

Self-attention mechanism

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 and $k_\varepsilon = \text{Sinkhorn}(C/\varepsilon)$ is

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

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

$$\text{Sequence } (\rho_k^\varepsilon, k \in [T/\varepsilon]) \rightarrow \text{Continuous curve } (\rho_t^\varepsilon, t \in [0, T])$$

$$\rho_k^\varepsilon = \left(T_{\rho_{k-1}^\varepsilon} \right)_\# \rho_{k-1}^\varepsilon \rightarrow \rho_t^\varepsilon = \rho_{\lfloor t/\varepsilon \rfloor}^\varepsilon, t \in [0, T]$$

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

Answer: Heat equation!

In fact, if $(\rho_t, t \in [0, T])$ is heat equation starting from ρ_0 ,

$$\limsup_{\varepsilon \rightarrow 0} \mathbb{W}_2^2(\rho_k^\varepsilon, \rho_{k\varepsilon}) = 0$$

Background

Entropy-regularized optimal transport (EOT)

Consider the EOT problem between measures μ and ν

$$\min_{\pi \in \Pi(\mu, \nu)} \left\{ \frac{1}{2} \|x - y\|^2 d\pi + \varepsilon \text{KL}(\pi \parallel \mu \otimes \nu) \right\}$$

The argmin $\pi_{\mu, \nu, \varepsilon}$ of EOT is the Schrödinger bridge. Consider same marginal setting $\mu = \nu = \rho$ with Schrödinger bridge $\pi_{\rho, \varepsilon}$. Then if $(X, Y) \sim \pi_{\rho, \varepsilon}$ and $\pi_{\rho, \varepsilon}$ admits the disintegration $\pi_{\rho, \varepsilon} = \int \pi_{\rho, \varepsilon, x} d\rho(x)$, define barycentric projection

$$\int y k^\infty(x, y) d\rho(y) = \int y d\pi_{\rho, \varepsilon, x}(y) =: \mathcal{B}_{\rho, \varepsilon}(x)$$

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

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

Heat Equation

General Wasserstein gradient flows characterized by the continuity equation

$$\partial_t \rho_t + \nabla_x \cdot (v_t \rho_t) = 0$$

where $v_t : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is the velocity field.

Heat equation $(\rho_t, t \geq 0)$ is $\partial_t \rho_t = \frac{1}{2} \Delta \rho_t$, therefore, $v_t = -\frac{1}{2} \nabla_x \log \rho_t$.

Main contribution: Under suitable conditions,

$$2x - \mathcal{B}_{\rho, \varepsilon} \approx x - \frac{\varepsilon}{2} \nabla_x \log \rho(x)$$

And in general, $\frac{1}{\varepsilon}(\mathcal{B}_{\rho, \varepsilon} - I)$ approximates the score function $\frac{1}{2} \nabla_x \log \rho$.

Convergence

Claim: Iterations of self-attention mechanism $(\rho_k^\varepsilon, k \in [T/\varepsilon])$ is a discretely approximate the heat flow starting from ρ_0 .

Define $S_\varepsilon^1(\rho) = (id - 0.5\varepsilon \nabla_x \log \rho)_\# \rho$ to be the explicit Euler step (EE)

and $SB_\varepsilon^1(\rho) := (2I_d - \mathcal{B}_{\sigma_\varepsilon})_\# \rho$ be the Schrödinger bridge (SB) step.

One-step Convergence

Proposition 1 (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

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

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

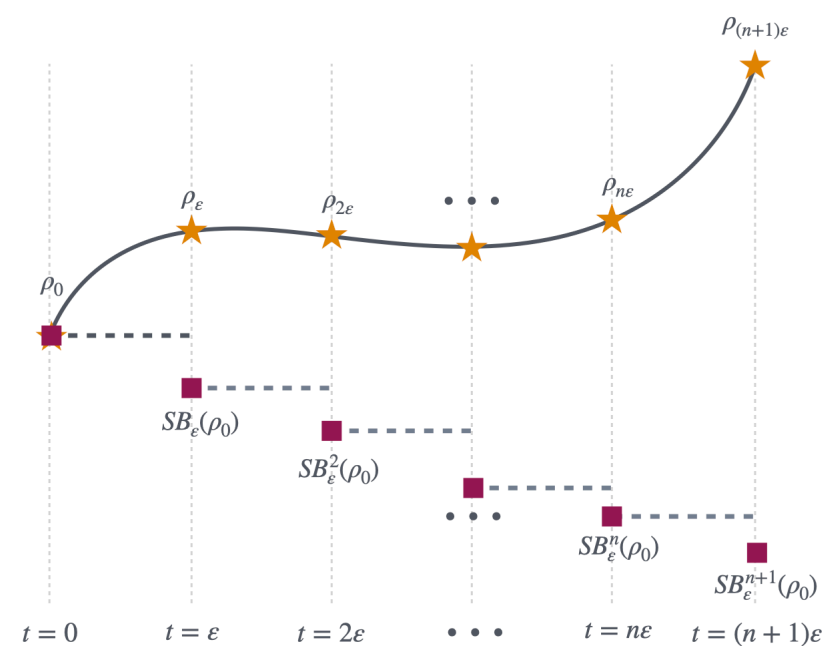
Uniform Convergence

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 accumulate.

Theorem (Uniform Convergence)

Under some regularity assumptions on the iterates $(S_\varepsilon^k, k \geq [T/\varepsilon])$ and $(SB_\varepsilon^k, k \in [T/\varepsilon])$, the sequence $(SB_\varepsilon^k, k \in [T/\varepsilon])$ is a first order approximation of heat flow $(\rho_t, t \in [0, T])$.

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



Preprint

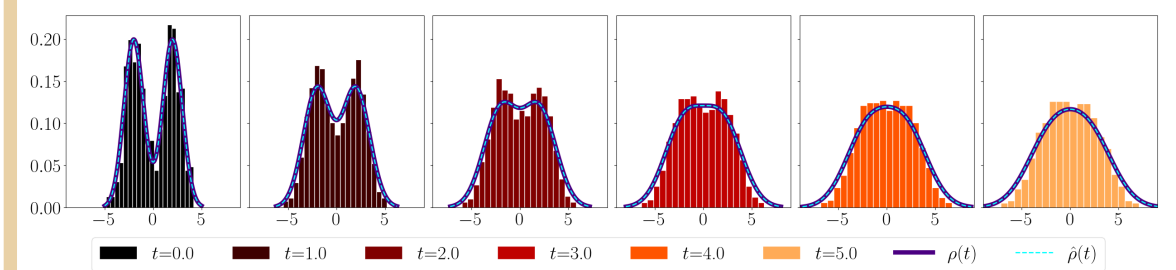


Code



Experiments

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



Heat flow with $\varepsilon = 10^2$ and $\rho_0 \sim \text{FashionMNIST}$

