

A Convergence Theory for SVGD in the Population Limit under Talagrand's Inequality T1



Adil Salim



Lukang Sun



Peter Richtárik

Microsoft Research, USA. KAUST, Saudi Arabia

ICML 2022

Sampling framework

- ▶ F smooth and nonconvex

$$\mu^*(x) \propto \exp(-F(x))$$

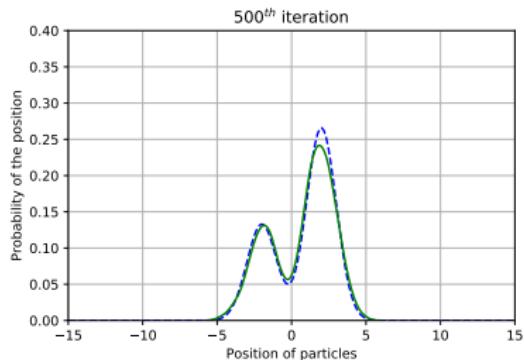
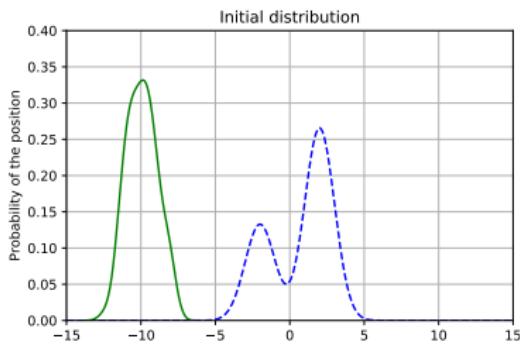


Figure: Simulation from [KSA⁺20] (Code from Q. Liu)

Sampling as optimization

$$\mu^* = \arg \min_{\mu \in \mathcal{P}(\mathbb{R}^d)} \text{KL}(\mu | \mu^*),$$

where

$$\text{KL}(\mu | \mu^*) := \int \log \left(\frac{d\mu}{d\mu^*}(x) \right) d\mu(x) \text{ if } \mu \ll \mu^*, +\infty \text{ else,}$$

satisfies $\text{KL}(\mu | \mu^*) \geq 0$ and $\text{KL}(\mu | \mu^*) = 0$ iff $\mu = \mu^*$.

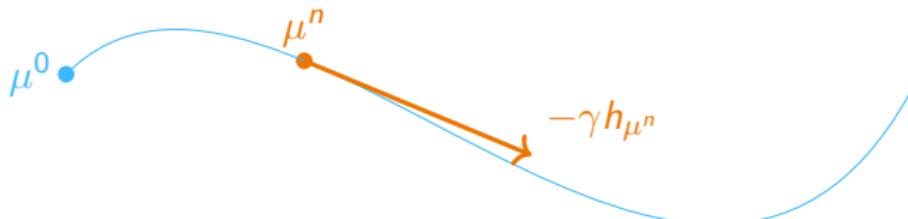
Gradient descent for $\text{KL}(\cdot | \mu^*)$?

Stein Variational Gradient Descent (SVGD)

$$\mu^{n+1} = (I - \gamma h_{\mu^n}) \# \mu^n, \quad (1)$$

where $h_\mu := \int \nabla F(x) k(\cdot, x) - \nabla_2 k(\cdot, x) d\mu(x)$ and $k(x, y)$ kernel of a RKHS H .

- ▶ When μ^0 is discrete then μ^n is discrete and (1) is called SVGD algorithm [Liu and Wang, 2016].
- ▶ When μ^0 has a density then μ^n has a density [SSR22] and (1) is called SVGD in the population limit.



Kernelized Stein Discrepancy

Remark: $\text{KSD}(\mu|\mu^*) := \|h_\mu\|_H$

[Liu et al., 2016, Chwialkowski et al., 2016, Oates et al., 2019, Gorham and Mackey, 2017].

If H rich enough, $\text{KSD}(\mu|\mu^*) = 0 \implies \mu = \mu^*$.

Analysis in the population limit

SVGD is a "Wasserstein gradient descent" in the metric of H
[Liu, 2017, Duncan et al., 2019, Chewi et al., 2020,
Nüsken and Renger, 2021, Shi et al., 2022, Gorham et al., 2020],
[KSA⁺20, SSR22].

Theorem 1 (Complexity and convergence of SVGD in the population limit [KSA⁺20, SSR22])

Convergence: If μ^0 has a density and if

$$\int \exp(\beta \|x - a\|^2) d\mu^*(x) < \infty, \text{ then } W_1(\mu^n, \mu^*) \rightarrow 0.$$

Moreover, complexity $n = \mathcal{O}(d^{3/2}/\varepsilon)$ to output μ s.t.

$$\text{KSD}^2(\mu|\mu^*) < \varepsilon.$$

Proof

$$\text{KL}(\mu^{n+1}|\mu^*) \leq \text{KL}(\mu^n|\mu^*) - \gamma(1 - \gamma C) \underbrace{\|h_{\mu^n}\|_H^2}_{=\text{KSD}^2(\mu^n|\mu^*)} .$$

References I

- [Chewi et al., 2020] Chewi, S., Gouic, T. L., Lu, C., Maunu, T., and Rigollet, P. (2020).
Svsgd as a kernelized wasserstein gradient flow of the chi-squared divergence.
In *Advances in Neural Information Processing Systems (NeurIPS)*.
- [Chwialkowski et al., 2016] Chwialkowski, K., Strathmann, H., and Gretton, A. (2016).
A kernel test of goodness of fit.
In *International Conference on Machine Learning (ICML)*, pages 2606–2615.
- [Duncan et al., 2019] Duncan, A., Nuesken, N., and Szpruch, L. (2019).
On the geometry of stein variational gradient descent.
arXiv preprint arXiv:1912.00894.
- [Gorham and Mackey, 2017] Gorham, J. and Mackey, L. (2017).
Measuring sample quality with kernels.
In *International Conference on Machine Learning (ICML)*, pages 1292–1301.
- [Gorham et al., 2020] Gorham, J., Raj, A., and Mackey, L. (2020).
Stochastic stein discrepancies.
Advances in Neural Information Processing Systems (NeurIPS), 33:17931–17942.
- [Liu, 2017] Liu, Q. (2017).
Stein variational gradient descent as gradient flow.
In *Advances in Neural Information Processing Systems (NIPS)*, volume 30.
- [Liu et al., 2016] Liu, Q., Lee, J., and Jordan, M. (2016).
A kernelized Stein discrepancy for goodness-of-fit tests.
In *International Conference on Machine Learning (ICML)*, pages 276–284.
- [Liu and Wang, 2016] Liu, Q. and Wang, D. (2016).
Stein variational gradient descent: A general purpose Bayesian inference algorithm.
In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 2378–2386.

References II

- [Nüsken and Renger, 2021] Nüsken, N. and Renger, D. (2021).
Stein variational gradient descent: many-particle and long-time asymptotics.
arXiv preprint arXiv:2102.12956.
- [Oates et al., 2019] Oates, C. J., Cockayne, J., Briol, F.-X., and Girolami, M. (2019).
Convergence rates for a class of estimators based on stein's method.
Bernoulli, 25(2):1141–1159.
- [Shi et al., 2022] Shi, J., Liu, C., and Mackey, L. (2022).
Sampling with mirrored stein operators.
In *International Conference on Learning Representations (ICLR)*.

Selected publications I

- [KSA⁺20] Anna Korba, Adil Salim, Michael Arbel, Giulia Luise, and Arthur Gretton.
A non-asymptotic analysis for Stein variational gradient descent.
In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- [SSR22] Adil Salim, Lukang Sun, and Peter Richtárik.
A convergence theory for SVGD in the population limit under Talagrand's inequality T1.
In *International Conference on Machine Learning (ICML)*, 2022.