

# On the global convergence of gradient descent using optimal transport

L. Chizat, F. Bach, NeurIPS 2018

Matilde Dolfato

*Introduction to Real Analysis II*, Prof. G. Savaré

January 16, 2026

# Overview

Introduction

Problem re-formulation

Apply gradient flow theory

- Particle gradient flow

- Generalize to infinite-dimensional gradient flow

- Global convergence guarantees

References

# Introduction

Classical task in machine learning:

$$\min_{\mu \in \mathcal{P}(\Theta)} j(\mu) \quad j(\mu) = \ell \left( \int \phi \, \mathrm{d}\mu \right) + g(\mu) \quad (1)$$

$\ell : \mathcal{F} \rightarrow \mathbb{R}_+$  smooth, convex *loss function*

data lives in a large parametrized set with parameters  $p \in \Theta \subset \mathbb{R}^d$

$g : \mathcal{P}(\Theta) \rightarrow \mathbb{R}$  optional convex *regularizer*

we look for a linear combination  $\phi$  mapping features to labels

$\phi(p) : x \mapsto \sigma(\sum_i p_i x_i + p_d), p \in \Theta$

# Introduction

Classical task in machine learning:

$$\min_{\mu \in \mathcal{P}(\Theta)} j(\mu) \quad j(\mu) = \ell\left(\int \phi \, \mathrm{d}\mu\right) + g(\mu) \quad (1)$$

$\ell : \mathcal{F} \rightarrow \mathbb{R}_+$  smooth, convex *loss function*

data lives in a large parametrized set with parameters  $p \in \Theta \subset \mathbb{R}^d$

$g : \mathcal{P}(\Theta) \rightarrow \mathbb{R}$  optional convex *regularizer*

we look for a linear combination  $\phi$  mapping features to labels

$\phi(p) : x \mapsto \sigma(\sum_i p_i x_i + p_d), p \in \Theta$

★ Minimize over the space of measures  $\mathcal{P}(\Theta)$  ★

# Introduction

$$\min_{\mu \in \mathcal{P}(\Theta)} j(\mu) = \ell\left(\int \phi \, \mathrm{d}\mu\right) + g(\mu) \quad (1)$$

Convex problem, but intractable

# Introduction

$$\min_{\mu \in \mathcal{P}(\Theta)} j(\mu) = \ell\left(\int \phi \, \mathrm{d}\mu\right) + g(\mu) \quad (1)$$

Convex problem, but intractable

$\Rightarrow$  **discretize** the measure into  $m$  particles parametrized by *weights* and *positions*

$$\mu_m := \frac{1}{m} \sum_{i=1}^m w_i \delta_{p_i}$$

$$\min_{\substack{w \in \mathbb{R} \\ p \in \Theta^m}} j_m(w, p) \quad j_m(w, p) := j(\mu_m) \quad (2)$$

# Introduction

$$\min_{\mu \in \mathcal{P}(\Theta)} j(\mu) = \ell\left(\int \phi \, d\mu\right) + g(\mu) \quad (1)$$

Convex problem, but intractable

$\Rightarrow$  discretize the measure into  $m$  particles parametrized by *weights* and *positions*

$$\mu_m := \frac{1}{m} \sum_{i=1}^m w_i \delta_{p_i}$$

$$\min_{\substack{w \in \mathbb{R}^m \\ p \in \Theta^m}} j_m(w, p) \quad j_m(w, p) := j(\mu_m) \quad (2)$$

Non-convex problem

★ Idea: use tools from optimal transport to exploit convexity of (1)  
to study global convergence in (2) ★

# Conceptual key points

1. Problem reformulation (lifting)  $j \rightarrow f$
2. To be able to study gradient flow of  $f_m$ ,  $(\mu_{m,t})_t$
3. Look at it as a particular case of infinite-dimensional case

$$\text{as } m \rightarrow \infty, (\mu_{m,t})_t \rightarrow (\mu_t)_t$$

4. Exploit convex structure of inf-dim case

$$\lim_{m,t \rightarrow \infty} f(\mu_{m,t}) = \min_{\mu_t \in \mathcal{P}(\Omega)} f(\mu_t)$$



# Overview

Introduction

Problem re-formulation

Apply gradient flow theory

- Particle gradient flow

- Generalize to infinite-dimensional gradient flow

- Global convergence guarantees

References

# Lifting

$$\min_{\nu \in \mathcal{P}(\Theta)} j(\nu) = \ell \left( \int \phi \, d\nu \right) + g(\nu) \quad (1)$$

$$\downarrow h_1^{-1}$$

$$\min_{\mu \in \mathcal{P}(\Omega)} f(\mu) = \ell \left( \int \varphi \, d\mu \right) + \int v \, d\mu \quad (3)$$

Recovering the lifting:

- ▶  $\Omega = \Theta \times \mathbb{R}$
- ▶  $\varphi(p, w) = w\phi(p)$
- ▶ projection map:  $h_1 : \Theta \times \mathbb{R} \rightarrow \Theta$ ,  $h_1(\mu)(p) = \int_{\mathbb{R}} w\mu(dw, p)$
- ▶  $g(\nu) = \inf_{\mu \in h_1^{-1}(\nu)} \int v \, d\mu$
- ▶ then,  $\inf_{\nu} j = \inf_{\mu} f$

## New problem

$$\min_{\mu \in \mathcal{P}(\Omega)} f(\mu) = \ell\left(\int \varphi \, d\mu\right) + \int v \, d\mu$$

Discretized problem:

$$\min_{u \in \Omega} f_m(u) := f\left(\underbrace{\frac{1}{m} \sum_{i=1}^m \delta_{u_i}}_{\mu_m}\right) = \ell\left(\frac{1}{m} \sum_{i=1}^m \varphi(u_i)\right) + \frac{1}{m} \sum_{i=1}^m v(u_i)$$

★ Weights  $w$  are another coordinate of a particle position  $u \in \Omega$   
 $\Rightarrow$  study gradient flow of  $f_m$  ★

# New problem

$$\min_{\mu \in \mathcal{P}(\Omega)} f(\mu) = \ell\left(\int \varphi \, d\mu\right) + \int v \, d\mu$$

Discretized problem:

$$\min_{u \in \Omega} f_m(u) := f\left(\underbrace{\frac{1}{m} \sum_{i=1}^m \delta_{u_i}}_{\mu_m}\right) = \ell\left(\frac{1}{m} \sum_{i=1}^m \varphi(u_i)\right) + \frac{1}{m} \sum_{i=1}^m v(u_i)$$

★ Weights  $w$  are another coordinate of a particle position  $u \in \Omega$   
⇒ study gradient flow of  $f_m$  ★

Assumptions (high-level):

- ▶  $\ell$  is smooth, with Lipschitz and bounded differential  $d\ell$
- ▶  $\phi$  is differentiable and  $v$  is semiconvex
- ▶ there is a family of nested closed convex sets  $(Q_r)_r$  on which the (sub)derivatives of  $\phi, v$  are Lipschitz and grow sublinearly

# Overview

Introduction

Problem re-formulation

Apply gradient flow theory

- Particle gradient flow

- Generalize to infinite-dimensional gradient flow

- Global convergence guarantees

References

# Particle gradient flow I

## Definition (Particle gradient flow)

A gradient flow for  $f_m$  is an absolutely continuous function  $u : \mathbb{R}_+ \rightarrow \Omega^m$  such that

$$u'(t) \in -m \partial f_m(u(t))$$

for almost every  $t \in \mathbb{R}_+$ .

Properties:

- (i) (existence and uniqueness) for any  $u(0) \in \Omega^m$  starting point, there exists a unique gradient flow for  $f_m$
- (ii) (derivative of  $f_m$ ) for a.e.  $t \in \mathbb{R}_+$ , it holds

$$\left. \frac{d}{ds} f_m(u(s)) \right|_{s=t} = -|u'(t)|^2$$

## Particle gradient flow II

(iii) (form of the velocity) the velocity *of the  $i$ -th particle*  $u_i(t)$  is a vector field  $v_t : \Omega \rightarrow \mathbb{R}^d$  given by  $u_i'(t) = v_t(u_i(t))$  where [2]

$$v_t(u_i) = \tilde{v}_t(u_i) - \text{proj}_{\partial v(u_i)}(\tilde{v}_t(u_i))$$

$$\text{with } \tilde{v}_t(u_i) = - \left[ \left( \ell' \left( \int \varphi \, d\mu_{m,t} \right), \partial_j \varphi(u_i) \right) \right]_{j=1}^d$$

from [2], gradient flow selects subgradients of minimal norm  
Observations:

- ▶ velocity is the evaluation at each  $u_i$  of the same vector field  $v_t$
- ▶ given an initialization  $u(0) \in \Omega^m$ , this defines an atomic measure  $\mu_{m,0}$

→ makes sense to generalize to arbitrary measures  $\mu_t$

# Generalize to Wasserstein gradient flow I

Since

1. Evolution of  $(\mu_t)_t$  under  $(v_t)_t$  satisfies:

$$\partial_t \mu_t = -\operatorname{div}(v_t \mu_t)$$

2. Link between  $v_t$  and  $f$ :

$$v_t \in -\partial f'(\mu_{m,t})$$

$$\text{where } f'(\mu)(u) := \left( \ell \left( \int \varphi \, d\mu \right), \varphi(u) \right) + v(u)$$

$\Rightarrow$  we expect  $(\mu_t)_t$  is a gradient flow on the space  $\mathcal{P}_2(\Omega)$ :  
*Wasserstein gradient flow*



# Generalize to Wasserstein gradient flow II

## Definition (Wasserstein gradient flow)

A Wasserstein gradient flow for the functional  $f$  on a time interval  $[0, T[$  is an absolutely continuous path  $(\mu_t)_{t \in [0, T[}$  in  $\mathcal{P}_2(\Omega)$  that satisfies, distributionally on  $[0, T[ \times \Omega^d$ ,

$$\partial_t \mu_t = -\operatorname{div}(v_t \mu_t) \quad \text{where } v_t \in -\partial f'(\mu_t)$$

# Generalize to Wasserstein gradient flow II

## Definition (Wasserstein gradient flow)

A Wasserstein gradient flow for the functional  $f$  on a time interval  $[0, T[$  is an absolutely continuous path  $(\mu_t)_{t \in [0, T[}$  in  $\mathcal{P}_2(\Omega)$  that satisfies, distributionally on  $[0, T[ \times \Omega^d$ ,

$$\partial_t \mu_t = -\operatorname{div}(v_t \mu_t) \quad \text{where } v_t \in -\partial f'(\mu_t)$$

→ It is well-defined starting from  $\mu_0$  concentrated on a convex closed subset of  $\Omega$  [1]

# Generalize to Wasserstein gradient flow II

## Definition (Wasserstein gradient flow)

A Wasserstein gradient flow for the functional  $f$  on a time interval  $[0, T[$  is an absolutely continuous path  $(\mu_t)_{t \in [0, T[}$  in  $\mathcal{P}_2(\Omega)$  that satisfies, distributionally on  $[0, T[ \times \Omega^d$ ,

$$\partial_t \mu_t = -\operatorname{div}(v_t \mu_t) \quad \text{where } v_t \in -\partial f'(\mu_t)$$

→ It is well-defined starting from  $\mu_0$  concentrated on a convex closed subset of  $\Omega$  [1]

$W_2$  gradient flow generalizes particle

Whenever  $(u_t)_t$  is a gradient flow for  $f_m$ ,  $t \mapsto \mu_{m,t} := \sum_{i=1}^m \delta_{u_i(t)}$  is a Wasserstein gradient flow for  $f$ !

# Many-particle limit

## Theorem (Many particle limit)

*Consider  $(t \mapsto u_m(t))_{m \in \mathbb{N}}$  a sequence of classical gradient flows for  $f_m$  initialized in a closed convex set. If  $\mu_{m,0}$  converges to some  $\mu_0 \in \mathcal{P}_2(\Omega)$  for the Wasserstein distance  $W_2$ , then  $(\mu_{m,t})_t$  converges, as  $m \rightarrow \infty$ , to the unique Wasserstein gradient flow of  $f$  starting from  $\mu_0$ .*

In a nutshell:

- ▶ if  $u(0) = (u_1(0), \dots, u_m(0))$  are distributed according to  $\mu_0$ , then  $\mu_{m,0}$  converges to  $\mu_0$  by the law of large numbers
- ▶  $\lim_{m \rightarrow \infty} (\mu_{m,t})_t = (\mu_t)_t$

# Global convergence result

Structural assumptions:

- ▶  $\varphi, v$  are **2-homogeneous** ( $\phi$  1-homogeneous)
- ▶ the support of the initialization of the Wasserstein gradient flow satisfies a **“separation” property**:  $B(0, r_b) \subset \mathbb{S}^{d-1}$  that separates  $r_a \mathbb{S}^{d-1}$  and  $r_b \mathbb{S}^{d-1}$  for  $r_a < r_b$

## Theorem (Global convergence of particle gradient descent)

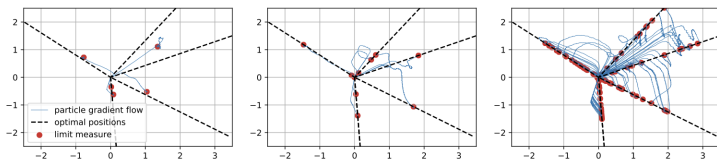
*Let  $(\mu_t)_{t \geq 0}$  be a Wasserstein gradient flow of  $f$  such that the support of  $\mu_0$  is  $S_0 \subset [-r_0, r_0] \times \Theta$  satisfies a separation property. If  $(\mu_t)_t$  converges to  $\mu_\infty$  in  $W_2$ , then  $\mu_\infty$  is a global minimizer of  $f$  over  $\mathcal{P}(\Omega)$ . In particular, if  $(u_m(t))_{m \in \mathbb{N}, t \geq 0}$  is a sequence of classical gradient flows initialized in  $[-r_0, r_0] \times \Theta$  such that  $\mu_{m,0}$  converges to  $\mu_0$  in  $W_2$  then*

$$\lim_{t, m \rightarrow \infty} f(\mu_{m,t}) = \min_{\mu \in \mathcal{P}(\Omega)} f(\mu).$$

# Application to ReLu neural networks

Setting:

- ▶ features live in  $\mathbb{R}^{d-2}$ , labels in  $\mathbb{R}$
- ▶  $\ell$  is either the square or logistic loss
- ▶ 2-homogeneous case
- ▶ domain  $\Theta$  is the disjoint union of 2 copies of  $\mathbb{R}^d$
- ▶  $\varphi(p) : x \mapsto \sigma(\sum_{i=1}^{d-1} s(p_i)x_i + s(p_d))$ ,  $s(p_i) = p_i|p_i|$
- ▶ regularizer:  $v(p) = |p|^2$  (as if  $w = |p|$ )



**Figure:** Training neural network with ReLU activation. Overfitting threshold  $m = 4$ . Failure for  $m = 5$ , success  $m = 10, 100$ .

# Final remarks

- ▶ importance of initialization, as confirmed by extensive empirical literature
- ▶ particle gradient flow corresponds to continuous-time gradient descent: what can we say about discrete-step case? (double descent)

Thank you!





# References I



Luigi Ambrosio, Nicola Gigli, and Giuseppe Savaré.

*Gradient flows: in metric spaces and in the space of probability measures.*

Springer, 2005.



Filippo Santambrogio.

{Euclidean, metric, and Wasserstein} gradient flows: an overview.

*Bulletin of Mathematical Sciences*, 7(1):87–154, 2017.

## Recovering the lifting of (3)

1. (equivalence of  $\mathcal{M}$  and  $\mathcal{P}$  under homogeneity)
2. surjectivity of  $h_1$
3. define  $\varphi, g$  as above and prove equality of  $f, j$

# Prerequisites

## Definition (Subgradient and subdifferential)

Given a function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , the *subgradient* of  $f$  at  $x_0$ ,  $x_0 \in \mathbb{R}^d$ , is  $p \in \mathbb{R}^d$  :

$$f(x) \geq f(x_0) + p \cdot (x - x_0) + o(x - x_0) \quad \forall x \in \mathbb{R}^d$$

The set of all such  $ps$  is called the *subdifferential* of  $f$  at  $x_0$  and we write  $\partial f(x_0)$ . The subdifferential is a closed and convex set.

## Definition (Gradient flow)

A function  $x : \mathbb{R}_+ \rightarrow \text{Dom}(f)$  is a *gradient flow* of  $f : \text{Dom}(f) = \mathbb{R}^d \rightarrow \mathbb{R}$  if  $x$  is absolutely continuous over  $\mathbb{R}_+$  and

$$x'(t) \in -\partial f(x(t)) \quad \text{for a.e. } t \in \mathbb{R}_+$$