π by a quadratic form on the tangent space of $\mathcal{P}_2(\mathbb{R}^d)$ at $\pi \ (\subset L^2(\pi))$?

$$\|\boldsymbol{\mathcal{S}}_{\pi,\boldsymbol{k}_{\pi}}\mathcal{L}_{\pi}\psi\|_{\mathcal{H}_{\boldsymbol{k}_{\pi}}}^{2} = \operatorname{Hess}_{\pi}\mathcal{F}(\psi,\psi) \geq \lambda \|\nabla\psi\|_{L^{2}(\pi)}^{2} ?$$

That would imply exponential decay of \mathcal{F} near π .

Curvature near equilibrium - negative result

The previous inequality

$$\|\boldsymbol{S}_{\pi,\boldsymbol{k}_{\pi}}\mathcal{L}_{\pi}\psi\|_{\mathcal{H}_{\boldsymbol{k}_{\pi}}}^{2} \geq \lambda \|\nabla\psi\|_{\boldsymbol{L}^{2}(\pi)}^{2}$$

can be seen as a kernelized version of the Poincaré inequality for π :

$$\|\mathcal{L}_{\pi}\psi\|_{L_{2}(\pi)}^{2} \geq \lambda_{\pi}\|\nabla\psi\|_{L_{2}(\pi)}^{2}.$$

can be written:

$$\langle \psi, \mathcal{P}_{\pi,k_{\pi}}\psi \rangle_{L_{2}(\pi)} \geq \lambda \langle \psi, \mathcal{L}_{\pi}^{-1}\psi \rangle_{L_{2}(\pi)},$$

where $\mathcal{P}_{\pi,k_{\pi}}: L^{2}(\pi) \rightarrow L^{2}(\pi), f \mapsto \int k_{\pi}(x,.)f(x)d\pi(x).$

Theorem : Let $\pi \propto e^{-V}$. Assume that $V \in C^2(\mathbb{R}^d)$, ∇V is Lipschitz and \mathcal{L}_{π} has discrete spectrum. Then exponential decay near equilibium does not hold.

Contributions

Pros:

- KSD Descent is a very simple algorithm, and can be used with L-BFGS [Liu and Nocedal, 1989] (fast, and does not require the choice of a step-size as in SVGD)
- works well on log-concave targets (unimodal gaussian, Bayesian logistic regression with gaussian priors) or "nice" distributions (banana)

Cons:

- ► KSD is not convex w.r.t. W₂, and no exponential decay near equilibrium holds
- does not work well on non log-concave targets (mixture of isolated gaussians, Bayesian ICA)

Open questions

- explain the convergence of KSD Descent when π is log-concave?
- quantify propagation of chaos ? (KSD for a finite number of particles vs infinite - but non uniformly Lipschitz vector field)
- how good is KSD quantisation?

Code

Python package to try KSD descent yourself: pip install ksddescent

- website: pierreablin.github.io/ksddescent/
- It also features pytorch/numpy code for SVGD.

```
>>> import torch
>>> from ksddescent import ksdd_lbfgs
>>> n, p = 50, 2
>>> x0 = torch.rand(n, p) # start from uniform distribution
>>> score = lambda x: x # simple score function
>>> x = ksdd_lbfgs(x0, score) # run the algorithm
```

Ongoing work - quantization of these methods - gaussian target



Figure: (a)-(c) Final states of the algorithms for 1000 particles, after 1e4 iterations. The kernel bandwidth for all algorithms is set to 1.

Ongoing work - quantization of these methods - gaussian target



Figure: target distribution: $\pi = \mathcal{N}(0, 1/dI_d)$.

Ongoing work - quantization of these methods - gaussian target



Figure: (same target π) Importance of the choice of the bandwidth in the MMD evaluation metric when evaluating the final states, in 2D. From Left to Right: (evaluation) MMD bandwidth = 1, 0.7, 0.3.

References I

 Alquier, P. and Ridgway, J. (2017). Concentration of tempered posteriors and of their variational approximations. *arXiv preprint arXiv:1706.09293*.

Ambrosio, L., Gigli, N., and Savaré, G. (2008). Gradient flows: in metric spaces and in the space of probability measures. Springer Science & Business Media.

 Arbel, M., Korba, A., Salim, A., and Gretton, A. (2019).
 Maximum mean discrepancy gradient flow.
 In Advances in Neural Information Processing Systems, pages 6481–6491.

References II

- Carrillo, J. A., Craig, K., and Patacchini, F. S. (2019).
 A blob method for diffusion.
 Calculus of Variations and Partial Differential Equations, 58(2):1–53.
- Chen, W. Y., Barp, A., Briol, F.-X., Gorham, J., Girolami, M., Mackey, L., Oates, C., et al. (2019).
 Stein point Markov Chain Monte Carlo.
 Proceedings of the 36th International Conference on Machine Learning,.
- Chen, W. Y., Mackey, L., Gorham, J., Briol, F.-X., and Oates, C. J. (2018). Stein points.

Proceedings of the 35th International Conference on Machine Learning,.

References III

Chu, C., Minami, K., and Fukumizu, K. (2020). The equivalence between stein variational gradient descent and black-box variational inference. arXiv preprint arXiv:2004.01822.

- Chwialkowski, K., Strathmann, H., and Gretton, A. (2016).
 A kernel test of goodness of fit.
 In International conference on machine learning.
- Dalalyan, A. S. (2017).

Further and stronger analogy between sampling and optimization: Langevin monte carlo and gradient descent. *arXiv preprint arXiv:1704.04752*.

References IV

 Detommaso, G., Cui, T., Marzouk, Y., Spantini, A., and Scheichl, R. (2018).
 A stein variational newton method. In *Advances in Neural Information Processing Systems*, pages 9169–9179.
 Domingo-Enrich, C., Bietti, A., Vanden-Eijnden, E., and Bruna, J. (2021).

On energy-based models with overparametrized shallow neural networks.

arXiv preprint arXiv:2104.07531.



References V

- Durmus, A., Majewski, S., and Miasojedow, B. (2019). Analysis of langevin monte carlo via convex optimization. *Journal of Machine Learning Research*, 20(73):1–46.
- Durmus, A. and Moulines, E. (2016). Sampling from strongly log-concave distributions with the unadjusted langevin algorithm. arXiv preprint arXiv:1605.01559, 5.
- Feng, Y., Wang, D., and Liu, Q. (2017). Learning to draw samples with amortized stein variational gradient descent.

arXiv preprint arXiv:1707.06626.

References VI

 Fisher, M. A., Nolan, T., Graham, M. M., Prangle, D., and Oates, C. J. (2020).
 Measure transport with kernel Stein discrepancy. arXiv preprint arXiv:2010.11779.

 Gorham, J. and Mackey, L. (2017).
 Measuring sample quality with kernels.
 In Proceedings of the 34th International Conference on Machine Learning-Volume 70, pages 1292–1301. JMLR. org.

Hodgkinson, L., Salomone, R., and Roosta, F. (2020). The reproducing Stein kernel approach for post-hoc corrected sampling. arXiv preprint arXiv:2001.09266.

References VII

 Kanagawa, H., Jitkrittum, W., Mackey, L., Fukumizu, K., and Gretton, A. (2020).
 A kernel Stein test for comparing latent variable models. arXiv preprint arXiv:1907.00586.

- Karimi, H., Nutini, J., and Schmidt, M. (2016).
 Linear convergence of gradient and proximal-gradient methods under the polyak-lojasiewicz condition.
 In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 795–811.
 Springer.
 - Korba, A., Salim, A., Arbel, M., Luise, G., and Gretton, A. (2020).

A non-asymptotic analysis for stein variational gradient descent.

arXiv preprint arXiv:2006.09797.

References VIII

Liu, C. and Zhu, J. (2018). Riemannian stein variational gradient descent for bayesian inference.

In Thirty-second aaai conference on artificial intelligence.

Liu, D. C. and Nocedal, J. (1989).
 On the limited memory BFGS method for large scale optimization.

Mathematical programming, 45(1-3):503–528.

Liu, Q. (2017).

Stein variational gradient descent as gradient flow.

In *Advances in neural information processing systems*, pages 3115–3123.

References IX

Liu, Q., Lee, J., and Jordan, M. (2016).
 A kernelized stein discrepancy for goodness-of-fit tests.
 In *International conference on machine learning*, pages 276–284.

Liu, Q. and Wang, D. (2016).

Stein variational gradient descent: A general purpose bayesian inference algorithm.

In *Advances in neural information processing systems*, pages 2378–2386.

Liu, Y., Ramachandran, P., Liu, Q., and Peng, J. (2017). Stein variational policy gradient. *arXiv preprint arXiv:1704.02399*.

References X

Lu, J., Lu, Y., and Nolen, J. (2019). Scaling limit of the stein variational gradient descent: The mean field regime. SIAM Journal on Mathematical Analysis, 51(2):648–671. Nocedal, J. and Wright, S. (2006). Numerical optimization. Springer Science & Business Media. Oates, C. J., Girolami, M., and Chopin, N. (2017). Control functionals for monte carlo integration.

Journal of the Royal Statistical Society: Series B (Statistical Methodology), 79(3):695–718.

References XI

 Pu, Y., Gan, Z., Henao, R., Li, C., Han, S., and Carin, L. (2017).
 Vae learning via stein variational gradient descent.

In Advances in Neural Information Processing Systems, pages 4236–4245.

 Riabiz, M., Chen, W., Cockayne, J., Swietach, P., Niederer, S. A., Mackey, L., Oates, C., et al. (2020).
 Optimal thinning of MCMC output. arXiv preprint arXiv:2005.03952.

Steinwart, I. and Christmann, A. (2008). Support vector machines. Springer Science & Business Media.

References XII

- Wang, D. and Liu, Q. (2016). Learning to draw samples: With application to amortized mle for generative adversarial learning. arXiv preprint arXiv:1611.01722.
- Xu, W. and Matsuda, T. (2020).
 A Stein goodness-of-fit test for directional distributions.
 In Proceedings of the 23rd International Conference on Artificial Intelligence and Statistics, volume 108, pages 320–330. PMLR.
 - Zhang, C., Butepage, J., Kjellstrom, H., and Mandt, S. (2018).

Advances in variational inference.

IEEE transactions on pattern analysis and machine intelligence.

1 - Bayesian Logistic regression

Datapoints $d_1, \ldots, d_q \in \mathbb{R}^p$, and labels $y_1, \ldots, y_q \in \{\pm 1\}$.

Labels y_i are modelled as $p(y_i = 1 | d_i, w) = (1 + \exp(-w^\top d_i))^{-1}$ for some $w \in \mathbb{R}^p$.

The parameters *w* follow the law $p(w|\alpha) = \mathcal{N}(0, \alpha^{-1}I_p)$, and $\alpha > 0$ is drawn from an exponential law $p(\alpha) = \text{Exp}(0.01)$.

The parameter vector is then $x = [w, \log(\alpha)] \in \mathbb{R}^{p+1}$, and we use KSD-LBFGS to obtain samples from $p(x|(d_i, y_i)_{i=1}^q)$ for 13 datasets, with N = 10 particles for each.



Accuracy of the KSD descent and SVGD on bayesian logistic regression for 13 datasets.

Both methods yield similar results. KSD is better by 2% on one dataset.

2 - Bayesian Independent Component Analysis

ICA: $x = W^{-1}s$, where x is an observed sample in \mathbb{R}^{p} , $W \in \mathbb{R}^{p \times p}$ is the unknown square unmixing matrix, and $s \in \mathbb{R}^{p}$ are the independent sources.

1)Assume that each component has the same density $s_i \sim p_s$. 2) The likelihood of the model is $p(x|W) = \log |W| + \sum_{i=1}^{p} p_s([Wx]_i)$. 3)Prior: *W* has i.i.d. entries, of law $\mathcal{N}(0, 1)$.

The posterior is $p(W|x) \propto p(x|W)p(W)$, and the score is given by $s(W) = W^{-\top} - \psi(Wx)x^{\top} - W$, where $\psi = -\frac{p'_s}{p_s}$. In practice, we choose p_s such that $\psi(\cdot) = \tanh(\cdot)$. We then use the presented algorithms to draw 10 particles $W \sim p(W|x)$ on 50 experiments.



Left: p = 2. Middle: p = 4. Right: p = 8.

Each dot = Amari distance between an estimated matrix and the true unmixing matrix.

KSD Descent is not better than random. Explanation: ICA likelihood is highly non-convex.

(descent lemma for SVGD) Sketch of proof - 1

Fix $n \ge 0$. Denote $g = P_{\mu_n} \nabla \log(\frac{\mu_n}{\pi})$, $\phi_t = I - tg$ for $t \in [0, \gamma]$ and $\rho_t = (\phi_t)_{\#} \mu_n$. We have $\frac{\partial \rho_t}{\partial t} = div(\rho_t w_t)$ with $w_t = -g \circ \phi_t^{-1}$.

Denote $\varphi(t) = KL(\rho_t | \pi)$. Using a Taylor expansion,

$$\varphi(\gamma) = \varphi(0) + \gamma \varphi'(0) + \int_0^{\gamma} (\gamma - t) \varphi''(t) dt.$$

Step 1. $\varphi(0) = KL(\mu_n | \pi)$ and $\varphi(\gamma) = KL(\mu_{n+1} | \pi)$. Step 2. Using the chain rule,

$$\varphi'(t) = \langle \nabla_{W_2} \operatorname{KL}(\rho_t | \pi), W_t \rangle_{L^2(\rho_t)}$$

Hence :

$$arphi'(\mathbf{0}) = - \langle
abla \log\left(rac{\mu_n}{\pi}
ight), oldsymbol{g}
angle_{L^2(\mu_n)} = - \left\|oldsymbol{S}_{\mu_n}
abla \log\left(rac{\mu_n}{\pi}
ight)
ight\|_{\mathcal{H}}^2$$

Sketch of proof - 2 Step 3.

$$\varphi''(t) = \langle w_t, Hess_{\mathsf{KL}(.|\pi)}(\rho_t)w_t \rangle_{L^2(\rho_t)} := \psi_1(t) + \psi_2(t),$$

$$\psi_1(t) = \mathbb{E}_{x \sim \rho_t} \left[\langle w_t(x), H_V(x)w_t(x) \rangle \right] \text{ and } \psi_2(t) = \mathbb{E}_{x \sim \rho_t} \left[\|Jw_t(x)\|_{HS}^2 \right]$$

where $\rho_t = (\phi_t)_{\#} \mu_n, w_t = -g \circ (\phi_t)^{-1}.$

Sketch of proof - 2 Step 3.

$$\varphi''(t) = \langle w_t, Hess_{\mathsf{KL}(.|\pi)}(\rho_t)w_t \rangle_{L^2(\rho_t)} := \psi_1(t) + \psi_2(t),$$

$$\psi_1(t) = \mathbb{E}_{x \sim \rho_t} \left[\langle w_t(x), H_V(x)w_t(x) \rangle \right] \text{ and } \psi_2(t) = \mathbb{E}_{x \sim \rho_t} \left[\|Jw_t(x)\|_{HS}^2 \right]$$

where $\rho_t = (\phi_t)_{\#} \mu_n, w_t = -g \circ (\phi_t)^{-1}.$
Step 3.a. Assuming $\|H_V\| \leq M$ and $k(.,.) \leq B$:

$$\psi_1(t) \leq M \|g\|_{L^2(\mu_n)}^2 \leq MB^2 \left\|S_{\mu_n} \nabla \log\left(\frac{\mu_n}{\pi}\right)\right\|_{\mathcal{H}}^2$$

•

Sketch of proof - 2 Step 3.

$$\varphi''(t) = \langle w_t, Hess_{\mathsf{KL}(.|\pi)}(\rho_t)w_t \rangle_{L^2(\rho_t)} := \psi_1(t) + \psi_2(t),$$

$$\psi_1(t) = \mathbb{E}_{x \sim \rho_t} \left[\langle w_t(x), H_V(x)w_t(x) \rangle \right] \text{ and } \psi_2(t) = \mathbb{E}_{x \sim \rho_t} \left[\|Jw_t(x)\|_{HS}^2 \right]$$

where $\rho_t = (\phi_t)_{\#} \mu_n, w_t = -g \circ (\phi_t)^{-1}.$
Step 3.a. Assuming $\|H_V\| \leq M$ and $k(.,.) \leq B$:

$$\psi_1(t) \leq M \|g\|_{L^2(\mu_n)}^2 \leq MB^2 \left\|S_{\mu_n} \nabla \log\left(\frac{\mu_n}{\pi}\right)\right\|_{\mathcal{H}}^2$$

Step 3.b. Since $\rho_t = (\phi_t)_{\#} \mu_n$, $w_t = -g \circ (\phi_t)^{-1}$,

$$egin{aligned} \psi_2(t) &= \mathbb{E}_{x \sim \mu_n} [\|J w_t \circ \phi_t(x)\|_{HS}^2] \leq \|Jg(x)\|_{HS}^2 \|(J \phi_t)^{-1}(x)\|_{op}^2 \ &\leq B^2 \left\|S_{\mu_n}
abla \log\left(rac{\mu_n}{\pi}
ight)
ight\|_{\mathcal{H}}^2 lpha^2, \end{aligned}$$

assuming $\|\nabla k(.,.)\| \leq B$ and choosing $\gamma \leq f(\alpha)$ with $\alpha > 1$.

•

From:

$$\varphi(\gamma) = \varphi(0) + \gamma \varphi'(0) + \int_0^{\gamma} (\gamma - t) \varphi''(t) dt$$

we have:

$$egin{aligned} \mathsf{KL}(\mu_{n+1}|\pi) - \mathsf{KL}(\mu_n|\pi) &\leq -\gamma \|m{S}_{\mu_n}
abla \log\left(rac{\mu_n}{\pi}
ight)\|_{\mathcal{H}}^2 \ &+ rac{\gamma^2}{2} (lpha^2 + m{M}) m{B}^2 \|m{S}_{\mu_n}
abla \log\left(rac{\mu_n}{\pi}
ight)\|_{\mathcal{H}}^2. \end{aligned}$$

Choosing γ small enough yields a descent lemma :

$$\mathsf{KL}(\mu_{n+1}|\pi) - \mathsf{KL}(\mu_n|\pi) \leq -c_\gamma \underbrace{\left\| \mathcal{S}_{\mu_n}
abla \log\left(rac{\mu_n}{\pi}
ight)
ight\|_{\mathcal{H}}^2}_{\mathsf{KSD}^2(\mu_n|\pi)}.$$