



Iterated Schrödinger Bridge Approximation to Wasserstein Gradient Flows



Medha Agarwal, Zaid Harchaoui, Garrett Mulcahy, Soumik Pal



Background

For initial d -dimensional particles/tokens $(x_1, \dots, x_n) \sim \rho_0$, the Transformer update is

$$x_i \leftarrow \sum_{j=1}^n K_{ij}^\infty W_V x_j, \text{ where } K^\infty = \text{Sinkhorn}(C) \text{ and } C_{ij} = (W_Q x_i)^\top (W_K x_j)$$

Value Matrix

Query Matrix

Key Matrix

Under $(W_Q)^\top W_K = -W_V = I_d$, the infinite particles counterpart with bandwidth/temperature ε is

$$\rho \leftarrow (T_{\rho, \varepsilon})_\# \rho \text{ where } T_{\rho, \varepsilon}(x) = 2x - \int k_\varepsilon^\infty(x, y) y d\rho(y) \text{ and } k_\varepsilon = \text{Sinkhorn}(c/\varepsilon)$$

Iteratively, with $\rho_0^\varepsilon = \rho_0$,

Sequence $(\rho_k^\varepsilon, k \in [T/\varepsilon])$

$$\rho_k^\varepsilon = (T_{\rho_{k-1}^\varepsilon, \varepsilon})_\# \rho_{k-1}^\varepsilon$$

Piece-wise continuous curve $(\rho_t^\varepsilon, t \in [0, T])$

$$\rho_t^\varepsilon = \rho_{\lfloor t/\varepsilon \rfloor}^\varepsilon, t \in [0, T]$$

As $\varepsilon \rightarrow 0+$, what is the limit of the curve $(\rho_t^\varepsilon, t \in [0, T])$? Answer: Heat equation!

In fact, if $(\rho, t \in [0, T])$ is heat equation starting from ρ_0 , $\lim_{\varepsilon \rightarrow 0} \sup_{k \in [T/\varepsilon]} \mathbb{W}_2^2(\rho_k^\varepsilon, \rho_{k\varepsilon}) = 0$

Setting

Consider the entropy regularized optimal transport (EOT) problem between measures μ and ν

$$\min_{\pi \in \Pi(\mu, \nu)} \{ \|x - y\|^2 d\pi + \varepsilon \text{KL}(\pi \mid \mu \otimes \nu) \}$$

argmin $\pi_{\mu, \nu, \varepsilon}$ is the Schrödinger bridge

Consider same marginal setting, i.e. $\mu = \nu = \rho$ with Schrödinger bridge $\pi_{\rho, \varepsilon}$, then

$$\mathcal{B}_{\rho, \varepsilon}(x) = \mathbb{E}_{\pi_{\rho, \varepsilon}}[Y \mid X = x] = \int y k_\varepsilon^\infty(x, y) d\rho(y)$$

Therefore, for infinite particles, the self-attention update is

$$\rho_k^\varepsilon = (T_{\rho_{k-1}^\varepsilon, \varepsilon})_\# \rho_{k-1}^\varepsilon = (2I_d - \mathcal{B}_{\rho, \varepsilon})_\# \rho_{k-1}^\varepsilon$$

Main contribution: Under suitable conditions, $2x - \mathcal{B}_{\rho, \varepsilon} \approx x - \frac{\varepsilon}{2} \nabla \log \rho(x)$

Claim: $(\rho_k^\varepsilon, k \in [T/\varepsilon])$ is a discrete approximate to the heat flow.

Wasserstein gradient flow is characterized by the continuity equation $\partial_t \rho_t + \nabla \cdot (v_t \rho_t) = 0$ where $v_t: \mathbb{R}^d \rightarrow \mathbb{R}^d$ is the velocity field. For heat flow, $v_t = -0.5 \nabla \log \rho_t$.

What about any gradient flow? Can we such approximations for them like self-attention for heat flow?

Key Contribution: (1) Under suitable conditions, self-attention iterations converge to the heat flow as $\varepsilon \rightarrow 0$. (2) We write the iterated scheme for a general functional \mathcal{F} . (3) We prove the uniform convergence of the scheme to the gradient flow of KL divergence with respect to a log-concave density.

Iterative Scheme

Consider absolutely continuous curve $(\rho_t, t \in [0, T])$ satisfying $\partial_t \rho_t + \nabla \cdot (v_t \rho_t) = 0$ where $v_t = v(\rho_t)$ is the velocity field. An example is gradient flow minimizing a functional $\mathcal{F}: \mathcal{P}(\mathbb{R}^d) \rightarrow \mathbb{R}$.

One-step Approximation

One discrete time approximation is analogous to explicit Euler scheme

$$S_\varepsilon^1(\rho) = (id + \varepsilon v)_\# \rho \text{ with } v = v(\rho).$$

Explicit Euler (EE) Step

Let $v = \nabla \psi$ for a smooth function ψ and there exists $\theta_\varepsilon \in \mathbb{R} \setminus \{0\}$ s.t. $\Lambda(\varepsilon) = \int_{\mathbb{R}^d} \exp(2\theta_\varepsilon \psi) < +\infty$.

Define the surrogate measure $\sigma_\varepsilon := \Lambda(\varepsilon)^{-1} \exp(2\theta_\varepsilon \psi)$ and the one-step update

$$SB_\varepsilon^1(\rho) := ((1 - \theta_\varepsilon^{-1}) id + \theta_\varepsilon^{-1} \mathcal{B}_{\sigma_\varepsilon, \varepsilon})_\# \rho$$

Schrödinger Bridge (SB) Step

Iterative Schemes

Explicit Euler Scheme

$$(\rho, S_\varepsilon^1(\rho), \dots, S_\varepsilon^{N_\varepsilon}(\rho)), N_\varepsilon = \lceil T/\varepsilon \rceil.$$

Schrödinger Bridge Scheme

$$(\rho, SB_\varepsilon^1(\rho), \dots, SB_\varepsilon^{N_\varepsilon}(\rho)), N_\varepsilon = \lceil T/\varepsilon \rceil.$$

Remark: For heat equation (gradient flow of entropy functional), $\sigma_\varepsilon = \rho$ and $\theta_\varepsilon = -1$. For Fokker Planck equation (gradient flow of KL divergence with respect to ν), $\sigma_\varepsilon = (\rho/\nu)^{-\theta_\varepsilon}$ where the sign of θ_ε depends on the integrability of σ_ε .

Convergence

Theorem 1 (Tight Approximation of Same Marginal Schrödinger Bridge)

Let $\rho = e^{-\mathcal{G}} \in \mathcal{P}(\mathbb{R}^d)$ with enough regularity such that there is strong solution to the Langevin SDE $dX_t = -\frac{1}{2} \nabla g(X_t) + dB_t$ with $X_0 \sim \rho$. Let $\ell_{\rho, \varepsilon} = \text{Law}(X_1, X_\varepsilon)$, then $H(\ell_{\rho, \varepsilon} \mid \pi_{\rho, \varepsilon}) + H(\pi_{\rho, \varepsilon} \mid \ell_{\rho, \varepsilon}) \leq \varepsilon^2 C(\varepsilon)$.

For Gaussian marginals, $C(\varepsilon) = \mathcal{O}(\varepsilon^2)$

One-step Convergence

Theorem 2 (Single Step Convergence)

If the surrogate measure σ_ε satisfies a set of regularity conditions, then there exists a constant $K > 0$ such that $\mathbb{W}_2(S_\varepsilon^1(\rho), SB_\varepsilon^1(\rho)) < K\varepsilon C(\varepsilon)$. Under assumptions on surrogate measure σ_ε , $C(\varepsilon) = o(1)$.

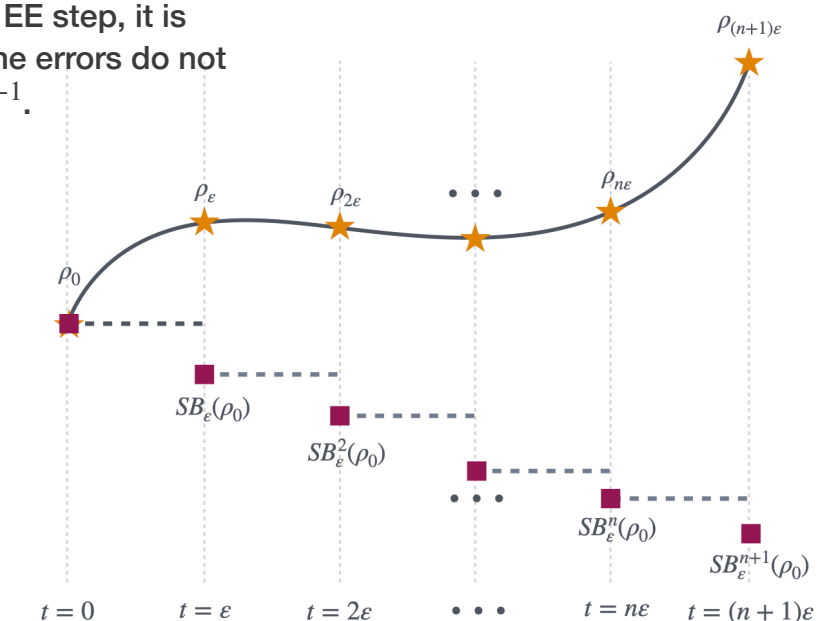
$$\lim_{\varepsilon \rightarrow 0} \varepsilon^{-1} \mathbb{W}_2(S_\varepsilon^1(\rho), SB_\varepsilon^1(\rho)) = 0$$

Remark: (1) The one step convergence relies on the close approximation of the same-marginal Schrödinger bridge by the Langevin diffusion.

Uniform Convergence

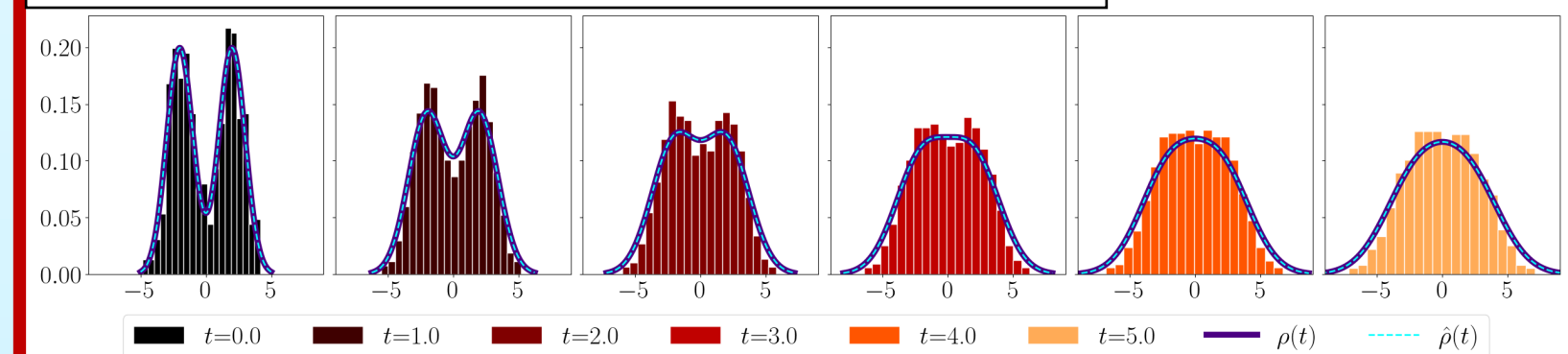
$$\mathbb{W}_2(\rho_{k\varepsilon}, SB_\varepsilon^k(\rho)) \leq \mathbb{W}_2(\rho_{k\varepsilon}, SB_\varepsilon^1(\rho_{(k-1)\varepsilon})) + \mathbb{W}_2(SB_\varepsilon^1(\rho_{(k-1)\varepsilon}), SB_\varepsilon^k(\rho))$$

Even though, we show that SB one step is $o(\varepsilon)$ approximation of the EE step, it is important to ensure that the errors do not culminate by more than ε^{-1} .



Experiments

Heat flow with $\varepsilon = 0.01$ and $\rho_0 = 0.5\mathcal{N}(-2, 1) + 0.5\mathcal{N}(2, 1)$



Ornstein-Uhlenbeck process with $\varepsilon = 0.01$ and $\rho_0 = \mathcal{N}(0, 4)$

