# Limitations of the Theory of Sampling with Kernelized Wasserstein Gradient Flows

Anna Korba
CREST, ENSAE, Institut Polytechnique de Paris

ICBINB seminar series

# Outline

# Sampling

**Sampling problem:** Sample (=generate new examples) from a target distribution $\pi$ over $\mathbb{R}^d$, given some information on $\pi$.

# Sampling

> **Sampling problem:** Sample (=generate new examples) from a target distribution $\pi$ over $\mathbb{R}^d$, given some information on $\pi$.

Two different settings:

1. $\pi$'s density w.r.t. Lebesgue measure is known up to an intractable normalisation constant $Z$ :

$$\pi(x) = \frac{\tilde{\pi}(x)}{Z}, \quad \tilde{\pi} \text{ known, } Z \text{ unknown.}$$

   Example: Bayesian inference.

2. one has access to a set of samples of $\pi$ : $x_1, \ldots, x_n \sim \pi$.

   Example: (some) Neural networks, generative modelling (GANS...).

We'll focus on the first setting.

# Bayesian inference

Let $\mathcal{D} = (w_i, y_i)_{i=1}^m$ a dataset of labelled examples $(w_i, y_i) \overset{i.i.d.}{\sim} P_{\text{data}}$.
Assume an underlying model parametrized by $\theta$, e.g. :

$$y = g(w, \theta) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, I)$$

**Goal: learn the best distribution over $\theta$ to fit the data.**

# Bayesian inference

Let $\mathcal{D} = (w_i, y_i)_{i=1}^m$ a dataset of labelled examples $(w_i, y_i) \overset{i.i.d.}{\sim} P_{\text{data}}$.
Assume an underlying model parametrized by $\theta$, e.g. :

$$y = g(w, \theta) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, I)$$

**Goal: learn the best distribution over $\theta$ to fit the data.**

1. Compute the Likelihood:

$$p(\mathcal{D}|\theta) = \prod_{i=1}^m p(y_i|\theta, w_i) \propto \exp\left(-\frac{1}{2}\sum_{i=1}^m \|y_i - g(w_i, \theta)\|^2\right).$$

# Bayesian inference

Let $\mathcal{D} = (w_i, y_i)_{i=1}^{m}$ a dataset of labelled examples $(w_i, y_i) \overset{i.i.d.}{\sim} P_{data}$.
Assume an underlying model parametrized by $\theta$, e.g. :

$$y = g(w, \theta) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, I)$$

**Goal: learn the best distribution over $\theta$ to fit the data.**

1. Compute the Likelihood:

$$p(\mathcal{D}|\theta) = \prod_{i=1}^{m} p(y_i|\theta, w_i) \propto \exp\left(-\frac{1}{2}\sum_{i=1}^{m} \|y_i - g(w_i, \theta)\|^2\right).$$

2. Choose a prior distribution on the parameter:

$$\theta \sim p, \quad \text{e.g. } p(\theta) \propto \exp\left(-\frac{\|\theta\|^2}{2}\right).$$

# Bayesian inference

Let $\mathcal{D} = (w_i, y_i)_{i=1}^m$ a dataset of labelled examples $(w_i, y_i) \overset{i.i.d.}{\sim} P_{\text{data}}$.
Assume an underlying model parametrized by $\theta$, e.g. :

$$y = g(w, \theta) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, I)$$

**Goal: learn the best distribution over $\theta$ to fit the data.**

1. Compute the Likelihood:

$$p(\mathcal{D}|\theta) = \prod_{i=1}^m p(y_i|\theta, w_i) \propto \exp\left(-\frac{1}{2}\sum_{i=1}^m \|y_i - g(w_i, \theta)\|^2\right).$$

2. Choose a prior distribution on the parameter:

$$\theta \sim p, \quad \text{e.g. } p(\theta) \propto \exp\left(-\frac{\|\theta\|^2}{2}\right).$$

3. Bayes' rule yields:

$$\pi(\theta) := p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{Z} \quad Z = \int_{\mathbb{R}^d} p(\mathcal{D}|\theta)p(\theta)d\theta$$

i.e. $\pi(\theta) \propto \exp\left(-V(\theta)\right), \quad V(\theta) = \frac{1}{2}\sum_{i=1}^m \|y_i - g(w_i, \theta)\|^2 + \frac{\|\theta\|^2}{2}.$

$\pi$ is needed both for

▶ prediction for a new input *w*:

$$y_{pred} = \int_{\mathbb{R}^d} g(w, \theta) d\pi(\theta)$$

▶ measure uncertainty on the prediction.

$\pi$ is needed both for

- prediction for a new input $w$:

$$y_{pred} = \int_{\mathbb{R}^d} g(w, \theta) d\pi(\theta)$$

- measure uncertainty on the prediction.

Given a discrete approximation $\mu_n = \frac{1}{n} \sum\limits_{j=1}^{n} \delta_{\theta_j}$ of $\pi$:

$$y_{pred} \approx \frac{1}{n} \sum_{j=1}^{n} g(w, \theta_j).$$

**Question: how can we build $\mu_n$?**

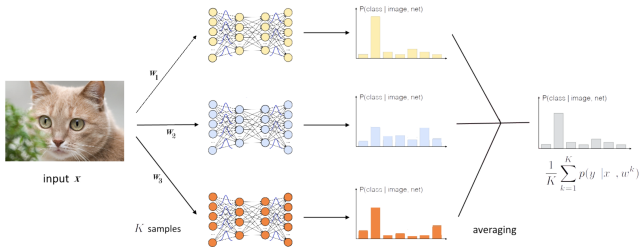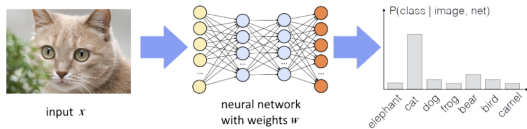Figure: Ensembling on deep neural networks.

# Sampling as optimisation

Notice that

$$\pi = \underset{\mu \in \mathcal{P}(\mathbb{R}^d)}{\operatorname{argmin}} \mathrm{KL}(\mu|\pi), \quad \mathrm{KL}(\mu|\pi) = \begin{cases} \int_{\mathbb{R}^d} \log\left(\frac{\mu}{\pi}(x)\right) d\mu(x) & \text{if } \mu \ll \pi \\ +\infty & \text{else.} \end{cases}$$

(does not depend on the normalisation constant $Z$ in $\pi(x) = \tilde{\pi}(x)/Z$ !)

# Sampling as optimisation

Notice that

$$\pi = \underset{\mu \in \mathcal{P}(\mathbb{R}^d)}{\operatorname{argmin}} \, \mathsf{KL}(\mu|\pi), \quad \mathsf{KL}(\mu|\pi) = \begin{cases} \int_{\mathbb{R}^d} \log\left(\frac{\mu}{\pi}(x)\right) d\mu(x) & \text{if } \mu \ll \pi \\ +\infty & \text{else.} \end{cases}$$

(does not depend on the normalisation constant $Z$ in $\pi(x) = \tilde{\pi}(x)/Z$ !)

Two ways to produce an approximation $\mu_n$:

1. Markov Chain Monte Carlo (MCMC) methods: generate a Markov chain whose law converges to $\pi \propto \exp(-V)$

   Example: discretize an overdamped Langevin diffusion

   $$d\theta_t = -\nabla V(\theta_t) + \sqrt{2}dB_t \Longrightarrow \theta_{l+1} = \theta_l - \gamma\nabla V(\theta_l) + \sqrt{2\gamma}\epsilon_l, \; \epsilon_l \sim \mathcal{N}(0, I_d)$$

   Its law corresponds to a Wasserstein gradient flow of the KL
   [Jordan et al., 1998].

# Sampling as optimisation

Notice that

$$\pi = \underset{\mu \in \mathcal{P}(\mathbb{R}^d)}{\operatorname{argmin}} \, \mathsf{KL}(\mu|\pi), \quad \mathsf{KL}(\mu|\pi) = \left\{ \begin{array}{ll} \int_{\mathbb{R}^d} \log\left(\frac{\mu}{\pi}(x)\right) d\mu(x) & \text{if } \mu \ll \pi \\ +\infty & \text{else.} \end{array} \right.$$

(does not depend on the normalisation constant $Z$ in $\pi(x) = \tilde{\pi}(x)/Z$ !)

Two ways to produce an approximation $\mu_n$:

1. Markov Chain Monte Carlo (MCMC) methods: generate a Markov chain whose law converges to $\pi \propto \exp(-V)$

   Example: discretize an overdamped Langevin diffusion

   $$d\theta_t = -\nabla V(\theta_t) + \sqrt{2}dB_t \Longrightarrow \theta_{l+1} = \theta_l - \gamma\nabla V(\theta_l) + \sqrt{2\gamma}\epsilon_l, \ \epsilon_l \sim \mathcal{N}(0, I_d)$$

   Its law corresponds to a Wasserstein gradient flow of the KL
   [Jordan et al., 1998].
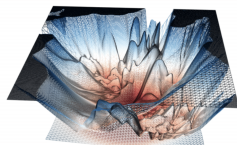
2. Interacting particle systems, e.g. by considering other metrics or functionals

# Difficult cases (in practice and in theory)

Recall that

$$\pi(\theta) \propto \exp\left(-V(\theta)\right), \quad V(\theta) = \underbrace{\sum_{i=1}^{m} \|y_i - g(w_i, \theta)\|^2}_{\text{loss}} + \frac{\|\theta\|^2}{2}.$$

- ▶ if $V$ is convex (e.g. $g(w, \theta) = \langle w, \theta \rangle$) many sampling methods are known to work quite well
- ▶ but if its not (e.g. $g(w, \theta)$ is a neural network), the situation is much more delicate



A highly nonconvex loss surface, as is common in deep neural nets.
From https://www.telesens.co/2019/01/16/
neural-network-loss-visualization.

# 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:

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

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 $\mu_0$ to $\pi$.

# Outline

# Euclidean gradient flow and continuity equation

Let $V : \mathbb{R}^d \to \mathbb{R}$. Consider the gradient flow

$$x'(t) = -\nabla V(x(t))$$

and assume $x(0)$ random with density $\mu_0$. What is the dynamics of the density $\mu_t$ of $x(t)$ ? Let $\phi : \mathbb{R}^d \to \mathbb{R}$ a test function.

$$\frac{d}{dt} \mathbb{E}(\phi(x(t))) = - \int \langle \nabla \phi, \nabla V \rangle \mu_t(x) dx = \int \phi(x) \boldsymbol{\nabla} \cdot (\mu_t \nabla V)(x) dx,$$

and

$$\frac{d}{dt} \mathbb{E}(\phi(x(t))) = \int \phi(x) \frac{\partial \mu_t}{\partial t}(x) dx.$$

Therefore,

$$\frac{\partial \mu_t}{\partial t} = \boldsymbol{\nabla} \cdot (\mu_t \nabla V).$$

# Setting - The Wasserstein space

Let $\mathcal{P}_2(\mathbb{R}^d)$ denote the space of probability measures on $\mathbb{R}^d$ with finite second moments, i.e.

$$\mathcal{P}_2(\mathbb{R}^d) = \{\mu \in \mathcal{P}(\mathbb{R}^d), \ \int \|x\|^2 d\mu(x) < \infty\}$$

# Setting - The Wasserstein space

Let $\mathcal{P}_2(\mathbb{R}^d)$ denote the space of probability measures on $\mathbb{R}^d$ with finite second moments, i.e.

$$\mathcal{P}_2(\mathbb{R}^d) = \{\mu \in \mathcal{P}(\mathbb{R}^d), \ \int \|x\|^2 d\mu(x) < \infty\}$$

# Setting - The Wasserstein space

Let $\mathcal{P}_2(\mathbb{R}^d)$ denote the space of probability measures on $\mathbb{R}^d$ with finite second moments, i.e.

$$\mathcal{P}_2(\mathbb{R}^d) = \{\mu \in \mathcal{P}(\mathbb{R}^d), \ \int \|x\|^2 d\mu(x) < \infty\}$$

$\mathcal{P}_2(\mathbb{R}^d)$ is endowed with the Wasserstein-2 distance from **Optimal transport** :

$$W_2^2(\nu, \mu) = \inf_{s \in \Gamma(\nu, \mu)} \int_{\mathbb{R}^d \times \mathbb{R}^d} \|x - y\|^2 ds(x, y) \qquad \forall \nu, \mu \in \mathcal{P}_2(\mathbb{R}^d)$$

where $\Gamma(\nu, \mu)$ is the set of possible couplings between $\nu$ and $\mu$ (joint distributions on $\mathbb{R}^d \times \mathbb{R}^d$ with first marginals $\nu$ and $\mu$).

## Setting - The Wasserstein space

Let $\mathcal{P}_2(\mathbb{R}^d)$ denote the space of probability measures on $\mathbb{R}^d$ with finite second moments, i.e.

$$\mathcal{P}_2(\mathbb{R}^d) = \{\mu \in \mathcal{P}(\mathbb{R}^d), \, \int \|x\|^2 d\mu(x) < \infty\}$$

$\mathcal{P}_2(\mathbb{R}^d)$ is endowed with the Wasserstein-2 distance from **Optimal transport** :

$$W_2^2(\nu, \mu) = \inf_{s \in \Gamma(\nu, \mu)} \int_{\mathbb{R}^d \times \mathbb{R}^d} \|x - y\|^2 \, ds(x, y) \qquad \forall \nu, \mu \in \mathcal{P}_2(\mathbb{R}^d)$$

where $\Gamma(\nu, \mu)$ is the set of possible couplings between $\nu$ and $\mu$ (joint distributions on $\mathbb{R}^d \times \mathbb{R}^d$ with first marginals $\nu$ and $\mu$).

Can also be written:

$$W_2^2(\nu, \mu) = \inf_{(\rho_t, v_t)_{t \in [0,1]}} \left\{ \int_0^1 \|v_t(x)\|_{L^2(\rho_t)}^2 dt(x) : \frac{\partial \rho_t}{\partial t} = \boldsymbol{\nabla} \cdot (\rho_t v_t), \rho_0 = \nu, \rho_1 = \mu \right\}$$

**Definition :** Let $\mu \in \mathcal{P}_2(\mathbb{R}^d)$, $T : \mathbb{R}^d \to \mathbb{R}^d$. The pushforward measure $T_{\#}\mu$ is characterized by:

- $\forall$ B meas. set, $T_{\#}\mu(B) = \mu(T^{-1}(B))$
- $x \sim \mu$, $T(x) \sim T_{\#}\mu$

**Definition :** Let $\mu \in \mathcal{P}_2(\mathbb{R}^d)$, $T : \mathbb{R}^d \to \mathbb{R}^d$. The pushforward measure $T_{\#}\mu$ is characterized by:

- $\forall$ B meas. set, $T_{\#}\mu(B) = \mu(T^{-1}(B))$
- $x \sim \mu$, $T(x) \sim T_{\#}\mu$

**(Brenier's theorem):** Let $\mu, \nu \in \mathcal{P}_2(\mathbb{R}^d)$ s.t. $\mu \ll$ *Leb*. Then, there exists $T_{\mu}^{\nu} : \mathbb{R}^d \to \mathbb{R}^d$ such that
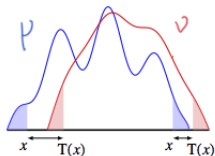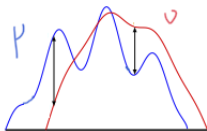
- $T_{\mu\#}^{\nu}\mu = \nu$
- $W_2^2(\mu, \nu) = \|I - T_{\mu}^{\nu}\|_{L_2(\mu)}^2 = \int \|x - T_{\mu}^{\nu}(x)\|^2 d\mu(x)$

$W_2$ geodesics?

$\rho(0) = \mu, \rho(1) = \nu$.

$\rho(t) = ((1-t)I + tT_{\mu}^{\nu})_{\#}\mu$

$\neq \rho(t) = \underbrace{(1-t)\mu + t\nu}_{\text{mixture}}$

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

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

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

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

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

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

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

$$\frac{\partial \mu_t}{\partial t} = \boldsymbol{\nabla} \cdot \left( \mu_t \nabla_{W_2} \mathcal{F}(\mu_t) \right),$$

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

# WGF of Free energies

In particular, if the functional $\mathcal{F}$ is a free energy:

$$\mathcal{F}(\mu) = \underbrace{\int H(\mu(x))dx}_{\text{internal energy}} + \underbrace{\int V(x)d\mu(x)}_{\text{potential energy}} + \underbrace{\int W(x,y)d\mu(x)d\mu(y)}_{\text{interaction energy}}$$

$$\text{Then}: \ \frac{\partial \mu_t}{\partial t} = \boldsymbol{\nabla} \cdot \left( \mu_t \underbrace{\nabla(H'(\mu_t) + V + W * \mu_t)}_{\nabla_{W_2}\mathcal{F}(\mu)} \right). \quad (1)$$

For instance, if $H = 0$ then (1) rules the density $\mu_t$ of particles $x_t \in \mathbb{R}^d$ driven by :

$$\frac{dx_t}{dt} = -\nabla V(x_t) - \int_{\mathbb{R}^d} \nabla W(x, x_t)d\mu_t(x)$$

$\mu_t = Law(x_t)$.

# (Some) unbiased time discretizations

For a step-size $\gamma > 0$:

1. Backward (expensive) [Jordan et al., 1998] :

$$\mu_{l+1} = \text{JKO}_{\gamma \mathcal{F}}(\mu_l)$$

where $\text{JKO}_{\gamma \mathcal{F}}(\mu_l) = \underset{\mu \in \mathcal{P}_2(\mathbb{R}^d)}{\text{argmin}} \left\{ \mathcal{F}(\mu) + \frac{1}{2\gamma} W_2^2(\mu, \mu_l) \right\}.$

## (Some) unbiased time discretizations

For a step-size $\gamma > 0$:

1. Backward (expensive) [Jordan et al., 1998] :

$$\mu_{l+1} = \text{JKO}_{\gamma \mathcal{F}}(\mu_l)$$

where $\text{JKO}_{\gamma \mathcal{F}}(\mu_l) = \underset{\mu \in \mathcal{P}_2(\mathbb{R}^d)}{\text{argmin}} \left\{ \mathcal{F}(\mu) + \frac{1}{2\gamma} W_2^2(\mu, \mu_l) \right\}.$

2. Forward (cheap) :

$$\mu_{l+1} = exp_{\mu_l}(-\gamma \nabla_{W_2} \mathcal{F}(\mu_l)) = \left( I - \gamma \nabla_{W_2} \mathcal{F}(\mu_l) \right)_{\#} \mu_l$$

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

$$X_{l+1} = X_l - \gamma \nabla_{W_2} \mathcal{F}(\mu_l)(X_l) \sim \mu_{l+1}, \text{ if } X_l \sim \mu_l.$$

# Space discretization - Interacting particle system

**Problem:** the vector field depends on the **unknown** $\mu_l$, the density of the particle at time *l*.

# Space discretization - Interacting particle system

**Problem:** the vector field depends on the **unknown** $\mu_l$, the density of the particle at time $l$.

**Idea:** replace it by the **empirical measure** of a system of $n$ interacting particles:

$$X_0^1, \ldots, X_0^n \sim \mu_0$$

and for $j = 1, \ldots, n$:

$$\begin{aligned}
X_{l+1}^j &= X_l^j - \gamma \nabla_{W_2} \mathcal{F}(\hat{\mu}_l)(X_l^j) \\
&= X_l^j - \frac{1}{\gamma} \left[ \nabla V(X_l^j) + \frac{1}{n} \sum_{i=1}^n \nabla W(X_l^j, X_l^j) \right]
\end{aligned}$$

where $\hat{\mu}_l = \frac{1}{n} \sum_{i=1}^n \delta_{X_l^j}$.

# Outline

**Goal:** Sample from a target distribution $\pi$, whose density w.r.t. Lebesgue measure is known up to an intractable normalisation constant $Z$ :

$$\pi(x) = \frac{\tilde{\pi}(x)}{Z}, \quad \tilde{\pi} \text{ known, } Z \text{ unknown.}$$

Remember that

$$\pi = \operatorname{argmin} \mathsf{KL}(\mu|\pi), \quad \mathsf{KL}(\mu|\pi) = \int \log\left(\frac{\mu}{\pi}\right) d\mu \text{ if } \mu \ll \pi$$

and that we can consider the Forward time discretisation:

$$x_{l+1} = x_l - \gamma \nabla_{W_2} \mathsf{KL}(\mu_l|\pi)(x_l), \quad x_l \sim \mu_l,$$

where $\nabla_{W_2} \mathsf{KL}(\mu_l|\pi) = \nabla \frac{\partial \mathsf{KL}(\mu_l|\pi)}{\partial \mu} = \nabla \log\left(\frac{\mu_l}{\pi}(.)\right)$.

**Problem:** $\mu_l$, hence $\nabla \log(\mu_l)$ is unknown and has to be estimated from a set of particles.

# 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
  ($(k(x_i, x_j)_{i=1}^n)$ is a p.s.d. matrix for all $x_1, \ldots, x_n \in \mathbb{R}^d$)

# Background on kernels and RKHS

- ▶ Let $k : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$ a positive, semi-definite kernel
    $((k(x_i, x_j)_{i=1}^n)$ is a p.s.d. matrix for all $x_1, \ldots, x_n \in \mathbb{R}^d)$
- ▶ examples:
    - ▶ the Gaussian kernel $k(x, y) = \exp\left(-\frac{\|x-y\|^2}{h}\right)$
    - ▶ the Laplace kernel $k(x, y) = \exp\left(-\frac{\|x-y\|}{h}\right)$
    - ▶ the inverse multiquadratic kernel
      $k(x, y) = (c + \|x - y\|)^{-\beta}$ with $\beta \in\ ]0, 1[$

# 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
  ($(k(x_i, x_j)_{i=1}^n)$ is a p.s.d. matrix for all $x_1, \ldots, x_n \in \mathbb{R}^d$)

- examples:

  - the Gaussian kernel $k(x, y) = \exp\left(-\frac{\|x-y\|^2}{h}\right)$
  - the Laplace kernel $k(x, y) = \exp\left(-\frac{\|x-y\|}{h}\right)$
  - the inverse multiquadratic kernel
    $k(x, y) = (c + \|x - y\|)^{-\beta}$ with $\beta \in ]0, 1[$

- $\mathcal{H}_k$ its corresponding RKHS (Reproducing Kernel Hilbert Space):

$$\mathcal{H}_k = \overline{\left\{ \sum_{i=1}^m \alpha_i k(\cdot, x_i); \; m \in \mathbb{N}; \; \alpha_1, \ldots, \alpha_m \in \mathbb{R}; \; x_1, \ldots, x_m \in \mathbb{R}^d \right\}}$$

# 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
  $((k(x_i, x_j)_{i=1}^n)$ is a p.s.d. matrix for all $x_1, \ldots, x_n \in \mathbb{R}^d)$

▶ examples:

  ▶ the Gaussian kernel $k(x, y) = \exp\left(-\frac{\|x-y\|^2}{h}\right)$

  ▶ the Laplace kernel $k(x, y) = \exp\left(-\frac{\|x-y\|}{h}\right)$

  ▶ the inverse multiquadratic kernel
     $k(x, y) = (c + \|x - y\|)^{-\beta}$ with $\beta \in ]0, 1[$

▶ $\mathcal{H}_k$ its corresponding RKHS (Reproducing Kernel Hilbert Space):

$$\mathcal{H}_k = \overline{\left\{\sum_{i=1}^m \alpha_i k(\cdot, x_i); \ m \in \mathbb{N}; \ \alpha_1, \ldots, \alpha_m \in \mathbb{R}; \ x_1, \ldots, x_m \in \mathbb{R}^d\right\}}$$

▶ $\mathcal{H}_k$ is a Hilbert space with inner product $\langle ., . \rangle_{\mathcal{H}_k}$ and norm $\|.\|_{\mathcal{H}_k}$.

# 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
  $((k(x_i, x_j)_{i=1}^n)$ is a p.s.d. matrix for all $x_1, \ldots, x_n \in \mathbb{R}^d)$
- ▶ examples:
  - ▶ the Gaussian kernel $k(x, y) = \exp\left(-\frac{\|x-y\|^2}{h}\right)$
  - ▶ the Laplace kernel $k(x, y) = \exp\left(-\frac{\|x-y\|}{h}\right)$
  - ▶ the inverse multiquadratic kernel
    $k(x, y) = (c + \|x - y\|)^{-\beta}$ with $\beta \in ]0, 1[$
- ▶ $\mathcal{H}_k$ its corresponding RKHS (Reproducing Kernel Hilbert Space):

$$\mathcal{H}_k = \overline{\left\{ \sum_{i=1}^m \alpha_i k(\cdot, x_i); \ m \in \mathbb{N}; \ \alpha_1, \ldots, \alpha_m \in \mathbb{R}; \ x_1, \ldots, x_m \in \mathbb{R}^d \right\}}$$

- ▶ $\mathcal{H}_k$ is a Hilbert space with inner product $\langle ., . \rangle_{\mathcal{H}_k}$ and norm $\|.\|_{\mathcal{H}_k}$.
- ▶ assume $\int_{\mathbb{R}^d \times \mathbb{R}^d} k(x, x) d\mu(x) < \infty$ for any $\mu \in \mathcal{P}(\mathbb{R}^d), \Longrightarrow \mathcal{H}_k \subset L^2(\mu)$.

# 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
$((k(x_i, x_j)_{i=1}^n)$ is a p.s.d. matrix for all $x_1, \ldots, x_n \in \mathbb{R}^d)$

▶ examples:

    ▶ the Gaussian kernel $k(x, y) = \exp\left(-\frac{\|x-y\|^2}{h}\right)$

    ▶ the Laplace kernel $k(x, y) = \exp\left(-\frac{\|x-y\|}{h}\right)$

    ▶ the inverse multiquadratic kernel
    $k(x, y) = (c + \|x - y\|)^{-\beta}$ with $\beta \in ]0, 1[$

▶ $\mathcal{H}_k$ its corresponding RKHS (Reproducing Kernel Hilbert Space):

$$\mathcal{H}_k = \overline{\left\{\sum_{i=1}^m \alpha_i k(\cdot, x_i); \ m \in \mathbb{N}; \ \alpha_1, \ldots, \alpha_m \in \mathbb{R}; \ x_1, \ldots, x_m \in \mathbb{R}^d\right\}}$$

▶ $\mathcal{H}_k$ is a Hilbert space with inner product $\langle ., . \rangle_{\mathcal{H}_k}$ and norm $\|.\|_{\mathcal{H}_k}$.

▶ assume $\int_{\mathbb{R}^d \times \mathbb{R}^d} k(x, x) d\mu(x) < \infty$ for any $\mu \in \mathcal{P}(\mathbb{R}^d), \Longrightarrow \mathcal{H}_k \subset L^2(\mu)$.

▶ It satisfies the reproducing property:

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

# Stein Variational Gradient Descent [Liu and Wang, 2016]

Consider the following metric depending on *k*

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

Then, the $W_k$ gradient flow of the KL writes as the PDE
[Liu, 2017], [Duncan et al., 2019]:

$$\frac{\partial \mu_t}{\partial t} + \boldsymbol{\nabla} \cdot \left( \mu_t P_{\mu_t} \nabla \log \left( \frac{\mu_t}{\pi} \right) \right) = 0, \quad P_\mu : f \mapsto \int k(x, .) f(x) d\mu(x).$$

It converges to $\pi \propto \exp(-V)$ under mild conditions on *k* and if
*V* grows at most polynomially [Lu et al., 2019].

# SVGD algorithm

**SVGD trick:** applying the kernel integral operator to the $W_2$ gradient of $\mathrm{KL}(\cdot|\pi)$ leads to

$$
\begin{aligned}
P_\mu \nabla \log\left(\frac{\mu}{\pi}\right)(\cdot) &= \int \nabla \log\left(\frac{\mu}{\pi}\right)(x) k(x,.) d\mu(x) \\
&= \int -\nabla \log(\pi(x)) k(x,.) d\mu(x) + \int \nabla(\mu(x)) k(x,.) dx \\
&\overset{I.P.P.}{=} -\int [\nabla \log \pi(x) k(x,\cdot) + \nabla_x k(x,\cdot)] d\mu(x),
\end{aligned}
$$

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

# SVGD algorithm

**SVGD trick:** applying the kernel integral operator to the $W_2$ gradient of $\text{KL}(\cdot|\pi)$ leads to

$$
\begin{aligned}
P_\mu \nabla \log \left( \frac{\mu}{\pi} \right) (\cdot) &= \int \nabla \log \left( \frac{\mu}{\pi} \right)(x) k(x, .) d\mu(x) \\
&= \int -\nabla \log(\pi(x)) k(x, .) d\mu(x) + \int \nabla(\mu(x)) k(x, .) dx \\
&\overset{I.P.P.}{=} -\int [\nabla \log \pi(x) k(x, \cdot) + \nabla_x k(x, \cdot)] d\mu(x),
\end{aligned}
$$

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

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

$$
\begin{aligned}
X_{l+1}^i &= X_l^i - \gamma \left[ \frac{1}{n} \sum_{j=1}^n k(X_l^i, X_l^j) \nabla_{X_l^j} \log \pi(X_l^j) + \nabla_{X_l^j} k(X_l^i, X_l^j) \right] \\
&= X_l^i - \gamma P_{\mu_l^n} \nabla \log \left( \frac{\mu_l^n}{\pi} \right) (X_l^i), \quad \text{with } \mu_l^n = \frac{1}{n} \sum_{j=1}^n \delta_{X_l^j}
\end{aligned}
$$

# SVGD algorithm

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

$$
\begin{aligned}
P_\mu \nabla \log \left( \frac{\mu}{\pi} \right) (\cdot) &= \int \nabla \log \left( \frac{\mu}{\pi} \right) (x) k(x,.) d\mu(x) \\
&= \int -\nabla \log(\pi(x)) k(x,.) d\mu(x) + \int \nabla(\mu(x)) k(x,.) dx \\
&\overset{I.P.P.}{=} -\int [\nabla \log \pi(x) k(x,\cdot) + \nabla_x k(x,\cdot)] d\mu(x),
\end{aligned}
$$

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

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

$$
\begin{aligned}
X_{l+1}^i &= X_l^i - \gamma \left[ \frac{1}{n} \sum_{j=1}^n k(X_l^i, X_l^j) \nabla_{X_l^j} \log \pi(X_l^j) + \nabla_{X_l^j} k(X_l^i, X_l^j) \right] \\
&= X_l^i - \gamma P_{\mu_l^n} \nabla \log \left( \frac{\mu_l^n}{\pi} \right) (X_l^i), \quad \text{with } \mu_l^n = \frac{1}{n} \sum_{j=1}^n \delta_{X_l^j}
\end{aligned}
$$

# SVGD in practice

- ▶ more than 600 citations for [Liu and Wang, 2016]
- ▶ Relative empirical success in Bayesian inference and more recently deep ensembles
- ▶ It can suffer for multimodal distributions [Wenliang and Kanagawa, 2020], underestimate the target variance [Ba et al., 2021], but still can be very efficient on difficult sampling problems.

| | | AUROC(H) | AUROC(MD) | Accuracy | $H_o/H_t$ | $MD_o/MD_t$ | ECE | NLL |
|---|---|---|---|---|---|---|---|---|
| FashionMNIST | Deep ensemble [38] | 0.958±0.001 | 0.975±0.001 | 91.122±0.013 | 6.257±0.005 | 6.394±0.001 | **0.012±0.001** | 0.129±0.001 |
| | SVGD [46] | 0.960±0.001 | 0.973±0.001 | 91.134±0.024 | 6.315±0.019 | 6.395±0.018 | 0.014±0.001 | 0.127±0.001 |
| | f-SVGD [67] | 0.956±0.001 | 0.975±0.001 | 89.884±0.015 | 5.652±0.009 | 6.531±0.005 | 0.013±0.001 | 0.150±0.001 |
| | kde-WGD (ours) | 0.960±0.001 | 0.970±0.001 | 91.238±0.019 | 6.587±0.019 | 6.379±0.018 | 0.014±0.001 | 0.128±0.001 |
| | sge-WGD (ours) | 0.960±0.001 | 0.970±0.001 | **91.312±0.016** | 6.562±0.007 | 6.363±0.009 | **0.012±0.001** | 0.128±0.001 |
| | ssge-WGD (ours) | 0.968±0.001 | 0.979±0.001 | 91.198±0.024 | 6.522±0.009 | 6.610±0.012 | **0.012±0.001** | 0.130±0.001 |
| | kde-fWGD (ours) | **0.971±0.001** | **0.980±0.001** | 91.260±0.011 | 7.079±0.016 | 6.887±0.015 | 0.015±0.001 | **0.125±0.001** |
| | sge-fWGD (ours) | 0.969±0.001 | 0.978±0.001 | 91.192±0.013 | 7.076±0.004 | 6.900±0.005 | 0.015±0.001 | **0.125±0.001** |
| | ssge-fWGD (ours) | **0.971±0.001** | **0.980±0.001** | 91.240±0.022 | **7.129±0.006** | **6.951±0.005** | 0.016±0.001 | 0.124±0.001 |
| CIFAR10 | Deep ensemble [38] | **0.843±0.004** | 0.736±0.005 | 85.552±0.076 | **2.244±0.006** | 1.667±0.008 | 0.049±0.001 | 0.277±0.001 |
| | SVGD [46] | 0.825±0.003 | 0.710±0.002 | 85.142±0.017 | 2.106±0.003 | 1.567±0.004 | 0.052±0.001 | 0.287±0.001 |
| | fSVGD [67] | 0.783±0.001 | 0.712±0.001 | 84.510±0.031 | 1.968±0.004 | 1.624±0.003 | 0.049±0.001 | 0.292±0.001 |
| | kde-WGD (ours) | 0.838±0.001 | 0.735±0.004 | **85.904±0.030** | 2.205±0.003 | 1.661±0.008 | 0.053±0.001 | **0.276±0.001** |
| | sge-WGD (ours) | 0.837±0.003 | 0.725±0.004 | 85.792±0.035 | 2.214±0.010 | 1.634±0.004 | 0.051±0.001 | **0.275±0.001** |
| | ssge-WGD (ours) | 0.832±0.003 | 0.731±0.005 | 85.638±0.038 | 2.182±0.015 | 1.655±0.001 | 0.049±0.001 | **0.276±0.001** |
| | kde-fWGD (ours) | 0.791±0.002 | **0.758±0.002** | 84.888±0.030 | 1.970±0.004 | **1.749±0.005** | **0.044±0.001** | 0.282±0.001 |
| | sge-fWGD (ours) | 0.795±0.001 | 0.754±0.002 | 84.766±0.060 | 1.984±0.003 | 1.729±0.002 | 0.047±0.001 | 0.288±0.001 |
| | ssge-fWGD (ours) | 0.792±0.002 | 0.752±0.005 | 84.762±0.034 | 1.970±0.006 | 1.723±0.005 | 0.046±0.001 | 0.286±0.001 |

# Continuous-time dynamics of SVGD

$$\frac{\partial \mu_t}{\partial t} + \boldsymbol{\nabla} \cdot \left( \mu_t P_{\mu_t} \nabla \log \left( \frac{\mu_t}{\pi} \right) \right) = 0, \quad P_\mu : f \mapsto \int k(x, .) f(x) d\mu(x).$$

# Continuous-time dynamics of SVGD

$$\frac{\partial \mu_t}{\partial t} + \boldsymbol{\nabla} \cdot \left( \mu_t P_{\mu_t} \nabla \log \left( \frac{\mu_t}{\pi} \right) \right) = 0, \quad P_\mu : f \mapsto \int k(x, .) f(x) d\mu(x).$$

How fast the KL decreases along SVGD dynamics? Apply the chain rule in the Wasserstein space:

$$\frac{d \, \mathsf{KL}(\mu_t | \pi)}{dt} = \left\langle V_t, \nabla \log \left( \frac{\mu_t}{\pi} \right) \right\rangle_{L^2(\mu_t)} = - \underbrace{\left\| P_{\mu_t} \nabla \log \left( \frac{\mu_t}{\pi} \right) \right\|^2_{\mathcal{H}_k}}_{\mathsf{KSD}^2(\mu_t | \pi)} \leq 0.$$

# Continuous-time dynamics of SVGD

$$\frac{\partial \mu_t}{\partial t} + \boldsymbol{\nabla} \cdot \left( \mu_t P_{\mu_t} \nabla \log \left( \frac{\mu_t}{\pi} \right) \right) = 0, \quad P_\mu : f \mapsto \int k(x,.)f(x)d\mu(x).$$

How fast the KL decreases along SVGD dynamics? Apply the chain rule in the Wasserstein space:

$$\frac{d\, \mathsf{KL}(\mu_t|\pi)}{dt} = \left\langle V_t, \nabla \log \left( \frac{\mu_t}{\pi} \right) \right\rangle_{L^2(\mu_t)} = - \underbrace{\left\| P_{\mu_t} \nabla \log \left( \frac{\mu_t}{\pi} \right) \right\|_{\mathcal{H}_k}^2}_{\mathsf{KSD}^2(\mu_t|\pi)} \leq 0.$$

On the r.h.s. we have the squared **Kernel Stein discrepancy (KSD)**
[Chwialkowski et al., 2016] or **Stein Fisher information** of $\mu_t$ relative to $\pi$:

$$\left\| P_{\mu,k} \nabla \log \left( \frac{\mu}{\pi} \right) \right\|_{\mathcal{H}_k}^2 = \langle P_{\mu,k} \nabla \log \left( \frac{\mu}{\pi} \right), P_{\mu,k} \nabla \log \left( \frac{\mu}{\pi} \right) \rangle_{\mathcal{H}_k}$$

$$= \iint \nabla \log \left( \frac{\mu}{\pi}(x) \right) \nabla \log \left( \frac{\mu}{\pi}(y) \right) k(x,y)d\mu(x)d\mu(y).$$

Recall that the Fisher divergence is defined as $\| \nabla \log \left( \frac{\mu}{\pi} \right) \|_{L^2(\mu)}^2$.

# Exponential decay?

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

$$\text{KL}(\mu|\pi) \leq \frac{1}{2\lambda} \text{KSD}^2(\mu|\pi).$$

# Exponential decay?

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

$$\text{KL}(\mu|\pi) \leq \frac{1}{2\lambda} \text{KSD}^2(\mu|\pi).$$

If it holds, we can conclude with Gronwall's lemma:

$$\frac{d \text{KL}(\mu_t|\pi)}{dt} = - \text{KSD}^2(\mu_t|\pi) \leq -2\lambda \text{KL}(\mu_t|\pi) \implies \text{KL}(\mu_t|\pi) \leq e^{-2\lambda t} \text{KL}(\mu_0|\pi).$$

**When is Stein log-Sobolev satisfied?** not so well understood

[Duncan et al., 2019]:

- ▶ it fails to hold if $k$ is too regular with respect to $\pi$ (e.g. $k$ bounded, $\pi$ Gaussian)
- ▶ some working examples in dimension 1, open question in greater dimensions...

# A descent lemma in discrete time for SVGD [Korba et al., 2020]

**Idea:** in optimisation, descent lemmas can be shown if the objective function has a bounded Hessian.

**Idea:** in optimisation, descent lemmas can be shown if the objective function has a bounded Hessian.

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

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

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

# A descent lemma in discrete time for SVGD [Korba et al., 2020]

**Idea:** in optimisation, descent lemmas can be shown if the objective function has a bounded Hessian.

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

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

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

**However:** In the case of SVGD, the descent directions $f$ are restricted to $\mathcal{H}_k$ (bounded functions for bounded $k$).

**Proposition:** Assume (boundedness of $k$ and $\nabla k$, $H_V$ and moments on the trajectory), then for $\gamma$ small enough:

$$KL(\mu_{l+1}|\pi) - KL(\mu_l|\pi) \leq -c_\gamma \underbrace{\left\| P_{\mu_l} \nabla \log \left( \frac{\mu_l}{\pi} \right) \right\|_{\mathcal{H}_k}^2}_{KSD^2(\mu_l|\pi)}.$$

# Rates in KSD

**Consequence of the descent lemma:** for $\gamma$ small enough,

$$\min_{l=1,\ldots,L} \mathsf{KSD}^2(\mu_l|\pi) \leq \frac{1}{L} \sum_{l=1}^{L} \mathsf{KSD}^2(\mu_l|\pi) \leq \frac{\mathsf{KL}(\mu_0|\pi)}{c_\gamma L}.$$

This result does not rely on:

▶ convexity of *V*

▶ nor on Stein log Sobolev inequality

▶ only on smoothness of *V*.

in contrast with many convergence results on LMC.

The KSD metrizes convergence for instance when
[Gorham and Mackey, 2017]:

▶ $\pi$ is distantly dissipative (log concave at infinity, e.g. mixture of Gaussians)

▶ *k* is the IMQ kernel defined by $k(x, y) = (c^2 + \|x - y\|_2^2)^\beta$ for $c > 0$ and $\beta \in (-1, 0)$.

# Open question 1: Rates in terms of the KL objective?

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

$$\text{KL}(\mu_{l+1}|\pi) - \text{KL}(\mu_l|\pi) \leq -c_\gamma \left\| S_{\mu_n} \nabla \log \left( \frac{\mu_l}{\pi} \right) \right\|_{\mathcal{H}_k}^2$$

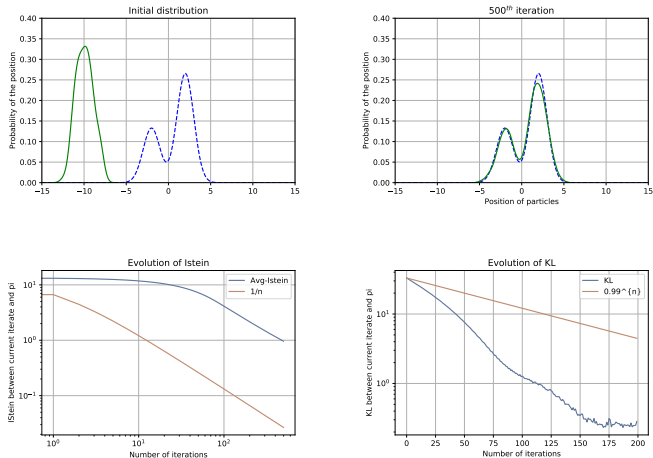and the Stein log-Sobolev inequality (2) with constant $\lambda$:

$$\text{KL}(\mu_{l+1}|\pi) - \text{KL}(\mu_l|\pi) \underbrace{\leq}_{(1)} -c_\gamma \left\| P_{\mu_l} \nabla \log \left( \frac{\mu_n}{\pi} \right) \right\|_{\mathcal{H}_k}^2 \underbrace{\leq}_{(2)} -c_\gamma 2\lambda \, \text{KL}(\mu_n|\pi).$$

Iterating this inequality yields $\text{KL}(\mu_l|\pi) \leq (1 - 2c_\gamma \lambda)^l \, \text{KL}(\mu_0|\pi)$.

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

**Problem:** not possible to combine both.

# First Experiments



Figure: The particle implementation of the SVGD algorithm illustrates the convergence of $\mathrm{KSD}^2(k \star \mu_l^n | \pi), \mathrm{KL}(k \star \mu_l^n | \pi)$ to 0.

# Not possible to combine both....

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

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

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

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

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
- $\pi$ with exploding $\beta$ moments with $\beta \geq 3$ (ex: a student distribution, which belongs to $\mathcal{P}_2$ the set of distributions with bounded second moment).

## Open question 2: SVGD quantisation

The quality of a set of points $(x^1, \ldots, x^n)$ can be measured by the integral approximation error:
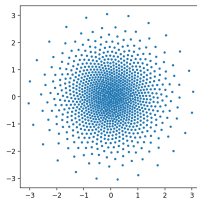
$$E(x_1, \ldots, x_n) = \left| \frac{1}{n} \sum_{i=1}^{n} f(x^i) - \int_{\mathbb{R}^d} f(x) d\pi(x) \right|. \tag{3}$$



(a) i.i.d.  (b) SVGD Gaussian $k$  (c) SVGD Laplace $k$

For i.i.d. points or MCMC iterates, (3) is of order $n^{-\frac{1}{2}}$. Can we bound (3) for SVGD final states?

Ongoing work with L. Xu and D. Slepcev.

# Outline

A lot of problems previously came from the fact that the KL is not defined for discrete measures $\mu_n$. Can we consider functionals that are well-defined for $\mu_n$?

A lot of problems previously came from the fact that the KL is not defined for discrete measures $\mu_n$. Can we consider functionals that are well-defined for $\mu_n$?

Remember the Kernel Stein discrepancy of $\mu$ relative to $\pi$:

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

With several integration by parts we have:

$$\begin{aligned}
\text{KSD}^2(\mu|\pi) &= \left\| P_{\mu,k} \nabla \log\left(\frac{\mu}{\pi}\right) \right\|_{\mathcal{H}_k}^2 \\
&= \int \int \nabla \log\left(\frac{\mu}{\pi}(x)\right) \nabla \log\left(\frac{\mu}{\pi}(y)\right) k(x,y) d\mu(x) d\mu(y) \\
&= \iint \nabla \log \pi(x)^T \nabla \log \pi(y) k(x,y) + \nabla \log \pi(x)^T \nabla_2 k(x,y) \\
&\quad + \nabla_1 k(x,y)^T \nabla \log \pi(y) + \boldsymbol{\nabla} \cdot_1 \nabla_2 k(x,y) d\mu(x) d\mu(y) \\
&:= \iint k_\pi(x,y) d\mu(x) d\mu(y).
\end{aligned}$$

**can be written in closed-form for discrete measures $\mu$.**

# 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 $\gamma$

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

---

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:

$$x_{n+1} = x_n - \gamma_n B_n^{-1} \nabla L(x_n) := x_n + \gamma_n d_n \tag{4}$$

where $B_n^{-1}$ is a p.s.d. matrix approximating the inverse Hessian at $x_n$.

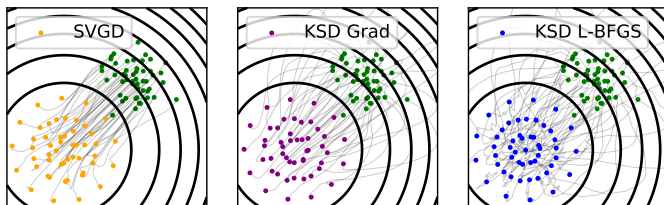Step1. (requires $\nabla L$) It computes a cheap version of $d_n$ based on BFGS recursion:

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

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

Step2. (requires $L$ and $\nabla L$) A line-search is performed to find the best step-size in (4) :

$$L(x_n + \gamma_n d_n) \leq L(x_n) + c_1 \gamma_n \nabla L(x_n)^T d_n$$
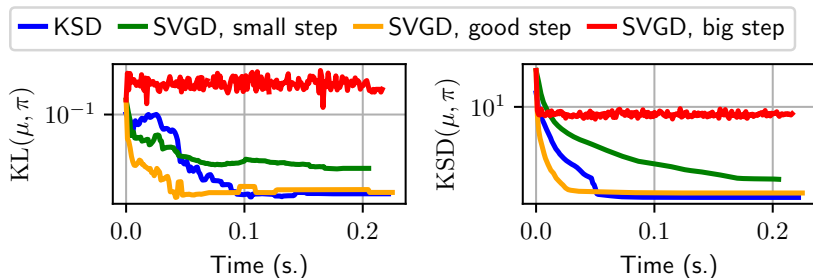$$\nabla L(x_n + \gamma_n d_n)^T d_n \geq c_2 \nabla L(x_n)^T d_n$$

# Toy experiments - 2D standard gaussian



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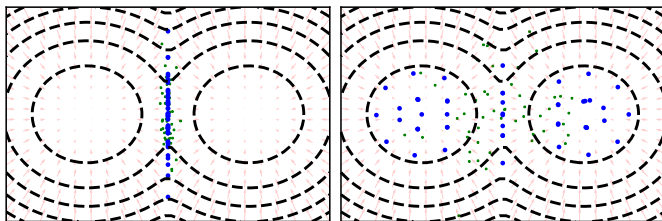The light grey curves correspond to their trajectories.

# SVGD vs KSD Descent - importance of the step-size



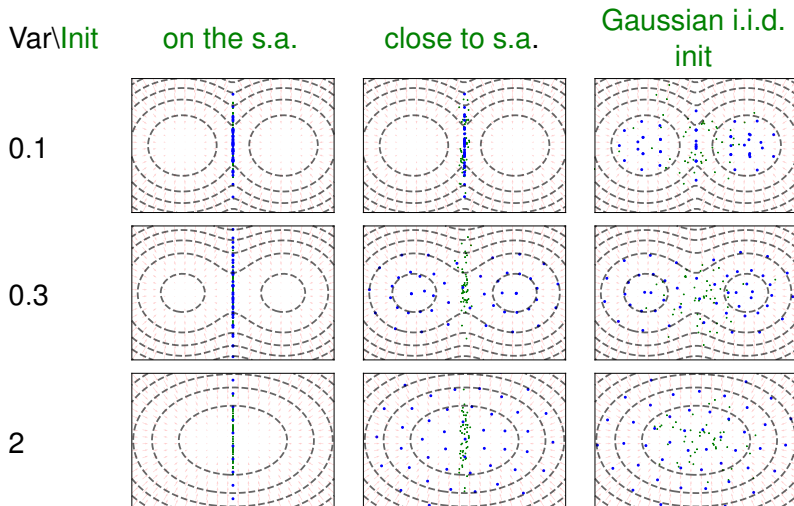Convergence speed of KSD and SVGD on a Gaussian problem in 1D, with 30 particles.

# 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



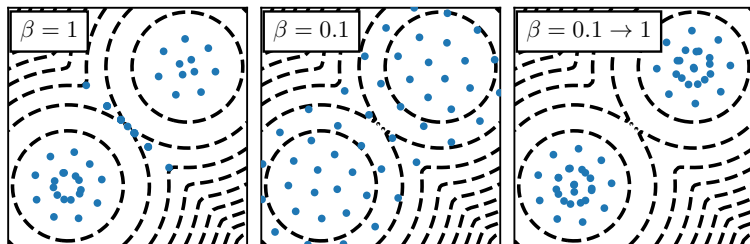| Var\Init | on the s.a. | close to s.a. | Gaussian i.i.d. init |
|----------|-------------|---------------|----------------------|
| 0.1 | | | |
| 0.3 | | | |
| 2 | | | |

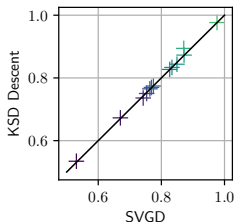Green crosses : initial particle positions
Blue crosses : final positions

# 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 $\beta$.)



This is a hard problem, even for Langevin diffusions, where tempering strategies also have been proposed [Lee et al., 2018].
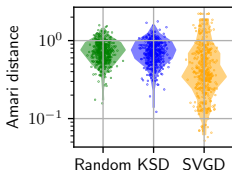
# Real world experiments (10 particles)



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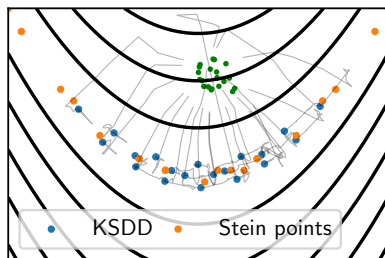
Bayesian ICA.
Each dot is the Amari distance between an estimated matrix and the true unmixing matrix ($d \leq 8$).
**KSD is not better than random.**
Hint: highly non-convex likelihood.

# 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.**

# Theoretical properties

Stationary measures:

▶ we show that if a stationary measure $\mu_\infty$ is full support, then $\mathcal{F}(\mu_\infty) = 0$.

▶ however, we also show that if $supp(\mu_0) \subset \mathcal{M}$, where $\mathcal{M}$ is a plane of symmetry of $\pi$, then for any time $t$ it remains true for $\mu_t$: $supp(\mu_t) \subset \mathcal{M}$.

# Theoretical properties

Stationary measures:

- ▶ we show that if a stationary measure $\mu_\infty$ is full support, then $\mathcal{F}(\mu_\infty) = 0$.

- ▶ however, we also show that if $supp(\mu_0) \subset \mathcal{M}$, where $\mathcal{M}$ is a plane of symmetry of $\pi$, then for any time $t$ it remains true for $\mu_t$: $supp(\mu_t) \subset \mathcal{M}$.

Explain convergence in the log-concave case? again an open question:

- ▶ the KSD is not geodesically convex

- ▶ it is not strongly geo convex near the global optimum $\pi$

- ▶ convergence of the continuous dynamics can be shown with a functional inequality, but which does not hold for discrete measures

# First strategy : obtain a functional inequality

How fast $\mathcal{F}(\mu_t)$ decreases along its WGF ?

$$\frac{\partial \mu_t}{\partial t} = \boldsymbol{\nabla} \cdot (\mu_t V_t), \quad V_t = \nabla_{W_2} \mathcal{F}(\mu_t)$$

$$\begin{aligned}
\frac{d\mathcal{F}(\mu_t)}{dt} &= \left\langle V_t, \nabla_{W_2} \mathcal{F}(\mu_t) \right\rangle_{L^2(\mu_t)} \\
&= -\left\| \nabla_{W_2} \mathcal{F}(\mu_t) \right\|_{L^2(\mu_t)}^2 \\
&= -\| \mathbb{E}_{x \sim \mu_t}[\nabla_2 k(x, y)] - \mathbb{E}_{x \sim \pi}[\nabla_2 k(x, y)] \|_{L^2(\mu_t)}^2 \\
&= -\underbrace{\| \nabla f_{\mu_t, \pi} \|_{L_2(\mu_t)}^2}_{\| f_{\mu_t, \pi} \|_{\dot{H}^{-1}(\mu_t)}}
\end{aligned}$$

where $f_{\mu_t, \pi} = \mathbb{E}_{x \sim \mu_t}[k(x, .)] - \mathbb{E}_{x \sim \pi}[k(x, .)]$.

# First strategy : obtain a functional inequality

How fast $\mathcal{F}(\mu_t)$ decreases along its WGF ?

$$\frac{\partial \mu_t}{\partial t} = \boldsymbol{\nabla} \cdot (\mu_t V_t), \quad V_t = \nabla_{W_2} \mathcal{F}(\mu_t)$$

$$\begin{aligned}
\frac{d\mathcal{F}(\mu_t)}{dt} &= \langle V_t, \nabla_{W_2} \mathcal{F}(\mu_t) \rangle_{L^2(\mu_t)} \\
&= -\|\nabla_{W_2} \mathcal{F}(\mu_t)\|^2_{L^2(\mu_t)} \\
&= -\|\mathbb{E}_{x \sim \mu_t}[\nabla_2 k(x, y)] - \mathbb{E}_{x \sim \pi}[\nabla_2 k(x, y)]\|^2_{L^2(\mu_t)} \\
&= -\underbrace{\|\nabla f_{\mu_t, \pi}\|^2_{L_2(\mu_t)}}_{\|f_{\mu_t, \pi}\|_{\dot{H}^{-1}(\mu_t)}}
\end{aligned}$$

where $f_{\mu_t, \pi} = \mathbb{E}_{x \sim \mu_t}[k(x, .)] - \mathbb{E}_{x \sim \pi}[k(x, .)]$.

It can be shown that:

$$\|f_{\mu_t, \pi}\|^2_{\mathcal{H}_k} \leq \|f_{\mu_t, \pi}\|_{\dot{H}(\mu_t)} \underbrace{\|\mu_t - \pi\|_{\dot{H}^{-1}(\mu_t)}}_{\sup_{\|g\|^2_{\dot{H}(\mu_t)} \leq 1} |\int g d\mu_t - \int g d\pi|}$$

Hence, if $\|\mu_t - \pi\|_{\dot{H}^{-1}(\mu_t)} \leq C$ for all $t \geq 0$, we have

$$\frac{d\mathcal{F}(\nu_t)}{dt} \leq -C\mathcal{F}(\nu_t)^2, \text{ hence}$$

$$\mathcal{F}(\mu_t) \leq \frac{1}{\mathcal{F}(\mu_0) + 4C^{-1}t}$$

where $\mathcal{F}(\mu_0) = \frac{1}{2}\text{MMD}^2(\mu_t, \pi)$.

Problems:

▶ depends on the whole sequence $(\mu_t)_{t \geq 0}$ (not only $\pi$)

▶ hard to verify in practice

▶ we observed convergence issues in practice (more for the MMD than the KSD)

# 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 $\text{Hess}_\mu \mathcal{F}(\psi, \psi) := \frac{d^2}{dt^2}\Big|_{t=0} \mathcal{F}(\rho_t)$,
which is related to the $W_2$ **Hessian of** $\mathcal{F}$ **at** $\mu$.

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

$$
\begin{aligned}
\text{Hess}_\mu \mathcal{F}(\psi, \psi) = \mathbb{E}_{x,y\sim\mu} & \left[ \nabla\psi(x)^T \nabla_1\nabla_2 k_\pi(x, y) \nabla\psi(y) \right] \\
+ \mathbb{E}_{x,y\sim\mu} & \left[ \nabla\psi(x)^T H_1 k_\pi(x, y) \nabla\psi(x) \right].
\end{aligned}
$$

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 $\pi$? the second term vanishes due to the Stein property of $k_\pi$ and :

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

where

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

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

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

$$\|S_{\pi, k_\pi} \mathcal{L}_\pi \psi\|^2_{\mathcal{H}_{k_\pi}} = \mathsf{Hess}_\pi \, \mathcal{F}(\psi, \psi) \geq \lambda \|\nabla \psi\|^2_{L^2(\pi)} \, ?$$

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

# Curvature near equilibrium - negative result

The previous inequality

$$\|S_{\pi,k_\pi}\mathcal{L}_\pi\psi\|^2_{\mathcal{H}_{k_\pi}} \geq \lambda\|\nabla\psi\|^2_{L^2(\pi)}$$

▶ can be seen as a kernelized version of the Poincaré inequality for $\pi$ :

$$\|\mathcal{L}_\pi\psi\|^2_{L_2(\pi)} \geq \lambda_\pi\|\nabla\psi\|^2_{L_2(\pi)}.$$

▶ can be written:

$$\langle\psi, P_{\pi,k_\pi}\psi\rangle_{L_2(\pi)} \geq \lambda\langle\psi, \mathcal{L}_\pi^{-1}\psi\rangle_{L_2(\pi)},$$

where $P_{\pi,k_\pi} : L^2(\pi) \to L^2(\pi), f \mapsto \int k_\pi(x,.)f(x)d\pi(x).$

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

# Conclusion

- ▶ Mixing kernels and Wasserstein gradient flows enable to design deterministic interacting particle systems

# Conclusion

▶ Mixing kernels and Wasserstein gradient flows enable to design deterministic interacting particle systems

▶ They can provide a better approximation of the target for a finite number of particles

# Conclusion

- ▶ Mixing kernels and Wasserstein gradient flows enable to design deterministic interacting particle systems

- ▶ They can provide a better approximation of the target for a finite number of particles

- ▶ Theory does not match practice yet

# Conclusion

- ▶ Mixing kernels and Wasserstein gradient flows enable to design deterministic interacting particle systems
- ▶ They can provide a better approximation of the target for a finite number of particles
- ▶ Theory does not match practice yet
- ▶ Numerics can be improved, via perturbed dynamics, change of geometry...

# Conclusion

▶ Mixing kernels and Wasserstein gradient flows enable to design deterministic interacting particle systems

▶ They can provide a better approximation of the target for a finite number of particles

▶ Theory does not match practice yet

▶ Numerics can be improved, via perturbed dynamics, change of geometry...

▶ Python package to try KSD descent:
  **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
```

# References I

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.

Ba, J., Erdogdu, M. A., Ghassemi, M., Sun, S., Suzuki, T., Wu, D., and Zhang, T. (2021).
Understanding the variance collapse of svgd in high dimensions.
In *International Conference on Learning Representations*.

# 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.

📄 Chwialkowski, K., Strathmann, H., and Gretton, A. (2016).
A kernel test of goodness of fit.
In *International conference on machine learning*.

📄 Duncan, A., Nüsken, N., and Szpruch, L. (2019).
On the geometry of stein variational gradient descent.
*arXiv preprint arXiv:1912.00894*.

# References III

📄 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.

📄 Jordan, R., Kinderlehrer, D., and Otto, F. (1998).
The variational formulation of the fokker–planck equation.
*SIAM journal on mathematical analysis*, 29(1):1–17.

📄 Karimi, H., Nutini, J., and Schmidt, M. (2016).
Linear convergence of gradient and proximal-gradient methods under the polyak-łojasiewicz condition.
In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 795–811. Springer.

# References IV

📑 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.*

📑 Lee, H., Risteski, A., and Ge, R. (2018).
Beyond log-concavity: Provable guarantees for sampling multi-modal distributions using simulated tempering langevin monte carlo.
*Advances in neural information processing systems*, 31.

📑 Liu, D. C. and Nocedal, J. (1989).
On the limited memory BFGS method for large scale optimization.
*Mathematical programming*, 45(1-3):503–528.

# References V

📄 Liu, Q. (2017).
Stein variational gradient descent as gradient flow.
In *Advances in neural information processing systems*,
pages 3115–3123.

📄 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.

📄 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.

# References VI

📄 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.

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

📄 Wenliang, L. K. and Kanagawa, H. (2020).
Blindness of score-based methods to isolated components and mixing proportions.
*arXiv preprint arXiv:2008.10087*.