

Maximum Mean Discrepancy Gradient Flow

Michael Arbel¹ Anna Korba¹ Adil Salim² Arthur Gretton¹

¹Gatsby Computational Neuroscience Unit, UCL, London

²Visual Computing Center, KAUST, Saudi Arabia

Stuttgart Worskhop on Statistical Learning 2019



Outline

1. Introduction and tools
2. Background on gradient flows
3. Maximum Mean Discrepancy Gradient Flow
4. Investigating MMD gradient flow convergence
5. A practical algorithm to descend the MMD flow
6. Applications
7. Conclusion

Outline

Introduction and tools

Background on gradient flows

Maximum Mean Discrepancy Gradient Flow

Investigating MMD gradient flow convergence

A practical algorithm to descend the MMD flow

Applications

Conclusion

Setting

- ▶ Let $\mathcal{X} \subset \mathbb{R}^d$ be the closure of a convex open set
- ▶ Let $\mathcal{P}_2(\mathcal{X})$ the set of probability measures on \mathcal{X} with finite second moment

The space $\mathcal{P}_2(\mathcal{X})$ is endowed with the Wasserstein-2 distance from
Optimal transport:

$$W_2^2(\nu, \mu) = \inf_{\pi \in \Pi(\nu, \mu)} \int \|x - y\|^2 d\pi(x, y) \quad \forall \nu, \mu \in \mathcal{P}_2(\mathcal{X})$$

where $\Pi(\nu, \mu)$ is the set of possible couplings between ν and μ .
In other words $\Pi(\nu, \mu)$ contains all possible distributions π on $\mathcal{X} \times \mathcal{X}$ such that if $(X, Y) \sim \pi$ then $X \sim \nu$ and $Y \sim \mu$.

Maximum Mean Discrepancy

- ▶ Let $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ a positive, semi-definite kernel
- ▶ \mathcal{H} its corresponding RKHS. It is a Hilbert space with inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ and norm $\|\cdot\|_{\mathcal{H}}$.

Suppose k is characteristic, ie the map:

$$\begin{aligned} \mathcal{P}_2(\mathcal{X}) &\rightarrow \mathcal{H} \\ \nu &\mapsto \int_{\mathcal{X}} k(x, \cdot) d\nu(x) \end{aligned}$$

is injective.

Maximum Mean Discrepancy

Maximum Mean Discrepancy ([Gretton et al., 2012]) defines a distance on $\mathcal{P}_2(\mathcal{X})$:

$$MMD(\mu, \nu) = \|f_{\mu, \nu}\|_{\mathcal{H}}, \text{ where}$$
$$f_{\nu, \mu}(\cdot) = \int k(x, \cdot) d\nu(x) - \int k(x, \cdot) d\mu(x)$$

$f_{\mu, \nu}$ is called the **witness function** and is the difference between the mean embeddings of ν and μ .

Now fix the (target) distribution μ . We consider the functional:

$$\mathcal{F} : \mathcal{P}_2(\mathcal{X}) \rightarrow \mathbb{R}$$
$$\nu \mapsto \frac{1}{2} MMD^2(\mu, \nu)$$

Outline

Introduction and tools

Background on gradient flows

Maximum Mean Discrepancy Gradient Flow

Investigating MMD gradient flow convergence

A practical algorithm to descend the MMD flow

Applications

Conclusion

Problem considered

Transport probability mass from a starting distribution ν to a target distribution μ , by finding a *continuous* path $(\nu_t)_{t \geq 0}$ decreasing $\mathcal{F}(\nu_t)$.

\implies **Gradient flows over the space of distributions** $\mathcal{P}_2(\mathcal{X})$

Problem considered

Transport probability mass from a starting distribution ν to a target distribution μ , by finding a *continuous* path $(\nu_t)_{t \geq 0}$ decreasing $\mathcal{F}(\nu_t)$.

⇒ **Gradient flows over the space of distributions** $\mathcal{P}_2(\mathcal{X})$

This talk: Establish conditions for convergence of MMD gradient flow to its global optimum

- ▶ novel flow over the space of distributions
- ▶ can model the optimization of some overparameterized neural networks models
- ▶ we propose a trick to improve convergence

Continuous time flows

In a **euclidean** setting, a curve $x : [0, \infty] \rightarrow \mathbb{R}^d$ is the gradient flow, or **steepest descent** of a differentiable function $F : \mathbb{R}^d \rightarrow \mathbb{R}$ if:

$$\frac{dx_t}{dt} = -\nabla F(x_t)$$

- ▶ **Initial value problem:** given x_0 , find the gradient flow $(x_t)_{t \geq 0}$.

Continuous time flows

In a **euclidean** setting, a curve $x : [0, \infty] \rightarrow \mathbb{R}^d$ is the gradient flow, or **steepest descent** of a differentiable function $F : \mathbb{R}^d \rightarrow \mathbb{R}$ if:

$$\frac{dx_t}{dt} = -\nabla F(x_t)$$

- ▶ **Initial value problem:** given x_0 , find the gradient flow $(x_t)_{t \geq 0}$.

By analogy, one can interpret the **gradient flow of a functional** $\mathcal{F} : \mathcal{P}_2(\mathcal{X}) \rightarrow \mathbb{R}$ to be a curve $\nu : [0, \infty] \rightarrow \mathcal{P}_2(\mathcal{X})$ that satisfies:

$$\frac{\partial \nu_t}{\partial t} = -\nabla_{W_2} \mathcal{F}(\nu_t)$$

for some generalized notion of gradient ∇_{W_2} , w.r.t. the W_2 metric.

Wassertein-2 gradient flows ([Ambrosio et al., 2008])

For a sufficiently regular \mathcal{F} and ν , we can write:

$$-\nabla_{W_2} \mathcal{F}(\nu) = \operatorname{div}\left(\nu \nabla \frac{\partial \mathcal{F}}{\partial \nu}\right)$$

where $\frac{\partial \mathcal{F}}{\partial \nu}$ denotes the **first variation of \mathcal{F} at ν** .

If it exists, it is the unique function such that for any $\nu, \nu' \in \mathcal{P}_2(\mathcal{X})$:

$$\lim_{\varepsilon \rightarrow 0} \frac{1}{\varepsilon} (\mathcal{F}(\nu + \varepsilon(\nu' - \nu)) - \mathcal{F}(\nu)) = \int_{\mathcal{X}} \frac{\partial \mathcal{F}}{\partial \nu}(\nu)(d\nu' - d\nu)$$

Wassertein gradient flows

Since $-\nabla_{W_2} \mathcal{F}(\mathbf{v}) = \operatorname{div}(\mathbf{v} \nabla \frac{\partial \mathcal{F}}{\partial \mathbf{v}})$, all Wassertein gradient flows are of the form:

$$\frac{\partial \mathbf{v}_t}{\partial t} + \operatorname{div}(\mathbf{v}_t V_t) = 0$$

continuity equation

Ruling the density ρ_t of particles driven by a velocity field V_t ($-\nabla \frac{\partial \mathcal{F}}{\partial \mathbf{v}}$).

Wassertein gradient flows

Since $-\nabla_{W_2} \mathcal{F}(\mathbf{v}) = \operatorname{div}(\mathbf{v} \nabla \frac{\partial \mathcal{F}}{\partial \mathbf{v}})$, all Wassertein gradient flows are of the form:

$$\frac{\partial \mathbf{v}_t}{\partial t} + \operatorname{div}(\mathbf{v}_t V_t) = 0$$

continuity equation

Ruling the density ρ_t of particles driven by a velocity field V_t ($-\nabla \frac{\partial \mathcal{F}}{\partial \mathbf{v}}$).

In particular, if the functional \mathcal{F} is a free energy:

$$\mathcal{F}(\mathbf{v}) = \underbrace{\int U(\mathbf{v}(x)) \mathbf{v}(x) dx}_{\text{internal potential } \mathcal{U}} + \underbrace{\int V(x) \mathbf{v}(x) dx}_{\text{external potential } \mathcal{V}} + \underbrace{\int W(x, y) \mathbf{v}(x) \mathbf{v}(y) dx dy}_{\text{interaction energy } \mathcal{W}}$$

$$\text{Then: } \frac{\partial \mathbf{v}_t}{\partial t} = \operatorname{div}(\mathbf{v}_t \nabla \frac{\partial \mathcal{F}}{\partial \mathbf{v}}(\mathbf{v}_t)) = \operatorname{div}(\mathbf{v}_t \nabla (U'(\mathbf{v}_t) + V + W * \mathbf{v}_t)).$$

Outline

Introduction and tools

Background on gradient flows

Maximum Mean Discrepancy Gradient Flow

Investigating MMD gradient flow convergence

A practical algorithm to descend the MMD flow

Applications

Conclusion

MMD functional

For a target distribution μ (fixed), for any $\nu \in \mathcal{P}_2(\mathcal{X})$:

$$\mathcal{F}(\nu) = \frac{1}{2} \text{MMD}^2(\mu, \nu) = \frac{1}{2} \|f_{\mu, \nu}\|_{\mathcal{H}}^2$$

- ▶ Since $\mathcal{F}(\nu) = \frac{1}{2} (\int f_{\mu, \nu} d\mu - \int f_{\mu, \nu} d\nu)$, we have $\frac{\partial \mathcal{F}}{\partial \nu} = f_{\mu, \nu}$
- ▶ Then, \mathcal{F} can be written as a free energy:

$$\mathcal{F}(\nu) = \underbrace{\int \mathbf{V}(x) d\nu(x)}_{\mathcal{V}} + \frac{1}{2} \underbrace{\int \mathbf{W}(x, y) d\nu(x) d\nu(y)}_{\mathcal{W}} + C.$$

where \mathbf{V} is a confinement potential, \mathbf{W} an interaction potential and C a constant defined by:

$$V(x) = - \int k(x, x') d\mu(x'), \quad W(x, x') = k(x, x'), \quad C = \frac{1}{2} \int k(x, x') d\mu(x) d\mu(x')$$

MMD Gradient flow

The MMD gradient flow w.r.t. W_2 is thus given by:

$$\frac{\partial \nu_t}{\partial t} = \operatorname{div}(\nu_t \nabla f_{\mu, \nu_t}) = \operatorname{div}(\nu_t \nabla (V + W * \nu_t)) \quad (1)$$

where $\nabla f_{\mu, \nu_t}(z) = \int \nabla k(x, z) d\mu(x) - \int \nabla k(x, z) d\nu_t(x)$.

This type of equation is associated in the probability theory literature to the so-called **McKean-Vlasov process** [Kac, 1956]:

$$dX_t = -\nabla f_{\mu, \underbrace{\nu_t}_{\text{depends on the current distribution of the process!}}}(X_t) dt \quad X_0 \sim \nu_0$$

whose distribution satisfy (1).

Outline

Introduction and tools

Background on gradient flows

Maximum Mean Discrepancy Gradient Flow

Investigating MMD gradient flow convergence

A practical algorithm to descend the MMD flow

Applications

Conclusion

First strategy - Convexity on the space of distributions

Definition

A curve $\rho : [0, 1] \rightarrow \mathcal{P}(\mathcal{X})$ is a geodesic between ν and μ if:

$$\rho(0) = \nu, \rho(1) = \mu, \text{ and}$$

$$L(\rho) = \min \{L(\tilde{\rho}), \tilde{\rho}(0) = \nu, \tilde{\rho}(1) = \mu\} = W_2(\nu, \mu).$$

(λ)-Geodesic convexity: Convexity of the functional \mathcal{F} on geodesic curves of $\mathcal{P}_2(\mathcal{X})$.

$$\mathcal{F}(\rho(t)) \leq (1-t)\mathcal{F}(\rho(0)) + t\mathcal{F}(\rho(1)) - t(1-t)\frac{\lambda}{2}d(\rho(0), \rho(1))^2$$

(classic λ -convexity on \mathbb{R}^d : $F((1-t)x+ty) \leq (1-t)F(x) + tF(y) - t(1-t)\frac{\lambda}{2}|x-y|^2$)

First strategy - Convexity on the space of distributions

Definition

A curve $\rho : [0, 1] \rightarrow \mathcal{P}(\mathcal{X})$ is a geodesic between ν and μ if:

$$\rho(0) = \nu, \rho(1) = \mu, \text{ and}$$

$$L(\rho) = \min \{L(\tilde{\rho}), \tilde{\rho}(0) = \nu, \tilde{\rho}(1) = \mu\} = W_2(\nu, \mu).$$

(λ) -Geodesic convexity: Convexity of the functional \mathcal{F} on geodesic curves of $\mathcal{P}_2(\mathcal{X})$.

$$\mathcal{F}(\rho(t)) \leq (1-t)\mathcal{F}(\rho(0)) + t\mathcal{F}(\rho(1)) - t(1-t)\frac{\lambda}{2}d(\rho(0), \rho(1))^2$$

(classic λ -convexity on \mathbb{R}^d : $F((1-t)x+ty) \leq (1-t)F(x) + tF(y) - t(1-t)\frac{\lambda}{2}|x-y|^2$)

Our finding: The MMD is λ -convex with $\lambda < 0$.

Too bad... $\lambda > 0$ would have guaranteed that all gradient flows of \mathcal{F} would converge the **unique** minimizer of \mathcal{F} .

Second strategy - Obtain a Lojasiewicz inequality

$$\frac{d\mathcal{F}(\mathbf{v}_t)}{dt} \leq -C\mathcal{F}(\mathbf{v}_t)^2 \quad (2)$$

Applying Gronwall's lemma results in: $\Rightarrow \mathcal{F}(\mathbf{v}_t) = \mathcal{O}(\frac{1}{t})$.

Second strategy - Obtain a Lojasiewicz inequality

$$\frac{d\mathcal{F}(\mathbf{v}_t)}{dt} \leq -C\mathcal{F}(\mathbf{v}_t)^2 \quad (2)$$

Applying Gronwall's lemma results in: $\Rightarrow \mathcal{F}(\mathbf{v}_t) = \mathcal{O}(\frac{1}{t})$.

► on the left we have the **weighted Sobolev semi-norm**:

$$\frac{d\mathcal{F}(\mathbf{v}_t)}{dt} = - \int \|\nabla f_{\mu, \mathbf{v}_t}(x)\|^2 \mathbf{v}_t(x) = -\|f_{\mu, \mathbf{v}_t}\|_{\dot{H}(\mathbf{v}_t)}^2$$

► on the right the **RKHS norm**: $\mathcal{F}(\mathbf{v}_t) = \frac{1}{2} \|f_{\mu, \mathbf{v}_t}\|_{\mathcal{H}}^2$

Second strategy - Obtain a Lojasiewicz inequality

$$\frac{d\mathcal{F}(\mathbf{v}_t)}{dt} \leq -C\mathcal{F}(\mathbf{v}_t)^2 \quad (2)$$

Applying Gronwall's lemma results in: $\Rightarrow \mathcal{F}(\mathbf{v}_t) = \mathcal{O}(\frac{1}{t})$.

► on the left we have the **weighted Sobolev semi-norm**:

$$\frac{d\mathcal{F}(\mathbf{v}_t)}{dt} = - \int \|\nabla f_{\mu, \mathbf{v}_t}(x)\|^2 \mathbf{v}_t(x) = -\|f_{\mu, \mathbf{v}_t}\|_{\dot{H}(\mathbf{v}_t)}^2$$

► on the right the **RKHS norm**: $\mathcal{F}(\mathbf{v}_t) = \frac{1}{2} \|f_{\mu, \mathbf{v}_t}\|_{\mathcal{H}}^2$

Let $L_2(\mathbf{v}) = \{f, \int f(x)^2 d\mathbf{v}(x) < \infty\}$, and $\|\cdot\|_{\dot{H}^{-1}(\mathbf{v})}$ the **weighted negative Sobolev norm**, defined for $p, q \in \mathcal{P}_2(\mathcal{X})$ by:

$$\|p - q\|_{\dot{H}^{-1}(\mathbf{v})} = \sup_{f \in L_2(\mathbf{v}), \|f\|_{\dot{H}(\mathbf{v})} \leq 1} \left| \int f(x) dp(x) - \int f(x) dq(x) \right|.$$

linearizes the W_2 !

A condition for global convergence

It can be shown that:

$$\|f_{\mu, \nu_t}\|_{\mathcal{H}}^2 \leq \|f_{\mu, \nu_t}\|_{\dot{H}(\nu_t)} \|\mu - \nu_t\|_{\dot{H}^{-1}(\nu_t)}.$$

Proof. Take $g = \|f_{\mu, \nu_t}\|_{\dot{H}(\nu_t)}^{-1} f_{\mu, \nu_t}$

- ▶ by def, $|\int g d\nu_t - \int g d\mu| = \|f_{\mu, \nu_t}\|_{\dot{H}(\nu_t)}^{-1} |\int f_{\mu, \nu_t} d\nu_t - \int f_{\mu, \nu_t} d\mu| = \|f_{\mu, \nu_t}\|_{\dot{H}(\nu_t)}^{-1} \|f_{\mu, \nu_t}\|_{\mathcal{H}}^2$
- ▶ $g \in L_2(\nu_t)^1$ and $\|g\|_{\dot{H}(\nu_t)} \leq 1$, so $|\int g d\nu_t - \int g d\mu| \leq \|\nu_t - \mu\|_{\dot{H}^{-1}(\nu_t)}$

Hence, provided $\|\mu - \nu_t\|_{\dot{H}^{-1}(\nu_t)}^2 \leq \frac{4}{C}$, we obtain the Lojasiewicz inequality (2) and the rate $\mathcal{O}(1/t)$ (thus global convergence).

However in practice it is hard to guarantee this condition ([Peyre, 2018]).

¹Under a Lipschitz assumption on ∇k , for all $\nu, \mu \in \mathcal{P}_2(\mathcal{X})$, $f_{\mu, \nu} \in L_2(\nu)$

Outline

Introduction and tools

Background on gradient flows

Maximum Mean Discrepancy Gradient Flow

Investigating MMD gradient flow convergence

A practical algorithm to descend the MMD flow

Applications

Conclusion

Euler scheme (Time-discretization of the flow)

For any $T : \mathcal{X} \rightarrow \mathcal{X}$ a measurable map, and $\nu \in \mathcal{P}_2(\mathcal{X})$, we denote the pushforward measure by $T_{\#}\nu$.

$$T_{\#}\nu(A) = \nu(T^{-1}(A)) \text{ for every measurable set } A,$$

Euler scheme (Time-discretization of the flow)

For any $T : \mathcal{X} \rightarrow \mathcal{X}$ a measurable map, and $\nu \in \mathcal{P}_2(\mathcal{X})$, we denote the pushforward measure by $T_{\#}\nu$.

$$T_{\#}\nu(A) = \nu(T^{-1}(A)) \text{ for every measurable set } A,$$

Starting from $\nu_0 \in \mathcal{P}_2(\mathcal{X})$ and using a step-size $\gamma > 0$, a sequence $\nu_n \in \mathcal{P}_2(\mathcal{X})$ is given by iteratively applying

$$\nu_{n+1} = (I - \gamma \nabla f_{\mu, \nu_n})_{\#} \nu_n. \quad (3)$$

For all n , equation (3) is the distribution of the process defined by

$$X_{n+1} = X_n - \gamma \nabla f_{\mu, \nu_n}(X_n) \quad X_0 \sim \nu_0.$$

A noisy update as regularization

The condition we exhibited for global convergence may not hold and $(\mathcal{F}(v_n))_{n \in \mathbb{N}}$ might be stuck at a local minima.

$$\frac{d\mathcal{F}(v_t)}{dt} = - \int \|\nabla f_{\mu, v_t}(x)\|^2 dv_t(x) \text{ at equilibrium} \implies \int \|\nabla f_{\mu, v^*}(x)\|^2 dv^*(x) = 0$$

If v^* positive everywhere this implies $f_{\mu, v^*} = cte = 0$ as soon as $0 \notin \mathcal{H}$. But v^* might be singular...

Our proposal: Inject noise into the gradient during updates:

$$X_{n+1} = X_n - \gamma \nabla f_{\mu, v_n}(X_n + \beta_n U_n), \quad n \geq 0,$$

where $U_n \sim \mathcal{N}(0, 1)$ and β_n is the noise level at n .

⚠ Different from adding a noise outside the gradient (i.e. diffusion)!

Guarantees

Proposition

For a choice of β_n such that:

$$8\lambda^2\beta_n^2 \mathcal{F}(\mathbf{v}_n) \leq \int \|\nabla f_{\mu, \mathbf{v}_n}(x + \beta_n u)\|^2 d\mathbf{v}_n(x) dg(u) \quad (4)$$

the following inequality holds:

$$\mathcal{F}(\mathbf{v}_{n+1}) - \mathcal{F}(\mathbf{v}_n) \leq -\frac{\gamma}{2} \left(1 - \frac{3}{2}\gamma L\right) \int \|\nabla f_{\mu, \mathbf{v}_n}(x + \beta_n u)\|^2 d\mathbf{v}_n(x) dg(u)$$

where λ and L are Lipschitz constants on the first derivatives of k .

Guarantees

Proposition

For a choice of β_n such that:

$$8\lambda^2\beta_n^2 \mathcal{F}(\mathbf{v}_n) \leq \int \|\nabla f_{\mu, \mathbf{v}_n}(x + \beta_n u)\|^2 d\mathbf{v}_n(x) dg(u) \quad (4)$$

the following inequality holds:

$$\mathcal{F}(\mathbf{v}_{n+1}) - \mathcal{F}(\mathbf{v}_n) \leq -\frac{\gamma}{2} \left(1 - \frac{3}{2}\gamma L\right) \int \|\nabla f_{\mu, \mathbf{v}_n}(x + \beta_n u)\|^2 d\mathbf{v}_n(x) dg(u)$$

where λ and L are Lipschitz constants on the first derivatives of k .

Moreover under (4)

$$\mathcal{F}(\mathbf{v}_n) \leq \mathcal{F}(\mathbf{v}_0) e^{-\Gamma \sum_{i=0}^n \beta_i^2}.$$

where $\Gamma = 4\lambda^2\gamma(1 - \frac{3}{2}\gamma L)$.

So $\sum_{i=0}^n \beta_i^2 \rightarrow \infty$ with (4) implies global convergence.

The sample-based approximate scheme

How can we simulate

$$X_{n+1} = X_n - \gamma \nabla f_{\mu, v_n}(X_n + \beta_n U_n), \quad n \geq 0?$$

It depends on:

- ▶ the current distribution $v_n \implies$ approximate it by the empirical distribution of a system of N interacting particles
- ▶ the target distribution $\mu \implies$ replace it by the empirical distribution of the M samples that we have access to ($\hat{\mu}$)

The sample-based approximate scheme

How can we simulate

$$X_{n+1} = X_n - \gamma \nabla f_{\mu, v_n}(X_n + \beta_n U_n), \quad n \geq 0?$$

It depends on:

- ▶ the current distribution $v_n \implies$ approximate it by the empirical distribution of a system of N interacting particles
- ▶ the target distribution $\mu \implies$ replace it by the empirical distribution of the M samples that we have access to ($\hat{\mu}$)

$$\hat{v}_{n+1} \begin{cases} X_{n+1}^1 = X_n^1 - \gamma \nabla f_{\hat{\mu}, \hat{v}_n}(X_n^1 + \beta_n U_n^1) \\ \dots \\ X_{n+1}^N = X_n^N - \gamma \nabla f_{\hat{\mu}, \hat{v}_n}(X_n^N + \beta_n U_n^N) \end{cases}$$

The sample-based approximate scheme

How can we simulate

$$X_{n+1} = X_n - \gamma \nabla f_{\mu, v_n}(X_n + \beta_n U_n), \quad n \geq 0?$$

It depends on:

- ▶ the current distribution $v_n \implies$ approximate it by the empirical distribution of a system of N interacting particles
- ▶ the target distribution $\mu \implies$ replace it by the empirical distribution of the M samples that we have access to ($\hat{\mu}$)

$$\hat{v}_{n+1} \begin{cases} X_{n+1}^1 = X_n^1 - \gamma \nabla f_{\hat{\mu}, \hat{v}_n}(X_n^1 + \beta_n U_n^1) \\ \dots \\ X_{n+1}^N = X_n^N - \gamma \nabla f_{\hat{\mu}, \hat{v}_n}(X_n^N + \beta_n U_n^N) \end{cases}$$

Our guarantees: For any iteration $n \in \mathbb{N}$ and $T > 0$, if $\beta_n < B$:

$$\mathbb{E}[W_2(\hat{v}_n, v_n)] \leq \frac{C_1(v_0, B, T)}{\sqrt{N}} + \frac{C_2(\mu, T)}{\sqrt{M}}$$

Outline

Introduction and tools

Background on gradient flows

Maximum Mean Discrepancy Gradient Flow

Investigating MMD gradient flow convergence

A practical algorithm to descend the MMD flow

Applications

Conclusion

Overparameterized single-layer neural network

Single-layer neural network, parameterized by $\theta \in \mathcal{X}$. Let (x, y) denote the input/output data.

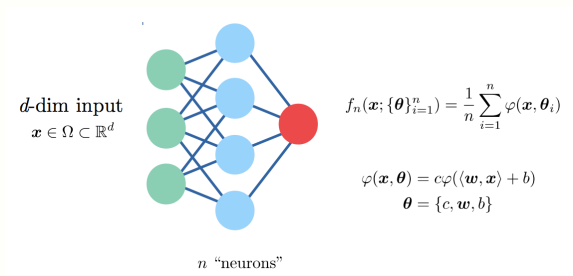


Figure: [Rotskoff et al., 2019]

Consider the supervised learning problem:

$$\min_{\underbrace{(\theta_1, \dots, \theta_n) \in \mathcal{X}}_{\text{Parameter space}}} \mathbb{E}_{(x, y) \sim p} \left[\left\| y - \frac{1}{n} \sum_{i=1}^n \phi(x, \theta_i) \right\|^2 \right]$$

Motivation

If $n \rightarrow \infty$, the previous problem can be rewritten:

$$\min_{\mathbf{v} \in \underbrace{\mathcal{P}_2(\mathcal{X})}_{\text{Distributions over the parameter space}}} \mathcal{L}(\mathbf{v}) \quad := \quad \mathbb{E}_{(x,y) \sim p} \left[\left\| y - \underbrace{\int \psi(x, \theta) d\mathbf{v}(\theta)}_{\Psi(x, \mathbf{v})} \right\|^2 \right]$$

If $\exists \mu \in \mathcal{P}_2(\mathcal{X})$ s.t. $\mathbb{E}_{y \sim p(\cdot|x)}[y] = \int \psi(x, \theta) d\mu(\theta)$,

$\implies \mathcal{L}(\mathbf{v}) = \text{MMD}^2(\mathbf{v}, \mu)$ with $k(\theta, \theta') = \mathbb{E}_{x \sim p} [\psi(x, \theta)^T \psi(x, \theta')]$.

[Chizat and Bach, 2018], [Rotskoff et al., 2019]: gradient descent on the parameters of a neural network can be seen as a particle transport problem.

Experiments - training a student-teacher network

- ▶ the teacher network $\Psi_T(x, \mu)$ is given by M particles $\mathcal{X} = (\xi_1, \dots, \xi_M)$ which are fixed during training \implies
 $\mu = \frac{1}{M} \sum_{j=1}^M \delta_{\xi_j}$
- ▶ the student network $\Psi_S(x, \nu_\Theta)$ has N particles $\Theta = (\theta_1, \dots, \theta_N)$ that are initialized randomly $\implies \nu_\Theta = \frac{1}{N} \sum_{i=1}^N \delta_{\theta_i}$

Experiments - training a student-teacher network

- ▶ the teacher network $\Psi_T(x, \mu)$ is given by M particles $\mathcal{X} = (\xi_1, \dots, \xi_M)$ which are fixed during training \implies
 $\mu = \frac{1}{M} \sum_{j=1}^M \delta_{\xi_j}$
- ▶ the student network $\Psi_S(x, \nu_\Theta)$ has N particles $\Theta = (\theta_1, \dots, \theta_N)$ that are initialized randomly $\implies \nu_\Theta = \frac{1}{N} \sum_{i=1}^N \delta_{\theta_i}$

Performing gradient descent to minimize

$$\min_{\Theta} \mathbb{E}_{x \sim p} [(\Psi_T(x, \mu) - \Psi_S(x, \nu_\Theta))^2]$$

can be seen as a particle version of the gradient flow of the MMD with a kernel given by $k(\theta, \theta') = \mathbb{E}_{x \sim p} [\psi(x, \theta') \psi(x, \theta)]$

Experiments - training a student-teacher network

- ▶ the teacher network $\Psi_T(x, \mu)$ is given by M particles $\mathcal{X} = (\xi_1, \dots, \xi_M)$ which are fixed during training \implies
 $\mu = \frac{1}{M} \sum_{j=1}^M \delta_{\xi_j}$
- ▶ the student network $\Psi_S(x, \nu_\Theta)$ has N particles $\Theta = (\theta_1, \dots, \theta_N)$ that are initialized randomly $\implies \nu_\Theta = \frac{1}{N} \sum_{i=1}^N \delta_{\theta_i}$

Performing gradient descent to minimize

$$\min_{\Theta} \mathbb{E}_{x \sim p} [(\Psi_T(x, \mu) - \Psi_S(x, \nu_\Theta))^2]$$

can be seen as a particle version of the gradient flow of the MMD with a kernel given by $k(\theta, \theta') = \mathbb{E}_{x \sim p} [\psi(x, \theta') \psi(x, \theta)]$

\implies approximated by

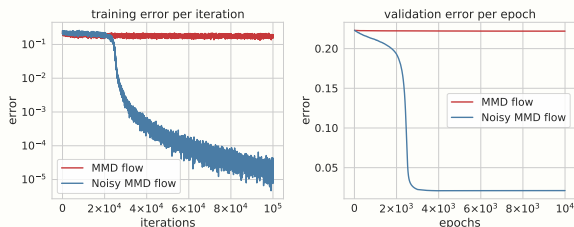
$$\hat{k}(\theta, \theta') = \frac{1}{n_b} \sum_{b=1}^{n_b} \psi(x_b, \theta)^T \psi(x_b, \theta').$$

where (x_1, \dots, x_{n_b}) are n_b samples from the data distribution.

Experiments

Leads to the approximate update:

$$\theta_{n+1}^i = \theta_n^i - \gamma \nabla \hat{f}_{\mu, v_n}(\theta_n^i)$$



⇒ adding noise to the gradient seems to lead to global convergence.

Outline

Introduction and tools

Background on gradient flows

Maximum Mean Discrepancy Gradient Flow

Investigating MMD gradient flow convergence

A practical algorithm to descend the MMD flow

Applications

Conclusion

Summary and openings

What we have done:

- ▶ novel flow over the space of distributions
- ▶ theoretical results on the MMD flow
- ▶ trick to improve convergence

Summary and openings

What we have done:

- ▶ novel flow over the space of distributions
- ▶ theoretical results on the MMD flow
- ▶ trick to improve convergence

Future work:

- ▶ Deeper understanding of the regularization proposed (continuous formulation?) and of the choice of the kernel
- ▶ Other regularizations to improve convergence?
- ▶ Other gradient flows? (here: $\nabla_{W_2} \mathcal{F}$ vs [Rotskoff et al., 2019] who get a global convergence for $\rightarrow \nabla_{WFR} \mathcal{F}$).

-  Ambrosio, L., Gigli, N., and Savaré, G. (2008).
Gradient flows: in metric spaces and in the space of probability measures.
Springer Science & Business Media.
-  Chizat, L. and Bach, F. (2018).
On the global convergence of gradient descent for over-parameterized models using optimal transport.
NIPS.
-  Gretton, A., Borgwardt, K. M., Rasch, M. J., Schölkopf, B., and Smola, A. (2012).
A kernel two-sample test.
Journal of Machine Learning Research.
-  Jordan, R., Kinderlehrer, D., and Otto, F. (1998).
The variational formulation of the fokker–planck equation.
SIAM journal on mathematical analysis, 29(1):1–17.
-  Kac, M. (1956).
Foundations of kinetic theory.

In *Proceedings of The third Berkeley symposium on mathematical statistics and probability*, volume 3, pages 171–197. University of California Press Berkeley and Los Angeles, California.



Mroueh, Y., Sercu, T., and Raj, A. (2019).
Sobolev descent.
In *AISTATS*.



Peyre, R. (2018).
Comparison between w_2 distance and \dot{H}^{-1} norm, and
localization of wasserstein distance.
ESAIM: Control, Optimisation and Calculus of Variations,
24(4):1489–1501.



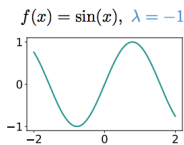
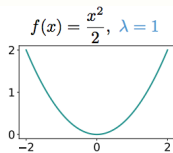
Rotskoff, G., Jelassi, S., Bruna, J., and Vanden-Eijnden, E.
(2019).
Global convergence of neuron birth-death dynamics.
In *ICML*.

Convexity on vector spaces

Existence, uniqueness results on gradient flows rely on the notion of **convexity**.

A function F defined on \mathbb{R}^d is λ -convex if $D^2F \geq \lambda I_{d \times d}$ or equivalently if for any $x, y \in \mathbb{R}^d$ and $t \in [0, 1]$:

$$F((1-t)x + ty) \leq (1-t)F(x) + tF(y) - t(1-t)\frac{\lambda}{2}|x-y|^2$$



- **Uniqueness when $\lambda > 0$:** any gradient flow $x(t)$ converges to some x^* .

W_2 distance

$$W_2^2(\nu, \mu) = \inf_{\pi \in \Pi(\nu, \mu)} \int \|x - y\|^2 d\pi(x, y) \quad \forall \nu, \mu \in \mathcal{P}_2(\mathcal{X})$$

where $\Pi(\nu, \mu)$ is the set of possible couplings between ν and μ . In other words $\Pi(\nu, \mu)$ contains all possible distributions π on $\mathcal{X} \times \mathcal{X}$ such that if $(X, Y) \sim \pi$ then $X \sim \nu$ and $Y \sim \mu$.

W_2 vs L_2 ?

L_2 geodesic: $\rho(t) = (1-t)\rho(0) + t\rho(1)$

W_2 geodesic: $\rho(t) = ((1-t)Id + tT_{\rho(0), \rho(1)})\# \rho(0)$

Informally, L^p distances are "vertical" (values of the distributions) whereas W_p distances are "horizontal" (mass of the distributions).

Gradient flows - comparison

	Euclidean	W_2
Metric (X, d)	(\mathbb{R}^d, \cdot)	$(\mathcal{P}_2(\mathbb{R}^d), W_2)$
Definition of ∇_X	$\langle \nabla F(x), v \rangle =$ $\lim_{h \rightarrow 0} \frac{F(x+hv) - F(x)}{h}$	$\langle \nabla \mathcal{F}(v), -\operatorname{div}(\xi v) \rangle_{\operatorname{Tan}_\mu \mathcal{P}_2(\mathbb{R}^d)} =$ $\lim_{h \rightarrow 0} \frac{\mathcal{F}((I+h\xi)_\# v) - \mathcal{F}(v)}{h}$
Formula for ∇_X	$\nabla_{\mathbb{R}^d} F(x) = \nabla F(x)$	$\nabla_{W_2} \mathcal{F} = -\operatorname{div}(v \nabla \frac{\partial \mathcal{F}}{\partial v})$

Comparison with Langevin

Seminal work of [Jordan et al., 1998] who revealed that the Fokker-Planck equation is a gradient flow of the Kullback-Leibler divergence:

$$\frac{\partial \nu}{\partial t} - \operatorname{div}(\nu V) = 0, \text{ where the vector field } V = \nabla_{W_2} KL(\nu) = \nabla \log\left(\frac{\nu}{\mu}\right).$$

Results in the Langevin Monte-Carlo algorithm (requires the knowledge of $\nabla \log(\mu)$):

$$X_{n+1} = X_n - \gamma \nabla \log(\mu)(X_n) + \varepsilon_n$$

where $\varepsilon_n \sim \mathcal{N}(0, 1)$.

Free energies

1. $\mathcal{F}(\mathbf{v}) = KL(\mu, \mathbf{v})$ admits a free-energy expression:

$$\mathcal{F}(\mathbf{v}) = \underbrace{\int U(\mathbf{v}(x))dx}_{\mathcal{U}} + \underbrace{\int V(x)\mathbf{v}(x)dx}_{\mathcal{V}}$$

with $U(s)$ the internal potential (entropy function) and V confinement potential defined as:

$$U(s) = s \log(s), \quad V(x) = -\log(\mu(x))$$

2. $\mathcal{F}(\mathbf{v}) = \frac{1}{2}MMD^2(\mu, \mathbf{v})$ also:

$$\mathcal{F}(\mathbf{v}) = \underbrace{\int V(x)d\mathbf{v}(x)}_{\mathcal{V}} + \underbrace{\frac{1}{2} \int W(x,y)d\mathbf{v}(x)d\mathbf{v}(y)}_{\mathcal{W}} + C.$$

where V is a confinement potential, W an interaction potential and C a constant defined by:

$$V(x) = - \int k(x,x')d\mu(x'), \quad W(x,x') = k(x,x'), \quad C = \frac{1}{2} \int k(x,x')d\mu(x)d\mu(x')$$