Accurate Quantization of Measures via Interacting Particle-based Optimization

Anna Korba ENSAE, CREST, IP Paris

INRIA Thoth seminar

Joint work with Lantian Xu, Dejan Slepčev (Carnegie Mellon University).



Problem and Motivation

Background on MMD/KSD Descent

MMD and KSD Quantization

Experiments

Quantization problem

Problem : approximate a target distribution $\pi \in \mathcal{P}(\mathbb{R}^d)$ by a finite set of *n* points x_1, \ldots, x_n , e.g. to compute functionals $\int_{\mathbb{R}^d} f(x) d\pi(x)$.

The quality of the set can be measured by the integral approximation error:

$$\operatorname{err}(x_1,\ldots,x_n) = \left| \frac{1}{n} \sum_{i=1}^n f(x_i) - \int_{\mathbb{R}^d} f(x) d\pi(x) \right|.$$

Several approaches, among which :

- MCMC methods : generate a Markov chain whose law converges to π, err(x₁,...,x_n) = O(n^{-1/2}) [Łatuszyński et al., 2013]
- deterministic particle systems, $err(x_1, \ldots, x_n)$?

Application: Bayesian statistics

- Let $\mathcal{D} = (x_i, y_i)_{i=1,...,m}$ a labelled dataset.
- Assume an underlying model parametrized by z ∈ ℝ^d, e.g. y ~ f(x, z) + ε (p(y|x, z) gaussian)

 \implies Compute the likelihood: $p(\mathcal{D}|z) = \prod_{i=1}^{m} p(y_i|x_i, z)$.

► Assume a prior distribution on the parameter *z* ~ *p*.

Application: Bayesian statistics

• Let $\mathcal{D} = (x_i, y_i)_{i=1,...,m}$ a labelled dataset.

Assume an underlying model parametrized by $z \in \mathbb{R}^d$, e.g. $y \sim f(x, z) + \epsilon$ (p(y|x, z) gaussian)

 \implies Compute the likelihood: $p(\mathcal{D}|z) = \prod_{i=1}^{m} p(y_i|x_i, z)$.

► Assume a prior distribution on the parameter *z* ~ *p*.

Bayes' rule :
$$\pi(z) := p(z|\mathcal{D}) = rac{p(\mathcal{D}|z)p(z)}{C}$$
 , $C = \int_{\mathbb{R}^d} p(\mathcal{D}|z)p(z)dz$.

 π is known up to a constant since *C* is intractable. How to sample from π then? e.g. to compute:

$$p(y|x, \mathcal{D}) = \int_{\mathbb{R}^d} p(y|x, z) d\pi(z)$$



Problem and Motivation

Background on MMD/KSD Descent

MMD and KSD Quantization

Experiments

Sampling as optimization over distributions

3 algorithms/particle systems at study:

- Maximum Mean Discrepancy Descent [Arbel et al., 2019]
- ► Kernel Stein Discrepancy Descent [Korba et al., 2021]
- Stein Variational Gradient Descent [Liu and Wang, 2016]

These particle systems are designed to minimize a loss.

Sampling as optimization over distributions

3 algorithms/particle systems at study:

- Maximum Mean Discrepancy Descent [Arbel et al., 2019]
- ► Kernel Stein Discrepancy Descent [Korba et al., 2021]
- Stein Variational Gradient Descent [Liu and Wang, 2016]

These particle systems are designed to minimize a loss.

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:

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

where D is a dissimilarity functional and \mathcal{F} "a loss".

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 .

Euclidean gradient flow and continuity equation

Let $V : \mathbb{R}^d \to \mathbb{R}$ and consider minimizing *V*. The gradient flow of *V* can be written

$$\frac{dx_t}{dt} = -\nabla V(x_t)$$

and assume x_0 random with density μ_0 . What is the dynamics of the density μ_t of x_t ? Let $\phi : \mathbb{R}^d \to \mathbb{R}$ a smooth function with compact support.

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

and applying the chain rule and using I.P.P.,

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

Therefore,

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

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

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

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

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

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 :

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

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

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 :

$$W_{2}^{2}(\nu,\mu) = \inf_{\boldsymbol{s}\in\Gamma(\nu,\mu)} \int_{\mathbb{R}^{d}\times\mathbb{R}^{d}} \left\|\boldsymbol{x}-\boldsymbol{y}\right\|^{2} d\boldsymbol{s}(\boldsymbol{x},\boldsymbol{y}) \qquad \forall \nu,\mu\in\mathcal{P}_{2}(\mathbb{R}^{d})$$

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

Can also be written (Benamou-Brenier formula):

$$W_{2}^{2}(\nu,\mu) = \inf_{(\rho_{t},v_{t})_{t\in[0,1]}} \left\{ \int_{0}^{1} \|v_{t}\|_{L^{2}(\rho_{t})}^{2} dt(x) : \frac{\partial \rho_{t}}{\partial t} = \nabla \cdot (\rho_{t}v_{t}), \rho_{0} = \nu, \rho_{1} = \mu \right\}.$$

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

The first variation of $\mu \mapsto \mathcal{F}(\mu)$ evaluated at $\mu \in \mathcal{P}(\mathbb{R}^d)$ is the unique function $\frac{\partial \mathcal{F}(\mu)}{\partial \mu} : \mathbb{R}^d \to \mathbb{R}$ s. t. for any $\mu, \nu \in \mathcal{P}(\mathbb{R}^d), \nu - \mu \in \mathcal{P}(\mathbb{R}^d)$: $\lim_{\epsilon \to 0} \frac{1}{\epsilon} (\mathcal{F}(\mu + \epsilon(\nu - \mu)) - \mathcal{F}(\mu)) = \int_{\mathbb{R}^d} \frac{\partial \mathcal{F}(\mu)}{\partial \mu} (x) (d\nu - d\mu) (x).$

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

The first variation of $\mu \mapsto \mathcal{F}(\mu)$ evaluated at $\mu \in \mathcal{P}(\mathbb{R}^d)$ is the unique function $\frac{\partial \mathcal{F}(\mu)}{\partial \mu} : \mathbb{R}^d \to \mathbb{R}$ s. t. for any $\mu, \nu \in \mathcal{P}(\mathbb{R}^d), \nu - \mu \in \mathcal{P}(\mathbb{R}^d)$: $\lim_{\epsilon \to 0} \frac{1}{\epsilon} (\mathcal{F}(\mu + \epsilon(\nu - \mu)) - \mathcal{F}(\mu)) = \int_{\mathbb{R}^d} \frac{\partial \mathcal{F}(\mu)}{\partial \mu} (x) (d\nu - d\mu)(x).$

The family $\mu : [0, \infty] \to \mathcal{P}_2(\mathbb{R}^d), t \mapsto \mu_t$ satisfies a Wasserstein gradient flow of \mathcal{F} if distributionnally:

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

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} .

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

The first variation of $\mu \mapsto \mathcal{F}(\mu)$ evaluated at $\mu \in \mathcal{P}(\mathbb{R}^d)$ is the unique function $\frac{\partial \mathcal{F}(\mu)}{\partial \mu} : \mathbb{R}^d \to \mathbb{R}$ s. t. for any $\mu, \nu \in \mathcal{P}(\mathbb{R}^d), \nu - \mu \in \mathcal{P}(\mathbb{R}^d)$: $\lim_{\epsilon \to 0} \frac{1}{\epsilon} (\mathcal{F}(\mu + \epsilon(\nu - \mu)) - \mathcal{F}(\mu)) = \int_{\mathbb{R}^d} \frac{\partial \mathcal{F}(\mu)}{\partial \mu} (x) (d\nu - d\mu)(x).$

The family $\mu : [0, \infty] \to \mathcal{P}_2(\mathbb{R}^d), t \mapsto \mu_t$ satisfies a Wasserstein gradient flow of \mathcal{F} if distributionnally:

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

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} .

It can be implemented by the deterministic process:

$$rac{dX_t}{dt} = -
abla_{W_2}\mathcal{F}(\mu_t)(X_t), \quad X_t \sim \mu_t$$

Particle system approximating the WGF

Euler time-discretization : in \mathbb{R}^d , move particles as:

$$X_{l+1} = X_l - \gamma \nabla_{W_2} \mathcal{F}(\mu_l)(X_l) \sim \mu_{l+1}, \quad X_0 \sim \mu_0.$$

Space discretization/particle system : Since μ_l is unknown, introduce a particle system X^1, \ldots, X^n where μ_l is replaced by $\hat{\mu}_l = \frac{1}{n} \sum_{i=1}^n \delta_{X_l^i}$:

$$\begin{aligned} X_{l+1}^{i} &= X_{l}^{i} - \gamma \nabla_{W_{2}} \mathcal{F}(\hat{\mu}_{l})(X_{l}^{i}) \quad \text{ for } i = 1, \dots, n, \\ X_{0}^{1}, \dots, X_{0}^{n} &\sim \mu_{0}. \end{aligned}$$

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||)^{-β} with β ∈]0,1[

► *H_k* its corresponding RKHS (Reproducing Kernel Hilbert Space):

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

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

It satisfies the reproducing property:

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

Maximum Mean Discrepancy [Gretton et al., 2012]

Assume $\mu \mapsto \int k(x, .) d\mu(x)$ injective.

Maximum Mean Discrepancy defines a distance on $\mathcal{P}_2(\mathbb{R}^d)$:

$$\begin{split} \mathsf{MMD}^2(\mu,\pi) &= \sup_{f \in \mathcal{H}_k, \|f\|_{\mathcal{H}_k} \le 1} \left| \int f d\mu - \int f d\pi \right|^2 \\ &= \|m_\mu - m_\pi\|_{\mathcal{H}_k}^2 \\ &= \iint_{\mathbb{R}^d} k(x,y) d\mu(x) d\mu(y) + \iint_{\mathbb{R}^d} k(x,y) d\pi(x) d\pi(y) \\ &- 2 \iint_{\mathbb{R}^d} k(x,y) d\mu(x) d\pi(y), \end{split}$$

by the reproducing property $\langle f, k(x, .) \rangle_{\mathcal{H}_k} = f(x)$ for $f \in \mathcal{H}_k$.

Maximum Mean Discrepancy [Gretton et al., 2012]

Assume $\mu \mapsto \int k(x, .) d\mu(x)$ injective.

Maximum Mean Discrepancy defines a distance on $\mathcal{P}_2(\mathbb{R}^d)$:

$$\begin{split} \mathsf{MMD}^2(\mu,\pi) &= \sup_{f \in \mathcal{H}_k, \|f\|_{\mathcal{H}_k} \le 1} \left| \int f d\mu - \int f d\pi \right|^2 \\ &= \|m_\mu - m_\pi\|_{\mathcal{H}_k}^2 \\ &= \iint_{\mathbb{R}^d} k(x,y) d\mu(x) d\mu(y) + \iint_{\mathbb{R}^d} k(x,y) d\pi(x) d\pi(y) \\ &- 2 \iint_{\mathbb{R}^d} k(x,y) d\mu(x) d\pi(y), \end{split}$$

by the reproducing property $\langle f, k(x, .) \rangle_{\mathcal{H}_k} = f(x)$ for $f \in \mathcal{H}_k$.

The differential of $\mu \mapsto \frac{1}{2} \text{MMD}^2(., \pi)$ evaluated at $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ is:

$$\int k(x,.)d\mu(x) - \int k(x,.)d\pi(x) : \mathbb{R}^d \to \mathbb{R}.$$

Hence, for *k* regular enough, $\nabla_{W_2} \frac{1}{2} \text{MMD}^2(\mu, \pi)$ is:

$$\int \nabla_2 k(x,.) d\mu(x) - \int \nabla_2 k(x,.) d\pi(x) : \mathbb{R}^d \to \mathbb{R}.$$
 12/33

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

If one does not have access to samples of π but only to its score, it is still possible to compute the KSD:

$$\mathsf{KSD}^{2}(\mu|\pi) = \iint k_{\pi}(x, y) d\mu(x) d\mu(y),$$

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 $k : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}, k \in C^2(\mathbb{R}^d)^1$

For $x, y \in \mathbb{R}^d$,

 $k_{\pi}(x, y) = s(x)^{T} s(y) k(x, y) + s(x)^{T} \nabla_{2} k(x, y)$ $+ \nabla_{1} k(x, y)^{T} s(y) + \nabla \cdot_{1} \nabla_{2} k(x, y)$ $= \sum_{i=1}^{d} \frac{\partial \log \pi(x)}{\partial x_{i}} \cdot \frac{\partial \log \pi(y)}{\partial y_{i}} \cdot k(x, y) + \frac{\partial \log \pi(x)}{\partial x_{i}} \cdot \frac{\partial k(x, y)}{\partial y_{i}}$

$$+\frac{\partial \log \pi(y)}{\partial y_i} \cdot \frac{\partial k(x,y)}{\partial x_i} + \frac{\partial^2 k(x,y)}{\partial x_i \partial y_i} \in \mathbb{R}.$$
¹e.g. : $k(x,y) = \exp(-||x-y||^2/h)$

13/33

KSD vs 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:

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

KSD as kernelized Fisher Divergence

Fisher Divergence:

$$\mathsf{FD}^{2}(\mu|\pi) = \left\| \nabla \log\left(\frac{\mu}{\pi}\right) \right\|_{L^{2}(\mu)}^{2} = \int \|\nabla \log\left(\frac{\mu}{\pi}(x)\right)\|^{2} d\mu(x)$$

"Kernelized" with k:

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

where
$$\mathcal{S}_{\mu,k}: L^2(\mu) o \mathcal{H}_k$$

 $f \mapsto \int k(x,.)f(x)d\mu(x).$

 \implies minimizing the KSD is close in spirit to score-matching [Hyvärinen and Dayan, 2005].

Recall that we want to study particle systems

$$\begin{aligned} X_{l+1}^{i} &= X_{l}^{i} - \gamma \nabla_{W_{2}} \mathcal{F}(\hat{\mu}_{l})(X_{l}^{i}) \quad \text{ for } i = 1, \dots, n, \end{aligned}$$

where $\hat{\mu}_{l} &= \frac{1}{n} \sum_{i=1}^{n} \delta_{X_{l}^{i}} \text{ and } \mathcal{F}(\mu) = \mathrm{D}(\mu | \pi). \end{aligned}$

Recall that we want to study particle systems

$$\begin{split} X_{l+1}^{i} &= X_{l}^{i} - \gamma \nabla_{W_{2}} \mathcal{F}(\hat{\mu}_{l})(X_{l}^{i}) \quad \text{for } i = 1, \dots, n, \\ \text{where } \hat{\mu}_{l} &= \frac{1}{n} \sum_{i=1}^{n} \delta_{X_{l}^{i}} \text{ and } \mathcal{F}(\mu) = \mathrm{D}(\mu | \pi). \end{split}$$

For discrete measures $\mu = \frac{1}{n} \sum_{i=1}^{n} \delta_{X^{i}}$, the MMD/KSD are well defined, hence we let $F(X^{1}, \dots, X^{n}) := \mathcal{F}(\mu)$.

Recall that we want to study particle systems

$$\begin{split} X_{l+1}^{i} &= X_{l}^{i} - \gamma \nabla_{W_{2}} \mathcal{F}(\hat{\mu}_{l})(X_{l}^{i}) \quad \text{ for } i = 1, \dots, n \\ \text{where } \hat{\mu}_{l} &= \frac{1}{n} \sum_{i=1}^{n} \delta_{X_{l}^{i}} \text{ and } \mathcal{F}(\mu) = \mathrm{D}(\mu | \pi). \end{split}$$

For discrete measures $\mu = \frac{1}{n} \sum_{i=1}^{n} \delta_{X^{i}}$, the MMD/KSD are well defined, hence we let $F(X^{1}, \dots, X^{n}) := \mathcal{F}(\mu)$.

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

$$\nabla_{x^i}F(X^1,\ldots,X^n)=\frac{1}{n}\sum_{j=1}^n\nabla_2k(X^j,X^j)-\int\nabla_2k(X^j,x)d\pi(x).$$

In contrast, if D is the KSD,

$$\nabla_{X^i}F(X^1,\ldots,X^n)=\frac{1}{n}\sum_{j=1}^n\nabla_2k_{\pi}(X^j,X^j).$$

Recall that we want to study particle systems

$$\begin{aligned} X_{l+1}^{i} &= X_{l}^{i} - \gamma \nabla_{W_{2}} \mathcal{F}(\hat{\mu}_{l})(X_{l}^{i}) \quad \text{ for } i = 1, \dots, n \\ \text{where } \hat{\mu}_{l} &= \frac{1}{n} \sum_{i=1}^{n} \delta_{X_{l}^{i}} \text{ and } \mathcal{F}(\mu) = \mathrm{D}(\mu | \pi). \end{aligned}$$

For discrete measures $\mu = \frac{1}{n} \sum_{i=1}^{n} \delta_{X^{i}}$, the MMD/KSD are well defined, hence we let $F(X^{1}, \dots, X^{n}) := \mathcal{F}(\mu)$.

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

$$\nabla_{x^i}F(X^1,\ldots,X^n)=\frac{1}{n}\sum_{j=1}^n\nabla_2k(X^j,X^j)-\int\nabla_2k(X^j,x)d\pi(x).$$

In contrast, if D is the KSD,

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

MMD/KSD Descent: at each time $l \ge 0$, for any i = 1, ..., n:

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

• The MMD/KSD/their W_2 gradient write as sums of integrals of μ and π

- The MMD/KSD/their W₂ gradient write as sums of integrals of μ and π
- Hence they can be evaluated in closed form for discrete μ and π ⇒ use L-BFGS to automatically select the best step-size

- The MMD/KSD/their W₂ gradient write as sums of integrals of μ and π
- Hence they can be evaluated in closed form for discrete μ and π ⇒ use L-BFGS to automatically select the best step-size
- depending on the information on π, choose the KSD (unnormalized density) or MMD (samples)

- The MMD/KSD/their W₂ gradient write as sums of integrals of μ and π
- Hence they can be evaluated in closed form for discrete μ and π ⇒ use L-BFGS to automatically select the best step-size
- depending on the information on π, choose the KSD (unnormalized density) or MMD (samples)
- The MMD upper bounds the integral approximation error for functions in the RKHS, since by the reproducing property and Cauchy-Schwartz:

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

Similarly for the KSD with $\mathcal{H}_{k_{\pi}}$.

Stein Variational Gradient Descent [Liu and Wang, 2016]

Stein Variational Gradient Descent (SVGD) performs gradient descent in $\mathcal{P}(\mathbb{R}^d)$ of the Kullback-Leibler (KL) divergence :

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

where the (W_2) gradient is smoothed through the kernel integral operator.

It corresponds to an Euler discretization of the gradient flow of the KL under a metric depending on k [Duncan et al., 2019]:

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

Stein Variational Gradient Descent [Liu and Wang, 2016]

Fix a reproducing kernel k. In continuous time, SVGD flow is defined by the continuity equation

$$\frac{\partial \mu_t}{\partial t} + \boldsymbol{\nabla} \cdot (\mu_t \boldsymbol{v}_{\mu_t}) = \boldsymbol{0}, \ \boldsymbol{v}_{\mu_t} = \boldsymbol{k} \star (\mu_t \nabla \log \pi) - \nabla \boldsymbol{k} \star \mu_t,$$

i.e. $\textit{v}_{\mu_t} = \textit{S}_{\mu_t,k}
abla \log\left(rac{\mu}{\pi}
ight)$ where

$$\blacktriangleright \nabla \log \left(\frac{\mu}{\pi}\right) = \nabla_{W_2} \operatorname{KL}(\mu|\pi),$$

$$\blacktriangleright \ S_{\mu,k}: L^2(\mu) \to \mathcal{H}_k, f \mapsto \int k(x,.)f(x)d\mu(x).$$

Stein Variational Gradient Descent [Liu and Wang, 2016]

Fix a reproducing kernel k. In continuous time, SVGD flow is defined by the continuity equation

$$\frac{\partial \mu_t}{\partial t} + \boldsymbol{\nabla} \cdot (\mu_t \boldsymbol{v}_{\mu_t}) = \boldsymbol{0}, \ \boldsymbol{v}_{\mu_t} = \boldsymbol{k} \star (\mu_t \nabla \log \pi) - \nabla \boldsymbol{k} \star \mu_t,$$

i.e. $\textit{v}_{\mu_t} = \textit{S}_{\mu_t,k}
abla \log\left(rac{\mu}{\pi}
ight)$ where

$$\blacktriangleright \nabla \log \left(\frac{\mu}{\pi}\right) = \nabla_{W_2} \operatorname{KL}(\mu|\pi),$$

$$\blacktriangleright \ S_{\mu,k}: L^2(\mu) \to \mathcal{H}_k, f \mapsto \int k(x,.)f(x)d\mu(x).$$

Let $\gamma > 0$ be a fixed step-size. Starting from $x_0^1, \ldots, x_0^n \sim \mu_0$, SVGD algorithm updates the *n* particles as follows at each iteration :

$$x_{l+1}^{i} = x_{l}^{i} - \frac{\gamma}{n} \sum_{j=1}^{n} \left[-\nabla \log \pi(x_{l}^{j}) k(x_{l}^{i}, x_{l}^{j}) + \nabla_{x_{l}^{j}} k(x_{l}^{i}, x_{l}^{j}) \right].$$

Remark: does not minimize a well-defined functional for discrete measures, it is only a discrete approximation of the flow. Hence, cannot be used with L-BFGS and measuring the quantization is tricky. 19/33



Problem and Motivation

Background on MMD/KSD Descent

MMD and KSD Quantization

Experiments

Motivation - Final states for a Gaussian target



Figure: Final states of the algorithms for 1000 particles, kernel bandwidth = 1. k_G is the Gaussian kernel and k_L is the Laplace kernel

MMD gradient is available in closed form for $\pi = \mathcal{N}(\mathbf{0}_d, \theta I_d)$

$$\dot{x}_{i} = -\frac{1}{nh^{2}(\sqrt{2\pi h^{2}})^{d}} \sum_{j=1}^{n} e^{-\frac{|x_{j}-x_{i}|^{2}}{2h^{2}}} (x_{j}-x_{i}) - \frac{1}{(h^{2}+\theta^{2})(\sqrt{2\pi (h^{2}+\theta^{2})})^{d}} e^{-\frac{|x_{i}|^{2}}{2(h^{2}+\theta^{2})}} x_{i}.$$

We are interested in establishing bounds on the quantization error

$$Q_n = \inf_{X_n = x_1, \dots, x_n} D(\pi, \mu_n), \quad \text{ for } \mu_n = \frac{1}{n} \sum_{i=1}^n \delta_{x_i},$$

where D is the MMD or KSD.

Remark: For $x_1, \ldots, x_n \sim \pi$ i.i.d., the rate is known to be $\mathcal{O}(n^{-1/2})$ [Gretton et al., 2006, Tolstikhin et al., 2017, Lu and Lu, 2020].

We are interested in establishing bounds on the quantization error

$$Q_n = \inf_{X_n = x_1, \dots, x_n} D(\pi, \mu_n), \quad \text{ for } \mu_n = \frac{1}{n} \sum_{i=1}^n \delta_{x_i},$$

where D is the MMD or KSD.

Remark: For $x_1, \ldots, x_n \sim \pi$ i.i.d., the rate is known to be $\mathcal{O}(n^{-1/2})$ [Gretton et al., 2006, Tolstikhin et al., 2017, Lu and Lu, 2020].

Assumption A1: Assume that the kernel is *d*-times continuously differentiable. Assume also that any mixed partial derivative of the kernel of order smaller than *d* has a RKHS norm bounded by a constant $C_{k,d} \ge 0$.

First result for the MMD

Theorem: Suppose A1 holds. Assume that (i) π is the Lebesgue measure or (ii) a probability measure on $[0, 1]^d$. Then, there exists a constant C_d , such that for all $n \ge 2$,

• if (i): there exist points x_1, \ldots, x_n such that

$$\mathrm{MMD}(\pi,\mu_n) \leq C_d \frac{(\log n)^{d-1}}{n}$$

• if (ii): there exist points x_1, \ldots, x_n such that

$$\mathrm{MMD}(\pi,\mu_n) \leq C_d \frac{(\log n)^{\frac{3d+1}{2}}}{n}$$

Proof: We use the well-known Koksma-Hlawka inequality [Aistleitner and Dick, 2015](Th1):

$$\left|\int_{[0,1]^d} f(x) d\pi(x) - \frac{1}{n} \sum_{i=1}^n f(x_i)\right| \leq \mathcal{D}(X_n, \pi) V(f),$$

• $\mathcal{D}(X_n, \pi) = 2^d \sup_{I = \prod_{i=1}^n [a_i, b_i]} |\pi(I) - \mu_n(I)|$ is the discrepancy of the point set X_n , can be bounded by [Aistleitner and Dick, 2015]

The variation of a function f : [0, 1]^d → ℝ with continuous mixed partial derivatives is defined as

$$V(f) = \sum_{\alpha \subseteq \{1,\ldots,d\}} \int_{[0,1]^{|\alpha|}} \left| \frac{\partial^{|\alpha|} f(x_{\alpha},1)}{\partial x_{\alpha}} \right| dx_{\alpha}.$$

Then, use the reproducing property on partial derivatives with Cauchy-Schwarz inequality, and **A1**:

$$\left|rac{\partial^{|lpha|} f(x_lpha,1)}{\partial x_lpha}
ight| \qquad \leq \qquad \left\|rac{\partial^{|lpha|} k((x_lpha,1),\cdot)}{\partial^{|lpha|} x_lpha}
ight\|_{\mathcal{H}_k} \|f\|_{\mathcal{H}_k} \qquad \leq \qquad \mathcal{C}_{k,d}.$$

Result for non compactly supported distributions π

Proposition 1: Suppose A1 holds and that *k* is bounded. Assume π is a light-tailed distribution on \mathbb{R}^d (i.e. which has a thinner tail than an exponential distribution). Then, for $n \ge 2$ there exist points $x_1, ..., x_n$ such that

$$\mathrm{MMD}(\pi,\mu_n) \leq C_d \frac{(\log n)^{\frac{5d+1}{2}}}{n}.$$

Result for non compactly supported distributions π

Proposition 1: Suppose A1 holds and that *k* is bounded. Assume π is a light-tailed distribution on \mathbb{R}^d (i.e. which has a thinner tail than an exponential distribution). Then, for $n \ge 2$ there exist points $x_1, ..., x_n$ such that

$$\mathrm{MMD}(\pi,\mu_n) \leq C_d \frac{(\log n)^{\frac{5d+1}{2}}}{n}.$$

Proof: Decompose $MMD(\pi, \mu_n) \leq MMD(\pi, \mu) + MMD(\mu, \mu_n)$, choosing μ compactly supported on $A_n = [-\log n, \log n]^d$.

As π is light-tailed, $\|\mu - \pi\|_{TV} \leq C_1/n$ distance, and we first get $MMD(\pi, \mu) \leq C_k \|\mu - \pi\|_{TV} \leq C/n$.

Then, we can take a discrete μ_n supported on A_n and bound MMD(μ, μ_n) using similar arguments as in the previous Theorem.

Result for the KSD

Theorem: Assume that *k* is a Gaussian kernel and that $\pi \propto \exp(-U)$ with $U \in C^{\infty}(\mathbb{R}^d)$. Assume furthermore that $U(x) > c_1 ||x||$ for large enough *x*, and that there exists a real-valued polynomial *V* of degree $m \ge 0$, such that for any multi-index β , $\left|\frac{\partial^{\beta} U(x)}{\partial^{\beta_1} x_1...\partial^{\beta_j} x_j}\right| \le V(x)$ for all $1 \le |\beta| \le d + 1$. Then there exist points $x_1, ..., x_n$ such that

$$ext{KSD}(\mu_n|\pi) \leq C_d rac{(\log n)^{rac{6d+2m+1}{2}}}{n}$$

Satisfied for gaussian mixtures π .

Result for the KSD

Theorem: Assume that *k* is a Gaussian kernel and that $\pi \propto \exp(-U)$ with $U \in C^{\infty}(\mathbb{R}^d)$. Assume furthermore that $U(x) > c_1 ||x||$ for large enough *x*, and that there exists a real-valued polynomial *V* of degree $m \ge 0$, such that for any multi-index β , $\left|\frac{\partial^{\beta} U(x)}{\partial^{\beta_1} x_1 \dots \partial^{\beta_j} x_j}\right| \le V(x)$ for all $1 \le |\beta| \le d + 1$. Then there exist points $x_1, ..., x_n$ such that

$$\mathrm{KSD}(\mu_n|\pi) \leq C_d \frac{(\log n)^{\frac{6d+2m+1}{2}}}{n}$$

Satisfied for gaussian mixtures π .

Proof: The proof relies on bounding the first and last term of the

$$\mathsf{KSD}(\mu_n, \pi) = 2 \iint \nabla \log(\pi)(x)^T \nabla_y k(x, y) d\mu(x) d\mu(y) + \underbrace{\iint \nabla \log(\pi)(x)^T \nabla \log(\pi)(y) k(x, y) d\mu(x) d\mu(y)}_{(1)} + \underbrace{\iint \nabla \cdot_x \nabla_y k(x, y) d\mu(x) d\mu(y)}_{(2)},$$

 $\mu = \mu_n - \pi$, as the cross terms can be upper bounded by the former ones by a simple computation.

(1) MMD(μ_n, π), with $k_1(x, y) = s(x)^T s(y) k(x, y)$, bounded by controlling $\|\nabla \log \pi\|_{H^d}$



Problem and Motivation

Background on MMD/KSD Descent

MMD and KSD Quantization

Experiments

Algorithms

we investigate numerically the quantization properties of :

- SVGD
- MMD descent
- KSD Descent
- Kernel Herding (KH) : greedy minimization of the MMD
- Stein points (SP) : greedy minimization of the KSD

Hyperparameters:

- kernel: Gaussian, Laplace...
- bandwith of the kernel
- step-size

Quantization rates of the algorithms, $\pi = \mathcal{N}(0, 1/dI_d)$



Averaged over 3 runs of each algorithm, run for 1e4 iterations, where the initial particles are i.i.d. samples of π . MMD/KSD Descent use bandwidth 1; Stein points use gridsize = 200 points in 2d, 50 in 3d; in 4d grid search was too slow.

| d | Eval. | SVGD | MMD-lbfgs | KSD-lbfgs | KH | SP |
|---|------------|----------------|----------------|----------------|----------------|----------------|
| 2 | KSD MMD | -0.98 -1.04 | -1.48 -1.60 | -1.46 -1.54 | -0.84 -0.93 | -0.77 -0.77 |
| 3 | KSD MMD | -0.91 -0.96 | -1.38 -1.51 | -1.44 -1.49 | -0.84 -0.92 | -0.78 -0.75 |
| 4 | KSD MMD | -0.91 -0.94 | -1.35 -1.46 | -1.39 -1.40 | -0.89 -0.95 | _ |
| 8 | KSD MMD | -0.84 -0.77 | -1.14 -1.25 | -1.16 -1.13 | _ | _ |

Some remarks:

- The slopes remain much steeper than the Monte Carlo rate, even when the dimension increases
- Their slopes are better than our theoretical upper bounds

Robustness to evaluation discrepancy



Figure: 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.

- if we measure the discrepancy using a kernel with smaller bandwidth, MMD and KSD results deteriorate significantly and SVGD/NSVGD perform the best.
- likely reason : Samples of MMD and KSD with Gaussian kernel have internal structures which can affect the discrepancy at lower bandwidths.

For $\nu, \mu \in \mathcal{P}_p(\mathbb{R}^d)$, the Sliced *p*-Wasserstein distance is defined as:

$$d_{sw,p}(\nu,\mu) = \int_{\mathbb{S}^{d-1}} W_p(P_{\theta\#}\nu, P_{\theta\#}\mu) d\theta,$$

where $P_{\theta} : x \mapsto x \cdot \theta$ and # is the pushforward operator.



Figure: Quantization rates measured in Sliced Wasserstein distance of the algorithms $\pi = \mathcal{N}(0, 1/dl_d)$. In practice, we use p = 1 and 50 random directions drawn uniformly on \mathbb{S}^{d-1} to discretize the integration.

The rates for NSVGD are approximately $n^{-0.72}$, $n^{-0.65}$, $n^{-0.63}$ in dimensions d = 2, 3, and 4, respectively. We note that these are quite close to the rate we theoretically predict for the distance between the 32/33

Conclusion

MMD/ KSD descent, SVGD can create "super samples"

Open questions/future work:

- improve our quantization bounds for MMD/KSD (dependence in dimension, Laplace kernel?)
- obtain quantization bounds for SVGD

Thank you !

References I

 Aistleitner, C. and Dick, J. (2015). Functions of bounded variation, signed measures, and a general Koksma-Hlawka inequality. *Acta Arith.*, 167(2):143–171.
 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

 Bach, F., Lacoste-Julien, S., and Obozinski, G. (2012). On the equivalence between herding and conditional gradient algorithms.
 In ICML 2012 International Conference on Machine Learning.

Chen, W. Y., Mackey, L., Gorham, J., Briol, F.-X., and Oates, C. J. (2018). Stein points.

International Conference on Machine Learning (ICML).

Chen, Y., Welling, M., and Smola, A. (2012). Super-samples from kernel herding. arXiv preprint arXiv:1203.3472.

References III

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.
- Gretton, A., Borgwardt, K., Rasch, M., Schölkopf, B., and Smola, A. (2006).
 A kernel method for the two-sample-problem.
 Advances in neural information processing systems, 19:513–520.
- Gretton, A., Borgwardt, K. M., Rasch, M. J., Schölkopf, B., and Smola, A. J. (2012).
 A kernel two-sample test. *JMLR*, 13.

References IV

Hyvärinen, A. and Dayan, P. (2005). Estimation of non-normalized statistical models by score matching. Journal of Machine Learning Research, 6(4).

 Korba, A., Aubin-Frankowski, P.-C., Majewski, S., and Ablin, P. (2021).
 Kernel Stein discrepancy descent. International Conference of Machine Learning.

Łatuszyński, K., Miasojedow, B., and Niemiro, W. (2013). Nonasymptotic bounds on the estimation error of mcmc algorithms.

Bernoulli, 19(5A):2033-2066.

References V

- 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.

Lu, Y. and Lu, J. (2020).

A universal approximation theorem of deep neural networks for expressing probability distributions.

Advances in Neural Information Processing Systems, 33.

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.

 Tolstikhin, I., Sriperumbudur, B. K., and Muandet, K. (2017).
 Minimax estimation of kernel mean embeddings. *The Journal of Machine Learning Research*, 18(1):3002–3048. Alternative assumption for the MMD bound:

A2. Let $k(x, y) = \eta(x - y)$ a translation invariant kernel on \mathbb{R}^d . Assume that $\eta \in C(\mathbb{R}^d) \cap L^1(\mathbb{R}^d)$, and that its Fourier transform verifies : $\exists C_{k,d} \ge 0$ such that $(1 + |\xi|^2)^d \le C_{k,d} |\hat{\eta}(\xi)|^{-1}$ for any $\xi \in \mathbb{R}^d$.

Laplace kernel $k(x, y) = \exp(-||x - y||)$ corresponds to j = (d + 1)/2.



Figure: Quantization rates of the algorithms at study when the target distribution is a 2D-Gaussian mixture distribution with variance 0.3, centred at [1,0] and [-1,0]. We evaluate them using MMD and KSD with bandwidth 1. We run algorithms under the same setting as the 2-4D experiments on Figure 29.

L-BFGS

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

$$x_{l+1} = x_l - \gamma_l B_l^{-1} \nabla F(x_l) := x_l + \gamma_l d_l$$
(1)

where B_l^{-1} is a p.s.d. matrix approximating the inverse Hessian at x_l . Step1. (requires ∇F) It computes a cheap version of d_l based on BFGS recursion:

$$B_{l+1}^{-1} = \left(I - \frac{\Delta x_l y_l^T}{y_l^T \Delta x_l}\right) B_l^{-1} \left(I - \frac{y_l \Delta x_l^T}{y_l^T \Delta x_l}\right) + \frac{\Delta x_l \Delta x_l^T}{y_l^T \Delta x_l}$$

where
$$\Delta x_l = x_{l+1} - x_l$$

 $y_l = \nabla F(x_{l+1}) - \nabla F(x_l)$

Step2. (requires *F* and ∇F) A line-search is performed to find the best step-size in (1) :

$$F(x_l + \gamma_l d_l) \leq F(x_l) + c_1 \gamma_l \nabla F(x_l)^T d_l$$
$$\nabla F(x_l + \gamma_l d_l)^T d_l \geq c_2 \nabla F(x_l)^T d_l$$

Kernel Herding (KH) and Stein Points (SP)

They attempt to solve MMD or KSD quantization in a greedy manner, i.e. by sequentially constructing μ_n , adding one new particle at each iteration to minimize MMD/KSD.

Kernel Herding (KH) for the MMD [Chen et al., 2012]:

$$x^{n+1} = \underset{x \in \mathbb{R}^d}{\operatorname{argmax}} \langle w_n, k(x, .) \rangle_{\mathcal{H}_k}$$
$$w_{n+1} = w_n + m_\pi - k(x_{n+1}, .)$$

[Bach et al., 2012] obtain a linear rate of convergence $\mathcal{O}(e^{-bn})$

- if the mean embedding m_π = E_{x∼π}[k(x,.)] lies in the relative interior of the marginal polytope *convexhull*({k(x,.), x ∈ ℝ^d}) with distance *b* away from the boundary
- however for infinite-dimensional kernels b = 0 and the rate does not hold.

Stein Points for the KSD [Chen et al., 2018] greedily minimizes the KSD similarly. The authors establish a $\mathcal{O}((\log(n)/n)^{\frac{1}{2}})$ rate, which seem slower than their empirical observations.

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.

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.

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.