

Info Geom & Grad Flow Reading Group

19/09/2024

Paper: Stein Variational Gradient Descent: A General Purpose Bayesian Inference Algorithm.
Lin Qiang, Wang Dili (NeurIPS 2016)

Motivation

Compute $\int f(x) p(x) dx$.
probability density function

- ① Intractable integral \rightarrow no (easy) analytic solution
- ② (Potentially) high-dimensional $x \rightarrow$ grid method fails, curse of dimensionality

Example: posterior mean in Bayesian inference.

Notice that $\int f(x) p(x) dx = \mathbb{E}_{x \sim p}[f(x)]$.

exact \swarrow
 $\approx \frac{1}{N} \sum_{i=1}^N f(x_i^e)$
with $x_1^e, x_2^e, \dots, x_N^e \sim p$.
(not necessarily independent)

e.g. MCMC, SMC

approximate \swarrow
 $\approx \frac{1}{N} \sum_{i=1}^N f(x_i^a)$
with $x_1^a, x_2^a, \dots, x_N^a \sim \tilde{p} \approx p$

e.g. ABC, VI

SVGd \subset Particle VI \subset VI.

Get some particles. Move them around. Use their empirical distribution \tilde{p} to approximate p .

↳ Key question.

Kernel Stein Discrepancy

Sometimes, we wish to measure the similarity of two distributions p, q .

e.g. $KL(p \parallel q)$, $TV(p \parallel q)$, $W^2(p, q)$.

① couldn't be computed easily. \rightarrow not practical

② require knowing the densities \rightarrow we often only have samples ...

Applications: Goodness of fit, Variational inference.

One computable measure of distribution distance is the Kernel Stein Discrepancy (KSD)

$$KSD^2(p \parallel q) = \sup_{f \in \mathcal{H}_K} \left| \mathbb{E}_{x \sim q} [S_p f(x)] \right|^2 = \mathbb{E}_{x, r \sim q} [k_p(x, r)] \approx \frac{1}{n} \sum_{i, j=1}^n k_p(x_i, x_j).$$

unit ball in RKHS

$$\begin{aligned} f &\in \mathcal{H}_K \\ \|f\|_{\mathcal{H}_K} &\leq 1 \end{aligned}$$

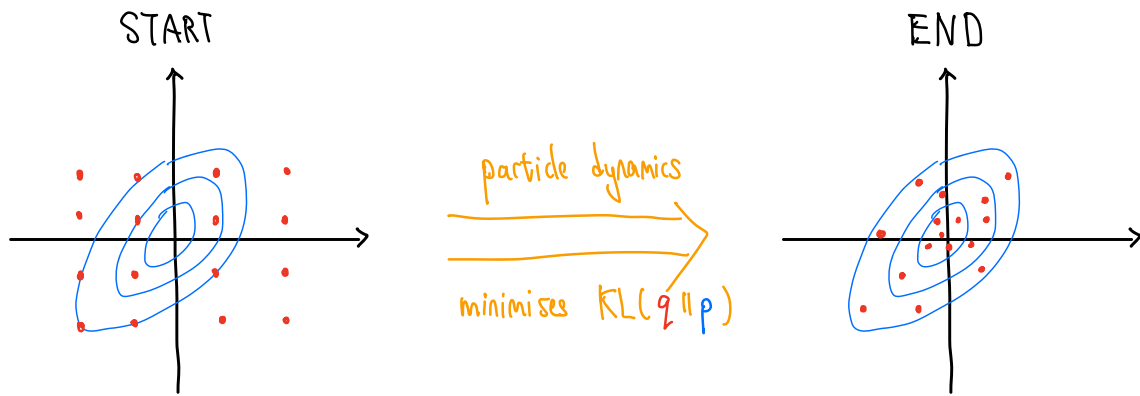
(Langevin) Stein operator

$$S_p f(x) = \langle \nabla_x \log p(x), f(x) \rangle + \langle \nabla, f(x) \rangle$$

$\nabla_x \log p(x)$, $\nabla_y \log p(y)$ (Stein) score function

$$k_p(x, y) = s_p(x) s_p(y) k(x, y) + s_p(x) \nabla_y k(x, y) + s_p(y) \nabla_x k(x, y) + \text{tr}(\nabla_x \nabla_y k(x, y)).$$

SVGD



Key Result (THM 3.1 & Lemma 3.2 of Liu & Wang (2016)).

Consider updates of the form $x_i \leftarrow T(x_i) = x_i + \varepsilon \phi(x_i)$. q denote the empirical distribution of the particles $\{x_i\}$. $q_{[T]}$ denote the empirical distribution after moving the particles by T . p is the target.

Then,

$$\nabla_{\varepsilon} \text{KL}(q_{[T]} \| p) \Big|_{\varepsilon=0} = - \mathbb{E}_{x \sim q} [\text{tr}(S_p \phi(x))].$$

If we force $\phi \in \mathcal{H}_k$ & $\|\phi\|_{\mathcal{H}_k} \leq 1$, i.e. ϕ is in the unit ball of RKHS, then we have

$$\begin{aligned} \arg \max_{\substack{\phi \in \mathcal{H}_k \\ \|\phi\|_{\mathcal{H}_k} \leq 1}} \nabla_{\varepsilon} \text{KL}(q_{[T]} \| p) \Big|_{\varepsilon=0} &= \phi^*(x') = \mathbb{E}_{x' \sim q} [S_p k(x, x')]. \\ &= \mathbb{E}_{x' \sim q} [k(x, x')^{\top} \nabla_{x'} \log p(x') + \nabla_{x'} k(x', x)]. \end{aligned}$$

Algorithm 1 SVGD

Require: Target distribution p . Initial particles $\{x_i^0\}_{i=1}^n$. Kernel k . Step size ε .

- 1: **for** $l = 0, 1, \dots$ **do**
- 2: **for** $i = 1, 2, \dots, n$ **do**
- 3:

$$x_i^{l+1} \leftarrow x_i^l + \frac{\varepsilon}{n} \sum_{j=1}^n \left[k(x_j^l, x_i^l) \nabla_{x_j^l} \log p(x_j^l) + \nabla_{x_j^l} k(x_j^l, x_i^l) \right].$$

- 4: **end for**
 - 5: **end for**
-