
Appendix of “Stein Variational Gradient Descent as Gradient Flow”

Qiang Liu
 Department of Computer Science
 Dartmouth College
 Hanover, NH 03755
 qiang.liu@dartmouth.edu

A Density Evolution of SVGD Dynamics

A.1 Proof of Lemma 3.1

Proof. Recall that $\mathbf{g}(x, x') = \mathcal{S}_p^{x'} \otimes k(x', x)$, and $\phi_{\mu,p}^*(x) = \mathbb{E}_{x' \sim \mu}[\mathbf{g}(x, x')]$, we have $\mathbf{T}_{\mu,p}(x) = x + \epsilon \mathbb{E}_{x' \sim \mu}[\mathbf{g}(x, x')]$. Therefore,

$$\begin{aligned}
 \|\mathbf{T}_{\mu,p}\|_{\text{Lip}} &= \max_{x \neq y} \frac{\|\mathbf{T}_{\mu,p}(x) - \mathbf{T}_{\mu,p}(y)\|_2}{\|x - y\|_2} \\
 &= \max_{x \neq y} \frac{\|x - y + \epsilon \mathbb{E}_{x' \sim \mu}[\mathbf{g}(x, x') - \mathbf{g}(y, x')]\|_2}{\|x - y\|_2} \\
 &\leq 1 + \epsilon \|\mathbf{g}\|_{\text{Lip}}, \tag{A.1}
 \end{aligned}$$

and for $\forall x$,

$$\|\mathbf{T}_{\mu,p}(x) - \mathbf{T}_{\nu,p}(x)\|_2 = \epsilon \|\mathbb{E}_{x' \sim \mu} \mathbf{g}(x, x') - \mathbb{E}_{x' \sim \nu} \mathbf{g}(x, x')\|_2 \leq \epsilon \|\mathbf{g}\|_{\text{BL}} \text{BL}(\mu, \nu). \tag{A.2}$$

For any h with $\|h\|_{\text{BL}} = \max(\|h\|_{\infty}, \|h\|_{\text{Lip}}) \leq 1$, we have

$$\begin{aligned}
 &|\mathbb{E}_{\Phi_p(\mu)}[h] - \mathbb{E}_{\Phi_p(\nu)}[h]| \\
 &= |\mathbb{E}_{\mu}[h \circ \mathbf{T}_{\mu,p}] - \mathbb{E}_{\nu}[h \circ \mathbf{T}_{\nu,p}]| \\
 &\leq |\mathbb{E}_{\mu}[h \circ \mathbf{T}_{\mu,p}] - \mathbb{E}_{\nu}[h \circ \mathbf{T}_{\mu,p}]| + |\mathbb{E}_{\nu}[h \circ \mathbf{T}_{\mu,p}] - \mathbb{E}_{\nu}[h \circ \mathbf{T}_{\nu,p}]|.
 \end{aligned}$$

We just need to bound these two terms. For the first term,

$$\begin{aligned}
 |\mathbb{E}_{\mu}[h \circ \mathbf{T}_{\mu,p}] - \mathbb{E}_{\nu}[h \circ \mathbf{T}_{\mu,p}]| &\leq \|h \circ \mathbf{T}_{\mu,p}\|_{\text{BL}} \text{BL}(\mu, \nu) \\
 &\leq \max(\|h\|_{\infty}, \|h\|_{\text{Lip}} \|\mathbf{T}_{\mu,p}\|_{\text{Lip}}) \text{BL}(\mu, \nu) \\
 &\leq (1 + \epsilon \|\mathbf{g}\|_{\text{Lip}}) \text{BL}(\mu, \nu), \quad \text{//by Equation A.1.}
 \end{aligned}$$

For the second term,

$$\begin{aligned}
 |\mathbb{E}_{\nu}[h \circ \mathbf{T}_{\mu,p}] - \mathbb{E}_{\nu}[h \circ \mathbf{T}_{\nu,p}]| &\leq \max_x |h \circ \mathbf{T}_{\mu,p}(x) - h \circ \mathbf{T}_{\nu,p}(x)| \\
 &\leq \|h\|_{\text{Lip}} \max_x \|\mathbf{T}_{\mu,p}(x) - \mathbf{T}_{\nu,p}(x)\|_2 \\
 &\leq \epsilon \|\mathbf{g}\|_{\text{BL}} \text{BL}(\mu, \nu), \quad \text{//by Equation A.2.}
 \end{aligned}$$

Therefore,

$$\text{BL}(\Phi_p(\mu), \Phi_p(\nu)) \leq (1 + \epsilon \|\mathbf{g}\|_{\text{Lip}} + \epsilon \|\mathbf{g}\|_{\text{BL}}) \text{BL}(\mu, \nu) \leq (1 + 2\epsilon \|\mathbf{g}\|_{\text{BL}}) \text{BL}(\mu, \nu).$$

□

A.2 Proof of Theorem 3.3

Proof. Denote by $\mu_\ell = \mu_\ell^\infty$ for notation convenience.

$$\begin{aligned}
& \text{KL}(\mu_{\ell+1} \parallel \nu_p) - \text{KL}(\mu_\ell \parallel \nu_p) \\
&= \text{KL}(\mathbf{T}_{\mu_\ell, p} \mu_\ell \parallel \nu_p) - \text{KL}(\mu_\ell \parallel \nu_p) \\
&= \text{KL}(\mu_\ell \parallel \mathbf{T}_{\mu_\ell, p}^{-1} \nu_p) - \text{KL}(\mu_\ell \parallel \nu_p) \quad // \text{by Lemma A.2} \\
&= -\mathbb{E}_{x \sim \mu_\ell} [\log p(\mathbf{T}_{\mu_\ell, p}(x)) + \log \det(\nabla \mathbf{T}_{\mu_\ell, p}(x)) - \log p(x)]. \tag{A.3}
\end{aligned}$$

Note that $\mathbf{T}_{\mu_\ell, p}(x) = x + \epsilon \phi_{\mu_\ell, p}^*(x)$. We have the follow version of Taylor approximation:

$$\log p(x) - \log p(\mathbf{T}_{\mu_\ell, p}(x)) \leq -\epsilon \nabla_x \log p(x)^\top \phi_{\mu_\ell, p}^*(x) + \frac{\epsilon^2}{2} \|\nabla \log p\|_{\text{Lip}} \cdot \|\phi_{\mu_\ell, p}^*\|_2^2. \tag{A.4}$$

This is because, defining $x_s = x + s\epsilon \phi_{\mu_\ell, p}^*(x)$, $\forall s \in [0, 1]$,

$$\begin{aligned}
& \log p(x) - \log p(\mathbf{T}_{\mu_\ell, p}(x)) \\
&= -\int_0^1 \nabla_s \log p(x_s) ds \\
&= -\int_0^1 \nabla_x \log p(x_s)^\top (\epsilon \phi_{\mu_\ell, p}^*(x)) ds \\
&= -\epsilon \nabla_x \log p(x)^\top \phi_{\mu_\ell, p}^*(x) - \int_0^1 (\nabla_x \log p(x_s) - \nabla_x \log p(x))^\top (\epsilon \phi_{\mu_\ell, p}^*(x)) ds \\
&\leq -\epsilon \nabla_x \log p(x)^\top \phi_{\mu_\ell, p}^*(x) + \epsilon^2 \|\nabla \log p\|_{\text{Lip}} \cdot \|\phi_{\mu_\ell, p}^*(x)\|_2^2 \int_0^1 s ds \\
&= -\epsilon \nabla_x \log p(x)^\top \phi_{\mu_\ell, p}^*(x) + \frac{\epsilon^2}{2} \|\nabla \log p\|_{\text{Lip}} \cdot \|\phi_{\mu_\ell, p}^*(x)\|_2^2.
\end{aligned}$$

where we used the fundamental theorem of calculus, which holds for weakly differentiable functions [20, Theorem 3.60, page 77]. In addition, Take $B = \nabla \phi_{\mu_\ell, p}^*(x)$ in bound (A.9) of Lemma A.1, and take $\epsilon < 1/(2\rho(B + B^\top))$, we have

$$\begin{aligned}
\log |\det(\nabla \mathbf{T}_{\mu_\ell, p}(x))| &\geq \epsilon \text{tr}(\nabla \phi_{\mu_\ell, p}^*(x)) - 2\epsilon^2 \|\nabla \phi_{\mu_\ell, p}^*(x)\|_F^2 \\
&= \epsilon \nabla \cdot \phi_{\mu_\ell, p}^*(x) - 2\epsilon^2 \|\nabla \phi_{\mu_\ell, p}^*(x)\|_F^2. \tag{A.5}
\end{aligned}$$

Combining (A.4) and (A.5) gives

$$\begin{aligned}
\text{KL}(\mu_{\ell+1} \parallel \nu_p) - \text{KL}(\mu_\ell \parallel \nu_p) &\leq -\epsilon \mathbb{E}_{\mu_\ell} [\mathcal{S}_p \phi_{\mu_\ell, p}^*] + \Delta \\
&= -\epsilon \mathbb{D}(\mu_\ell \parallel \nu_p)^2 + \Delta,
\end{aligned}$$

where Δ is a residual term:

$$\Delta = \epsilon^2 \mathbb{E}_{x \sim \mu_\ell} \left[\frac{1}{2} \|\nabla \log p\|_{\text{Lip}} \cdot \|\phi_{\mu_\ell, p}^*(x)\|_2^2 + 2 \|\nabla \phi_{\mu_\ell, p}^*(x)\|_F^2 \right]$$

We need to bound $\|\phi_{\mu_\ell, p}^*(x)\|_2$ and $\|\nabla \phi_{\mu_\ell, p}^*(x)\|_F$. This can be done using the reproducing property:

let $\phi_{\mu_\ell, p}^* = [\phi_1, \dots, \phi_d]^\top$; recall that $\phi_i \in \mathcal{H}_0$ and $\phi_{\mu_\ell, p}^* \in \mathcal{H} = \mathcal{H}_0^d$, then

$$\phi_i(x) = \langle \phi_i(\cdot), k(x, \cdot) \rangle_{\mathcal{H}_0}, \quad \partial_{x_j} \phi_i(x) = \langle \phi_i(\cdot), \partial_{x_j} k(x, \cdot) \rangle_{\mathcal{H}_0}, \quad \forall i, j = 1, \dots, d, x \in X.$$

Also note that $\|\phi_{\mu_\ell, p}^*\|_{\mathcal{H}}^2 = \sum_{i=1}^d \|\phi_i\|_{\mathcal{H}_0}^2 = \mathbb{D}(\mu_\ell \parallel \nu_p)^2$, we have by Cauchy-Swarchz inequality,

$$\begin{aligned}
\|\phi_{\mu_\ell, p}^*(x)\|_2^2 &= \sum_{i=1}^d \phi_i(x)^2 \\
&= \sum_{i=1}^d (\langle k(x, \cdot), \phi_i(\cdot) \rangle_{\mathcal{H}_0})^2 \\
&\leq \sum_i \|k(x, \cdot)\|_{\mathcal{H}_0}^2 \cdot \|\phi_i\|_{\mathcal{H}_0}^2 \\
&= k(x, x) \cdot \|\phi_{\mu_\ell, p}^*\|_{\mathcal{H}}^2 \\
&= k(x, x) \cdot \mathbb{D}(\mu_\ell \parallel \nu_p)^2,
\end{aligned}$$

and

$$\begin{aligned}
\|\nabla\phi_{\mu_\ell,p}^*(x)\|_F^2 &= \sum_{ij} \partial_{x_j}\phi_i(x)^2 \\
&= \sum_{ij} (\langle \partial_{x_j}k(x,\cdot), \phi_i(\cdot) \rangle_{\mathcal{H}_0})^2 \\
&\leq \sum_{ij} \|\partial_{x_j}k(x,\cdot)\|_{\mathcal{H}_0}^2 \cdot \|\phi_i\|_{\mathcal{H}_0}^2 \\
&= \sum_{ij} \partial_{x_j,x'_j}k(x,x')|_{x=x'} \cdot \|\phi_i\|_{\mathcal{H}_0}^2 \\
&= \nabla_{xx'}k(x,x) \cdot \|\phi_{\mu_\ell,p}^*\|_{\mathcal{H}}^2 \\
&= \nabla_{xx'}k(x,x) \cdot \mathbb{D}(\mu_\ell \parallel \nu_p)^2.
\end{aligned} \tag{A.6}$$

Therefore,

$$\begin{aligned}
\Delta &\leq \epsilon^2 \mathbb{D}(\mu_\ell \parallel \nu_p)^2 \left(\frac{1}{2} \mathbb{E}_{x \sim \mu_\ell} [\|\nabla \log p\|_{\text{Lip}} k(x,x) + 2\nabla_{xx'}k(x,x)] \right) \\
&= \epsilon^2 R \mathbb{D}(\mu_\ell \parallel \nu_p)^2.
\end{aligned}$$

This gives

$$\text{KL}(\mu_{\ell+1} \parallel \nu_p) - \text{KL}(\mu_\ell \parallel \nu_p) \leq -\epsilon(1 - \epsilon R) \mathbb{D}(\mu_\ell \parallel \nu_p)^2.$$

□

Lemma A.1. *Let B be a square matrix and $\|B\|_F = \sqrt{\sum_{ij} b_{ij}^2}$ its Frobenius norm. Let ϵ be a positive number that satisfies $0 \leq \epsilon < \frac{1}{\rho(B+B^\top)}$, where $\rho(\cdot)$ denotes the spectrum radius. Then $I + \epsilon(B + B^\top)$ is positive definite, and*

$$\log |\det(I + \epsilon B)| \geq \text{etr}(B) - \epsilon^2 \frac{\|B\|_F^2}{1 - \epsilon \rho(B + B^\top)}. \tag{A.7}$$

Therefore, take an even smaller ϵ such that $0 \leq \epsilon \leq \frac{1}{2\rho(B+B^\top)}$, we get

$$\log |\det(I + \epsilon B)| \geq \text{etr}(B) - 2\epsilon^2 \|B\|_F^2. \tag{A.8}$$

Proof. When $\epsilon < \frac{1}{\rho(B+B^\top)}$, we have $\rho(I + \epsilon(B + B^\top)) \geq 1 - \epsilon \rho(B + B^\top) > 0$, so $I + \epsilon(B + B^\top)$ is positive definite.

By the property of matrix determinant, we have

$$\begin{aligned}
\log |\det(I + \epsilon B)| &= \frac{1}{2} \log \det((I + \epsilon B)(I + \epsilon B)^\top) \\
&= \frac{1}{2} \log \det(I + \epsilon(B + B^\top) + \epsilon^2 BB^\top) \\
&\geq \frac{1}{2} \log \det(I + \epsilon(B + B^\top)),
\end{aligned} \tag{A.9}$$

where (A.10) holds because both $I + \epsilon(B + B^\top)$ and $\epsilon^2 BB^\top$ are positive semi-definite.

Let $A = B + B^\top$. We can establish

$$\log \det(I + \epsilon A) \geq \text{etr}(A) - \frac{\epsilon^2}{2} \frac{\|A\|_F^2}{1 - \epsilon \rho(A)}, \tag{A.10}$$

which holds for any symmetric matrix A and $0 \leq \epsilon < 1/\rho(A)$. This is because, assuming $\{\lambda_i\}$ are the eigenvalues of A ,

$$\begin{aligned}
\log \det(I + \epsilon A) - \epsilon \operatorname{tr}(A) &= \sum_i [\log(1 + \epsilon \lambda_i) - \epsilon \lambda_i] \\
&= \sum_i \left[\int_0^1 \frac{\epsilon \lambda_i}{1 + s \epsilon \lambda_i} ds - \epsilon \lambda_i \right] \\
&= - \sum_i \int_0^1 \frac{s \epsilon^2 \lambda_i^2}{1 + s \epsilon \lambda_i} ds \\
&\geq - \sum_i \frac{\epsilon^2 \lambda_i^2}{1 - \epsilon \max_i |\lambda_i|} \int_0^1 s ds \\
&\geq - \sum_i \frac{\epsilon^2 \lambda_i^2}{2(1 - \epsilon \max_i |\lambda_i|)} \\
&= - \frac{\epsilon^2}{2} \frac{\|A\|_F^2}{1 - \epsilon \rho(A)}.
\end{aligned}$$

Take $A = B + B^\top$ in (A.11) and combine it with (A.10), we get

$$\begin{aligned}
\log |\det(I + \epsilon B)| &\geq \frac{1}{2} \log \det(I + \epsilon(B + B^\top)) \\
&\geq \frac{\epsilon}{2} \operatorname{tr}(B + B^\top) - \frac{\epsilon^2}{4} \frac{\|B + B^\top\|_F^2}{1 - \epsilon \rho(B + B^\top)} \\
&\geq \epsilon \operatorname{tr}(B) - \epsilon^2 \frac{\|B\|_F^2}{1 - \epsilon \rho(B + B^\top)},
\end{aligned}$$

where we used the fact that $\operatorname{tr}(B) = \operatorname{tr}(B^\top)$ and $\|B + B^\top\|_F \leq \|B\|_F + \|B^\top\|_F = 2\|B\|_F$. \square

Lemma A.2. *Let T be a one-to-one map, and μ and ν two probability measures. We have*

$$\operatorname{KL}(T\mu \parallel \nu) = \operatorname{KL}(\mu \parallel \nu),$$

given that the KL divergence between μ and ν exists.

Proof. We prove this for f -divergence in general, which includes KL divergence as a special case. Given a convex function f such that $f(1) = 0$, the f -divergence is defined

$$D_f(\mu \parallel \nu) = \mathbb{E}_\nu \left[f \left(\frac{d\mu}{d\nu} \right) \right].$$

Assume f^* is the convex conjugate of f , we have a variational representation for f -divergence:

$$D_f(\mu \parallel \nu) = \sup_g \left\{ \mathbb{E}_\mu[g(x)] - \mathbb{E}_\nu[f^*(g(x))] \right\},$$

where g is over the set of all measurable functions. Therefore, we have

$$\begin{aligned}
D_f(T\mu \parallel \nu) &= \sup_g \left\{ \mathbb{E}_\mu[g \circ T(x)] - \mathbb{E}_\nu[f^*(g(x))] \right\} \\
&= \sup_{\tilde{g}} \left\{ \mathbb{E}_\mu[\tilde{g}(x)] - \mathbb{E}_\nu[f^*(\tilde{g} \circ T^{-1}(x))] \right\} \quad // \text{Define } \tilde{g} = g \circ T. \\
&= D_f(\mu \parallel T^{-1}\nu).
\end{aligned}$$

\square

A.3 Proof of Fokker-Planck Equation (13)

Proof. Recall that $T_{\mu,p}(x) = x + \epsilon \phi_{\mu,p}^*(x)$ and we denote by q the density of measure μ . Assume ϵ is sufficiently small so that $\nabla T_{\mu,p}(x) = I + \epsilon \nabla \phi_{\mu,p}^*(x)$ is positive definite (See Lemma A.1). By the implicit function theorem, we have

$$T_{\mu,p}^{-1}(x) = x - \epsilon \phi_{\mu,p}^*(x) + o(\epsilon).$$

Therefore We have

$$\begin{aligned}
\log q'(x) &= \log q(\mathbf{T}_{\mu,p}^{-1}(x)) + \log \det(\nabla_x \mathbf{T}_{\mu,p}^{-1}(x)) \\
&= \log q(x - \epsilon \cdot \phi_{\mu,p}^*(x)) + \log \det(I - \epsilon \nabla_x \phi_{\mu,p}^*(x)) + o(\epsilon) \\
&= \log q(x) - \epsilon \nabla_{x_i} \log q(x)^\top \phi_{\mu,p}^*(x) - \epsilon q(x) \cdot \text{tr}(\nabla_x \phi_{\mu,p}^*(x)) + o(\epsilon) \\
&= \log q(x) - \epsilon \mathcal{S}_q \phi_{\mu,p}^*(x) + o(\epsilon).
\end{aligned}$$

Therefore,

$$\begin{aligned}
\frac{q'(x) - q(x)}{\epsilon} &= \frac{q(\log q(x) - \log q(x))}{\epsilon} + o(\epsilon) \\
&= -q(x) \mathcal{S}_q \phi_{q\ell,p}^*(x) + o(\epsilon) \\
&= -\nabla \cdot (\phi_{q\ell,p}^*(x) q_\ell(x)) + o(\epsilon).
\end{aligned}$$

Taking $\epsilon \rightarrow 0$ gives the result. \square

A.4 Proof of Theorem 3.5

Proof. Since $q' = q + qf dt$ is equivalent to transforming the variable by $\mathbf{T}(x) = x + \psi_{q,f} dt$, the corresponding change on KL divergence is

$$\begin{aligned}
F(q + qf dt) &= F(q) + \mathbb{E}_q[\mathcal{S}_p \psi_{q,f}] dt \\
&= F(q) + \langle \phi_{q,p}^*, \psi_{q,f} \rangle_{\mathcal{H}} dt \\
&= F(q) + \langle \nabla \cdot (\phi_{q,p}^* q), \nabla \cdot (\psi_{q,f} q) \rangle_{q\mathcal{H}_q} dt \\
&= F(q) + \langle \nabla \cdot (\phi_{q,p}^* q), qf \rangle_{q\mathcal{H}_q} dt
\end{aligned}$$

This proves that $\nabla \cdot (\phi_{q,p}^* q)$ is the covariant functional gradient. \square

References

- [1] Q. Liu and D. Wang. Stein variational gradient descent: A general purpose bayesian inference algorithm. In *Advances in Neural Information Processing Systems*, 2016.
- [2] M. J. Wainwright, M. I. Jordan, et al. Graphical models, exponential families, and variational inference. *Foundations and Trends® in Machine Learning*, 1(1–2):1–305, 2008.
- [3] Y. Chen, M. Welling, and A. Smola. Super-samples from kernel herding. In *Conference on Uncertainty in Artificial Intelligence (UAI)*, 2010.
- [4] J. Dick, F. Y. Kuo, and I. H. Sloan. High-dimensional integration: the quasi-monte carlo way. *Acta Numerica*, 22:133–288, 2013.
- [5] B. Dai, N. He, H. Dai, and L. Song. Provable Bayesian inference via particle mirror descent. In *The 19th International Conference on Artificial Intelligence and Statistics*, 2016.
- [6] C. Stein. Approximate computation of expectations. *Lecture Notes-Monograph Series*, 7:i–164, 1986.
- [7] Q. Liu, J. D. Lee, and M. I. Jordan. A kernelized Stein discrepancy for goodness-of-fit tests and model evaluation. In *International Conference on Machine Learning (ICML)*, 2016.
- [8] K. Chwialkowski, H. Strathmann, and A. Gretton. A kernel test of goodness-of-fit. In *International Conference on Machine Learning (ICML)*, 2016.
- [9] C. J. Oates, M. Girolami, and N. Chopin. Control functionals for Monte Carlo integration. *Journal of the Royal Statistical Society, Series B*, 2017.
- [10] J. Gorham and L. Mackey. Measuring sample quality with kernels. In *International Conference on Machine Learning (ICML)*, 2017.
- [11] C. J. Oates, J. Cockayne, F.-X. Briol, and M. Girolami. Convergence rates for a class of estimators based on Stein’s identity. *arXiv preprint arXiv:1603.03220*, 2016.

- [12] W. Braun and K. Hepp. The Vlasov dynamics and its fluctuations in the $1/n$ limit of interacting classical particles. *Communications in mathematical physics*, 56(2):101–113, 1977.
- [13] A. A. Vlasov. On vibration properties of electron gas. *J. Exp. Theor. Phys*, 8(3):291, 1938.
- [14] H. Spohn. *Large scale dynamics of interacting particles*. Springer Science & Business Media, 2012.
- [15] P. Del Moral. *Mean field simulation for Monte Carlo integration*. CRC Press, 2013.
- [16] A. Berline and C. Thomas-Agnan. *Reproducing kernel Hilbert spaces in probability and statistics*. Springer Science & Business Media, 2011.
- [17] F. Otto. The geometry of dissipative evolution equations: the porous medium equation. 2001.
- [18] C. Villani. *Optimal transport: old and new*, volume 338. Springer Science & Business Media, 2008.
- [19] J. Han and Q. Liu. Stein variational adaptive importance sampling. In *Uncertainty in Artificial Intelligence*, 2017.
- [20] J. K. Hunter. Notes on partial differential equations. 2014. URL https://www.math.ucdavis.edu/~hunter/m218a_09/pde_notes.pdf.