

## Lecture 7: Mean-field neural networks

Summary of the class so far: (Mid-semester "sanity check")

2 central questions: { Why optimization is tractable?  
Why do NNs generalize?

### Two hypotheses

\* Tractability via overparametrization: the highly non-convex landscape simplify dramatically as # parameters  $\gg$  # samples so that local search methods (e.g., GD/SGD) succeed at finding global minima ("interpolating" solution)

\* Implicit regularization of training algorithm: algo itself bias the selection of a solution that generalizes well on new data DESPITE the model being overparametrized and interpolating noisy data.

Lectures 2/3/4: Overparametrization is not incompatible with generalization

Lecture 2: capacity of NNs (Rademacher complexity) depends on norms / spectral properties of the weights and not # parameters.

Lecture 3/4: solution selected by GD (min-norm interpolating solut°) in linear models can overfit benignly thanks to a self-induced regularization

This motivates the study of gradient-trained NNs that are highly-overparametrized, and as a first approximation  $\# \text{neurons} \rightarrow \infty$

$$\text{2-layer NN: } f(x; \Theta) = \frac{1}{\sqrt{M}} \sum_{j=1}^M a_j^\circ \sigma(\langle w_j^\circ, x \rangle) \quad \|x\|_2 = 1$$

$$a_j^\circ \stackrel{\text{iid}}{\sim} N(0, 1) \quad w_j^\circ \stackrel{\text{iid}}{\sim} N(0, I_d)$$

"Standard Xavier initialization" ("used in pytorch")

(This is more subtle than this  $\rightarrow$  see next week)

Lecture 5: As  $M \rightarrow \infty$ , GD on this NN effectively behave as a linear model

$$f(x; \theta^t) \approx f(x, \theta^0) + \langle \theta^t - \theta^0, \nabla_{\theta} f(x, \theta^0) \rangle \quad (3)$$

CV exponentially fast to global solution (morebly named efficiently)

↳ "Lazy regime"  $\|\theta^t - \theta^0\|_2 \ll 1$

(general mechanism that has nothing to do with NNs a priori)

Lecture 6: Linearized NNs are kernel methods which have limited adaptivity (adapt to "smoothness" but not to other low-dimensional structure)

→ "fixed feature" method       $\nabla_{\theta} f(x, \theta^0)$   
 "fixed representation"

→ suffer from the curse of dimensionality

If  $\theta^t$  moves away from  $\theta^0$  during dynamics, can hope to learn good representation  $\nabla_{\theta} f(x, \theta^t)$  and do much better

"feature learning", "representation learning", "learning a good kernel"

Example: Learning a single neuron

Consider  $y_i = \sigma(\langle w_*, x_i \rangle) + \varepsilon$

$$x_i \sim \text{Unif}(S^{d-1})$$

$$\|w_*\|_2 = \sqrt{d}$$

Fit with  $f(x; \theta) = \frac{1}{\sqrt{M}} \sum_{j=1}^M \alpha_j \sigma(\langle w_j, x \rangle)$

$$(1) \text{ on } \hat{R}_m(\theta) = \frac{1}{m} \sum_{i=1}^m (y_i - f(x_i; \theta))^2$$

Then [Montanari, Zhong, '22, Ghorbani, Mei, Mniakiewicz, Montanari '21]

Fix  $S > 0$  and take  $M d \geq m^{1+S}$ . Let  $\hat{\theta}$  be the solution found in the Lazy training regime ( $\alpha \geq \sqrt{\frac{m^2}{Md}}$ )

If  $m \leq d^{1+S}$ , then

$$\liminf_{m, d \rightarrow \infty} R(f(\cdot, \hat{\theta})) \geq \|P_{\geq d} \sigma\|_{L^2}^2$$

fit at most a  $d^{\circ}$  polynomial  
approximat<sup>o</sup> to  $\sigma$

This is disappointing given that  $\sigma(\langle w_*, x \rangle)$  is

(5)

extremely simple and we can fit with one neuron !!!

The problem is  $w_j^*$  don't move during optimization in the lazy regime and we are trying to fit  $\sigma(\langle w_*, \alpha \rangle)$  with  $\sigma(\langle w_j^*, \alpha \rangle)$

$$\hookrightarrow \text{in high dimensions } \sup_{j \in [M]} \frac{|\langle w_j^*, w_* \rangle|}{\|w_j^*\|_2 \|w_*\|_2} \leq \frac{1}{\sqrt{d}}$$

If instead  $M = 1$   $f(x; \theta) = \sigma(\langle x, w \rangle)$

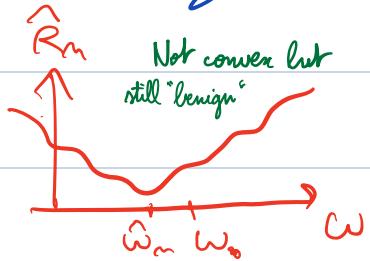
$$\hat{R}_m(w) = \frac{1}{m} \sum_{i=1}^m (y_i - \sigma(\langle x_i, w \rangle))^2$$

Thm: [Bai, Mei, Montanari, '18] Assume  $\sigma$  bounded and 3 times differentiable bounded derivatives, and  $\sigma'(\cdot)$  for all  $t \in \mathbb{R}$ . Then there exists a constant  $C > 0$  such that if  $m \geq C d \log d$ , we have

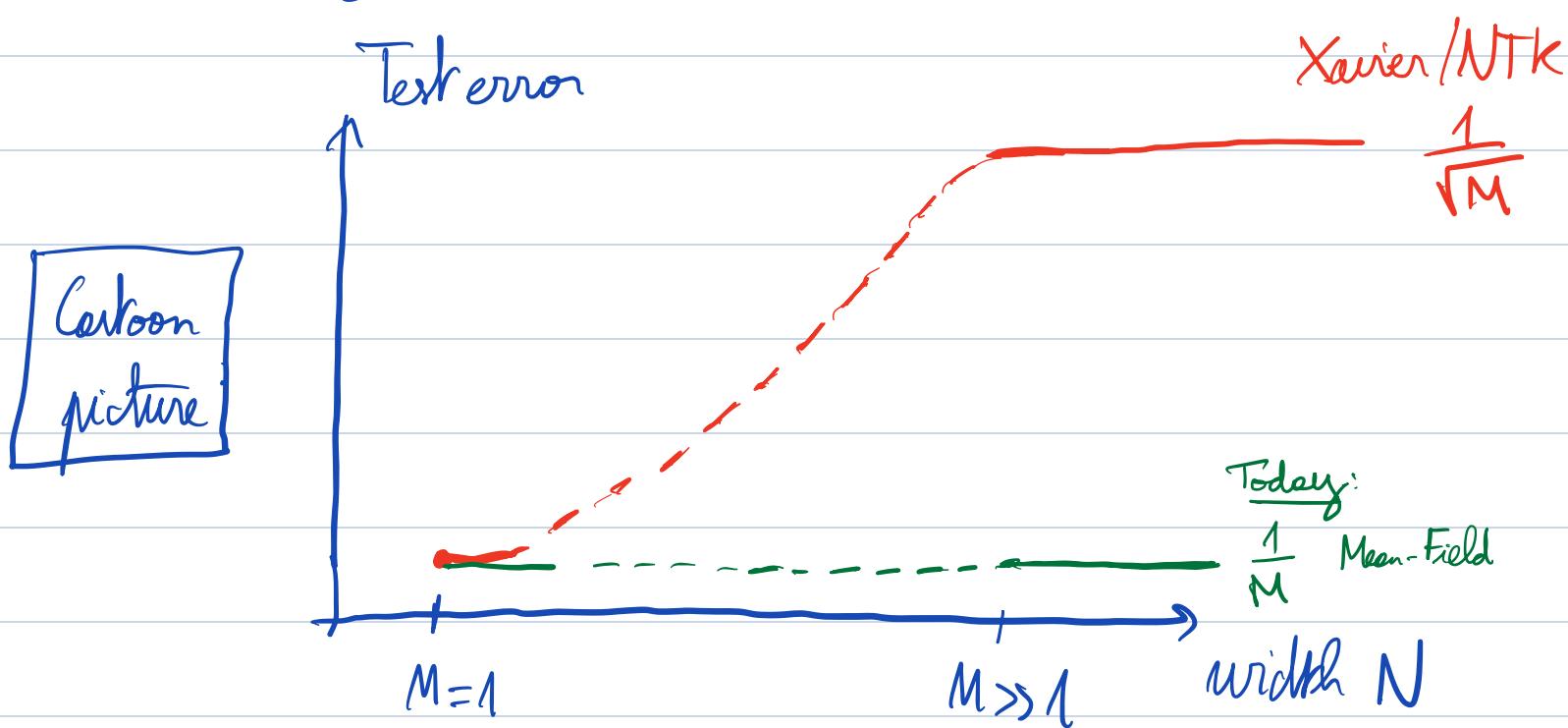
(i)  $\hat{R}_m(w)$  has unique local/global minimizer  $\hat{w}_m$

(ii) GD CV to  $\hat{w}_m$

(iii)  $R(\hat{w}_m) \leq C \sqrt{\frac{d \log m}{m}}$



Summary: learning a single neuron



Test error deteriorates as  $M \rightarrow \infty$ . What is going wrong?

→ This is only a specific infinite-width limit

↳ can chose different scalings and obtain other limits with different behavior.

Today: Mean-Field limit; scaling  $\frac{1}{M}$

## Mean-Field scaling

$$f(x; \Theta) = \frac{1}{M} \sum_{j=1}^M a_j \sigma(\langle w_j, x \rangle)$$

$\xrightarrow{\quad}$

$$\frac{\alpha_M}{\sqrt{M}} \quad \alpha_M := \frac{1}{\sqrt{M}}$$

Why is this a natural parametrization?

① Lecture 1: approximation theory of NNs

$$\frac{1}{M} \sum_{j=1}^M a_j \sigma(\langle w_j, x \rangle) \xrightarrow[M \rightarrow \infty]{} \int a(w) \sigma(\langle w, x \rangle) \mu(dw)$$

$f_*$  well approximated with  $a(w)\mu(dw)$  → then sample

② Lecture 2: Rademacher complexity  $\lesssim n_0 \sqrt{\frac{d}{m}} \quad \frac{\|a\|_1}{M} \leq n_0$   
with MF scaling.

$$\left[ \text{With NTK scaling} \lesssim n_0 \sqrt{\frac{Md}{m}} \quad \frac{\|a\|_1}{M} \leq n_0 \right]$$

③ Lecture 5:  $\frac{\alpha_M}{\sqrt{M}}$  lazy regime  $\alpha_M \gtrsim \sqrt{\frac{m^2}{Md}}$

To not be in the lazy regime as  $M \rightarrow \infty$  for  $m/d$  fixed

[we need  $\alpha_M \lesssim \frac{1}{\sqrt{M}}$  as  $M \rightarrow \infty$ ]

Rmk: Systematic study of infinite width limits of NNs through "ABC" parametrization [Greg Yang, 2021]

→ classify scalings that give stable and non-trivial limits as  $M \rightarrow \infty$

→ MF is special among these limits (for 2-layer) as it maximizes "feature learning"

(Might cover it next week... very useful (easy) computation)

(3)

## Stochastic Gradient Descent (SGD)

Let's rewrite

$$f(x; \Theta) = \frac{1}{M} \sum_{j=1}^M \sigma(x; \Theta_j) \quad \Theta_j \in \mathbb{R}^d$$

$$\Theta = (\Theta_1, \dots, \Theta_M) \in \mathbb{R}^{MD}$$

$$\text{e.g.: } \sigma(x; \Theta_j) = a_j \tilde{\sigma}(\langle x, w_j \rangle) \quad \Theta_j = (a_j, w_j) \in \mathbb{R}^{d+1}$$

$$\text{Consider squared loss: } l(y, \hat{y}) = \frac{1}{2} (y - \hat{y})^2$$

$$\text{Test error: } R_D(\Theta) = \frac{1}{2} E_{(y, x) \sim D} [(y - f(x; \Theta))^2]$$

$$\text{Train error: } \hat{R}_m(\Theta) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (y_i - f(x_i; \Theta))^2$$

$$= R_{\hat{D}_m}(\Theta) \quad \hat{D}_m = \frac{1}{m} \sum_{i=1}^m S_{y_i, x_i}$$

Simplest algorithm: Gradient Descent ( $\approx 1850$ )

$$\Theta^{k+1} = \Theta^k - \eta \nabla \hat{R}_m(\Theta^k)$$

$\approx 1950$ : we do not need to compute exact gradient

$$\mathbb{w}^{k+1} = \mathbb{w}^k - \eta \left( \nabla \hat{R}_m(\mathbb{w}^k) + z^k \right)$$

↑  
iid noise with zero mean

→ noise averages out and algo converges as GD (with suitable step sizes)

SGD (with batch size = 1)

\* Initializing:  $\mathbb{w}^0, \theta_j^0 \sim_{\text{iid}} \mathcal{C}_0$

\* At each step: sample  $(y_k, x_k)$

$$\mathbb{w}^{k+1} = \mathbb{w}^k - \underbrace{\eta M}_{\text{scaling that gives } \mathbb{w}^k(1)} \nabla_{\mathbb{w}} l(y_k, f(x_k; \mathbb{w}^k))$$

"multi-pass SGD"

Two cases (1) Finite training samples  $(y_m, x_m)$  drawn uniformly at random from  $\hat{D}_m$

as  $\eta \downarrow 0$   $\mathbb{w}^t$  follows GF on  $\hat{R}_m(\mathbb{w})$

$$\frac{d}{dt} \mathbb{w}^t = -M \nabla_{\mathbb{w}} \hat{R}_m(\mathbb{w}^t)$$

(2) Infinite training samples  $(y_k, x_k) \stackrel{iid}{\sim} D$

"one-pass SGD"  $\rightarrow$  each sample only used once

$$\text{as } \eta \downarrow 0 \quad \frac{d}{dt} \Theta^t = -M \nabla_{\Theta} R(\Theta^t) \quad \text{"population" GF}$$

Both cases can be treated in a unified way  $\hat{R}_m = R_{\hat{D}_m}$  (population dist is  $D_m$ )

$\hookrightarrow$  below "one pass SGD"

Dynamics:  $\Theta_j^t \stackrel{iid}{\sim} \rho_0$

$$\cdot \Theta_j^{k+1} = \Theta_j^k + \eta \nabla_{\Theta_j} \sigma(x_k; \Theta) \left( y_k - \frac{1}{M} \sum_{s=1}^M \sigma(x_s; \Theta_s) \right)$$

Goal: we want to produce a simplified description of the SGD dynamics in this setting

$$* \underline{\eta \downarrow 0}: \frac{d}{dt} \Theta^t = - M \nabla_{\Theta} R(\Theta^t)$$

$$* \underline{M \rightarrow \infty}: f(x; \Theta^t) = f(x; \hat{\rho}_t^{(M)})$$

$$f(x; \rho) = \int \sigma(x; \theta) \rho(\theta)$$

$$\hat{\rho}_t^{(M)} = \frac{1}{M} \sum_{j=1}^M S_{\Theta_j^t}$$

# Empirical weight distribution

Note that  $\Theta_j^0 \stackrel{iid}{\sim} \rho_0$  at  $t=0$   $\hat{\rho}_0^{(M)} \xrightarrow{P} \rho_0$

It is reasonable to expect

$\hat{\rho}_t^{(M)} \Rightarrow \rho_t$  for all  $t \geq 0$ . [i.e.,  $\forall$  bounded  $C^0$  fct  
 $\int f d\rho_t^0 \rightarrow \int f d\rho_t$ ]

where  $\rho_t$  is a deterministic measure following some PDE which we will describe next.

Remark: This is classically referred to as "propagation of chaos": as  $M \rightarrow \infty$ , a system of interacting particles can be well approximated by independent particles interacting with a "mean-field"  $\rho_t$ .

## Mean-field limit: static properties

$$\begin{aligned}
 \text{Risk: } R_M(\Theta) &= \frac{1}{2} \mathbb{E}_{y|x} \left[ (y - f(x; \Theta))^2 \right] \\
 &= \frac{1}{2} \left[ R_{\#} + \frac{2}{M} \sum_{j=1}^M V(\Theta_j) + \frac{1}{M^2} \sum_{i,j=1}^M U(\Theta_i, \Theta_j) \right]
 \end{aligned}$$

where

$$R_{\#} = \mathbb{E}[y^2] \quad V(\Theta) = -\mathbb{E}_{y|x} [y \sigma_x(x; \Theta)]$$

$$U(\Theta_1, \Theta_2) = \mathbb{E}_x [\sigma_x(x; \Theta_1) \sigma_x(x; \Theta_2)]$$

Interpretation:  $R_M(\Theta)$  is the energy of a system of  $M$  particles  $\Theta_j \in \mathbb{R}^D$ , interacting with an external potential  $V(\Theta)$  and via a pairwise potential  $U(\Theta_1, \Theta_2)$

In particular:  $U$  is P.S.D. i.e. for any bounded and compactly supported function  $h$ , we have

$$\int h(\Theta_1) U(\Theta_1, \Theta_2) h(\Theta_2) d\Theta_1 d\Theta_2 \geq 0$$

"repulsive interaction".

Gradient w.r.t one particle:

$$\begin{aligned} M \nabla_{\theta_j} R_M(\theta) &= \nabla_{\theta_j} \left( V(\theta_j) + \frac{1}{M} \sum_{i=1}^M U(\theta_j, \theta_i) \right) \\ &=: \nabla_{\theta_j} \Psi(\theta_j; \theta) \end{aligned}$$

As  $M$  is large, it is natural to replace  $\theta_1, \dots, \theta_M$  by a density  $\rho \in \mathcal{P}(R^D)$

$$R(\rho) = \frac{1}{2} R_{\#} + \int V(\theta) \rho(d\theta) + \int U(\theta_1, \theta_2) \rho(d\theta_1) \rho(d\theta_2)$$

$$\text{And } \Psi(\theta; \rho) = V(\theta) + \int U(\theta, \theta') \rho(d\theta')$$

$$" = \frac{\delta R(\rho)}{\delta \rho(\theta)} = \begin{array}{l} \text{"additional energy of adding} \\ \text{a single particle at } \theta \in R^D \end{array}$$

$$\underline{\text{Rmk: }} R(\hat{\rho}^{(M)}) = R_M(\theta)$$

$$\frac{d}{dt} \theta_j^t = -\nabla_{\theta} \Psi(\theta_j; \hat{\rho}_t^{(M)}) \quad \begin{array}{l} \text{"this is a completely equivalent} \\ \text{descript. of GF on M memory"} \end{array}$$

Hence, heuristically,  $\hat{\rho}_t^{(M)} \Rightarrow \rho_t$

Each particle

$$\theta^t \sim \rho_t \text{ satisfy } \frac{d}{dt} \theta^t = -\nabla_{\theta} \Psi(\theta^t; \rho_t)$$

"consistency equation"

## PDE characterization

We can describe evolution of  $e_t$  in terms of a PDE

Below is a heuristic derivation: assume  $e_t(d\theta) = \overset{\text{has a density}}{\underset{\text{for simplicity}}{\overset{\wedge}{e_t}(\theta)}} d\theta$

Consider a differentiable test function  $g: \mathbb{R}^D \rightarrow \mathbb{R}$

$$\begin{aligned} \frac{d}{dt} \mathbb{E}[g(\theta^t)] &= \frac{d}{dt} \int g(\theta) e_t(\theta) d\theta \\ &= \int g(\theta) \partial_t e_t(\theta) d\theta \end{aligned} \quad \downarrow \text{differentiating under integral}$$

On the other hand, exchanging derivative / expectation

$$\frac{d}{dt} \mathbb{E}[g(\theta^t)] = \mathbb{E}\left[\frac{d}{dt} g(\theta^t)\right] = \mathbb{E}\left[\langle \nabla_{\theta} g(\theta^t), \frac{d}{dt} \theta^t \rangle\right]$$

$$= -\mathbb{E}\left[\langle \nabla_{\theta} g(\theta^t), \nabla_{\theta} \psi(\theta; e_t) \rangle\right]$$

$$\begin{aligned} \text{Integral by part} \quad &= - \int \langle \nabla_{\theta} g(\theta), \nabla_{\theta} \psi(\theta; e_t) \rangle e_t(\theta) d\theta \\ &= \int g(\theta) \nabla \cdot [e_t(\theta) \nabla_{\theta} \psi(\theta; e_t)] d\theta \end{aligned}$$

$$\text{divergence } \nabla_{\theta} \cdot v(\theta) = \sum_{i=1}^D \frac{\partial}{\partial \theta_i} v_i(\theta) \quad (16)$$

Hence putting equalities together:

$$\int g(\theta) \left\{ \partial_t c_t(\theta) - \nabla \cdot [c_t(\theta) \nabla_{\theta} \Psi(\theta; c_t)] \right\} d\theta = 0$$

"Weak form" of the following PDE

$$\boxed{\partial_t c_t = \nabla \cdot [c_t \nabla_{\theta} \Psi(\theta; c_t)]}$$

This is often referred to as "Distributional Dynamics"

This is known as a McKean-Vlasov type PDE.

PDE independently proven in 4 concurrent work

[CB'18, MMN'18, RVE'18, SS'18]

If initializing a particle  $\theta_j^0 \sim c_0$  and

$$\frac{d}{dt} \theta_j^t = -\nabla_{\theta} \Psi(\theta_j^t; c_t)$$

then  $\theta_j^t \sim c_t$  for all  $t \geq 0$

## Approximation guarantees

How well  $f(x; \rho_t)$  tracks  $f(x; \Theta^k)$  ?

where •  $\Theta^k$  solut° of SGD with M neurons  
η step size

- $\rho_t$  solut° of above PDE.

Assumptions : (i)  $|g|, |y| \leq C$ ,  $\nabla_\theta g(x; \theta)$  C-subGaussian  
 (ii)  $V(\theta)$ ,  $V(\theta_1, \theta_2)$  have bounded Lipschitz gradient  
 (weaker assume in the paper)

Thm [Mei, Misiakiewicz, Montanari, '19]

There exists a constant K that only depends on assumptions  
 s.t. for all  $T \geq 0$ , with probability  $\geq 1 - e^{-\gamma}$

$$\sup_{k \in [0, \frac{T}{\eta}] \cap \mathbb{N}} |R_M(\Theta^k) - R(\rho_{k\eta})| \leq \frac{Ke^{KT}}{\sqrt{M}} (\sqrt{\log M} + \gamma)$$

error due to  
 $M \rightarrow \infty$  approx

$$+ Ke^{KT} (\sqrt{D + \log M} + \gamma) \sqrt{\eta}$$

error due to  $\eta \rightarrow 0$  approx.

Rank: •  $\hat{C}_k^{(M)} \approx C_k \eta$   $t = k\eta$

- For  $T$  fix,  $M \gtrsim 1$  and  $\eta \lesssim \frac{1}{D}$  for MF to be a good approximation.
  - In particular  $M$  does not need to depend on ① but only on properties of the target fct
  - similar to approximat<sup>o</sup> theory (e.g Barron's result)
- $e^{kT}$  is bad but in some sense unavoidable
  - ↳  $T = O(1)$        $\frac{T}{\eta} = O(D)$  SGD steps

↳ good approx limited to this linear scaling  
 ↳ will talk more about it in following lectures

## Wasserstein gradient flow description

The PDE has an interesting structure: it is a gradient flow in the space of distributions.

It is worth fleshing out this relation as it connects MF limit to a rich field of mathematics: optimal transport.

### What is a gradient flow?

Consider:

- \*  $(M, \rho)$  a metric space ( $\rho = \text{distance}$ )
- \*  $F: M \rightarrow \mathbb{R}$

Trajectory:  $z(0), z(\varepsilon), z(2\varepsilon), \dots, z(k\varepsilon), \dots$

$$z(k\varepsilon + \varepsilon) = \underset{z \in M}{\operatorname{argmin}} \left\{ F(z) + \frac{1}{2\varepsilon} \rho(z, z(k\varepsilon)) \right\}$$

As  $\varepsilon \downarrow 0$ , this defines a continuous trajectory  $z(t)$

"GF on  $F$  in  $(M, \rho)$ "

We minimize "movements" in  $(M, \rho)$  to descend on  $F(z)$

E.g. Euclidean space  $\frac{dz(t)}{dt} = -\nabla F(z(t))$  is indeed the limit of the above

Previous PDE is a Gradient flow on  $R(\rho)$  in the metric space  $(\mathcal{P}_2(\mathbb{R}^D), W_2)$

\*  $\mathcal{P}_2(\mathbb{R}^D)$ : space of probe  $\rho$  on  $\mathbb{R}^D$  with  $\int \| \theta \|^2 \rho(d\theta) < \infty$

\*  $W_2$ : Wasserstein distance

$$W_2(\mu, \nu) = \inf_{\gamma \in \Gamma(\mu, \nu)} \left( \int \|x - y\|_2^2 \gamma(dx, dy) \right)^{\frac{1}{2}}$$

where  $\Gamma(\mu, \nu)$  is the set of all couplings of  $\mu$  and  $\nu$   
i.e. joint probe on  $\mathbb{R}^{D \times D}$  with marginals  $\mu$  and  $\nu$

$$\rho_{t+\varepsilon} \approx \underset{\rho \in \mathcal{P}(\mathbb{R}^D)}{\operatorname{argmin}} \left\{ R(\rho) + \frac{1}{2\varepsilon} W_2(\rho, \rho_t)^2 \right\}$$

In particular  $\frac{d R(\rho_t)}{dt} = - \int \|\nabla \varphi(\theta; \rho_t)\|_2^2 \rho_t(d\theta).$

$\leq 0$

## Global convergence

$R(\rho_t)$  is monotonically decreasing. Can we show global convergence? (And have some success as lazy regime)

$$\lim_{t \rightarrow \infty} R(\rho_t) = \inf_{\rho} R(\rho)$$

[Here we consider one-pass SGD hence global CV is directly on test error, i.e., directly implies that we generalize well]

→ not quite: showing global CV of this Wasserstein flow is hard and there are many subtle difficulties.

Below I describe some basic properties.

### Static properties:

Minimizing  $R_M(\cdot)$  is not much different than minimizing  $R(\rho)$

Lemma: Assume  $\exists \varepsilon, k > 0$  such that

$$R(\rho) \leq \inf_{\theta} R(\rho) + \varepsilon \Rightarrow \int U(\theta, \theta) \rho(d\theta) \leq k$$

Then:

$$\left| \inf_{\theta \in \mathbb{R}^M} R_M(\theta) - \inf_{\theta} R(\rho) \right| \leq \frac{k}{M}$$

Proof: Take  $\rho_*$  s.t.  $R(\rho_*) \leq \inf_{\theta} R(\rho) + \delta$   $\delta \leq \varepsilon$

(in particular  $\int U(\theta, \theta) \rho_*(d\theta) \leq k$  by assumption)

Consider  $(\theta_j)_{j \in [M]} \sim_{\text{iid}} \rho_*$ , then simple calculation yields

$$\begin{aligned} \mathbb{E}_{\rho_*} [R_M(\theta)] - R(\rho_*) &= \frac{1}{M} \left\{ \int U(\theta, \theta) \rho_*(d\theta) - \underbrace{\int U(\theta_1, \theta_2) \rho_*(d\theta_1) \rho_*(d\theta_2)}_{\geq 0 \text{ by PSD.}} \right\} \\ &\leq \frac{1}{M} \int U(\theta, \theta) \rho_*(d\theta) \leq \frac{k}{M} \end{aligned}$$

Hence  $\inf_{\theta} R_M(\theta) \leq \inf_{\theta} R(\rho) + \frac{k}{M} + \delta$  for any  $\delta \square$

We further have the following characterization of minimizer

Prop:

Assume  $U, V$  are  $C^1$  and bounded from below.

Assume  $\exists \varepsilon, \gamma, K > 0$ , s.t.

$$R(\rho) \leq \inf_{\rho} R(\rho) + \varepsilon \Rightarrow \int \|U\|^\gamma \rho(d\theta) \leq K$$

Then (1)  $\inf_{\rho} R(\rho)$  is achieved at some  $\rho_*$ .

(2)  $\rho_*$  is a minimizer iff

$$\text{supp}(\rho_*) \subseteq \arg\min_{\theta} U(\theta; \rho_*)$$

(Proof is by doing a perturbation of  $\rho_*$ )

Corollary: Suppose  $\rho_*$  satisfy (i)  $\text{supp}(\rho_*) = \mathbb{R}^D$   
and (ii)  $U(\theta; \rho_*)$  is a constant

Then  $\rho_*$  is a minimizer.

Proof: [This is direct from previous proposition (2)]

## Dynamical properties:

$\rho$  is a fixed point of the PDE if and only if

$$(*) \quad \text{supp } (\rho) \subseteq \{ \theta : \nabla_\theta U(\theta; \rho) = 0 \}$$

(simply comes from  $\frac{d}{dt} R(\rho_t) = -\int \|\nabla_\theta U\|_2^2 \rho_t(d\theta)$ )

There will be  $\infty$  number of "bad" stationary points  
 (e.g. any stationary pt of landscape with finite # of neurons will be a stationary of the PDE with empirical measure  $\hat{\rho}^{(n)}$ )

It is hard & subtle to argue that  $\rho_t$  will not get trapped in these stationary points.

Here is a simple result:

Thm: Assume (a)  $\rho_t \rightarrow \rho_\infty \in \mathcal{P}_2(\mathbb{R}^D)$   
 (b)  $\text{supp } (\rho_\infty) = \mathbb{R}^D$

(+  $U, V$  differentiable with bounded gradient)

Then  $\rho_\infty$  is a minimizer, that is

$$R(\rho_t) \rightarrow \min_e R(e)$$

Proof:  $\rho_\infty$  stationary point, i.e  
 $\mathbb{R}^D = \text{supp}(\rho_\infty) \subseteq \{\theta: \nabla \Psi(\theta; \rho_\infty) = 0\}$

hence  $\Psi(\theta; \rho_\infty)$  cte and we conclude using  
the corollary above □

Assumption  $\text{supp}(\rho_\infty) = \mathbb{R}^D$  is very restrictive

[Chizat, Bach, '18] [Wojtowysch, '22]

[Nguyen, Pham, '20]

Always need to assume

$\rho_t \rightarrow \rho_\infty$  (\*)  
 $\rho_\infty$  has some abstract properties

[(\*) this is not a benign condition as  $\rho_t$  might oscillate while  
 $R(\rho_t) \downarrow R_*$ ]

Nonetheless: These proofs are instructive and show that it is hard to trap a distribution in a non-optimal stationary point  
(Kickability via overparametrized)

Summary: • No quantitative general CV guarantees  
• Even qualitative guarantees are under abstract assumptions

# Noisy SGD

26

Consider the following variant of SGD

→ add  $\frac{\lambda}{2} \|\Theta\|_2^2$  regularization to the risk

→ add at each step Gaussian noise to the gradient

$$\text{At each step: } \theta_j^{k+1} = (1 - \lambda) \theta_j^k + \eta \nabla_{\theta} \sigma(x; \theta_j^k) (y_k - f(x_k; \Theta^k))$$

$$+ \sqrt{\frac{\eta}{\beta}} g_j^k \xrightarrow{iid} g_j^k \sim N(0, I_D)$$

$$\text{Risk: } R_\lambda(\rho) = R(\rho) + \frac{\lambda}{2} \int \|\theta\|_2^2 \rho(d\theta)$$

$$\Psi_\lambda(\theta; \rho) = \Psi(\theta; \rho) + \frac{\lambda}{2} \|\theta\|_2^2$$

Then as  $M \rightarrow \infty$   $\eta \downarrow 0$

The resulting dynamics

Additional diffusion term

$$\partial_t \rho_t = \nabla_{\theta} \cdot [\rho_t \nabla_{\theta} \Psi_\lambda(\theta; \rho_t)] + \frac{1}{\beta} \Delta \rho_t$$

"inverse temperature"

Laplacian

$$\Delta \cdot f(\theta) = \sum_{i=1}^D \frac{\partial^2}{\partial \theta_i^2} f(\theta)$$

(27)

This is also a Wasserstein gradient flow but now on the free energy

$$F_{\lambda, \beta}(e) = R_\lambda(e) - \frac{1}{\beta} S(e)$$

with  $S(e) = - \int e(\theta) \log e(\theta) d\theta$  the entropy of  $e$

The interesting property when adding entropic regularizer is that the Free energy is strongly convex and there is only one stationary point which is the global minimizer of  $F_{\lambda, \beta}(e)$

Lemma: For any  $\lambda \geq 0$  and  $\beta > 0$ ,  $F_{\lambda, \beta}(e)$  has at most one fixed point  $F_{\lambda, \beta}(e_*) < \infty$ .

Furthermore  $e_*$  has density satisfying self-consistent Boltzmann equation

$$e_*(\theta) = \frac{1}{Z(\lambda, \beta)} \exp(-\beta \Psi_\lambda(\theta; e_*))$$

Proof: This follows from showing that the dissipative equation

$$\frac{d F_{\alpha, \beta}(\rho_t)}{dt} = - \int_{\mathbb{R}^D} \left\| \nabla_\theta \left[ \Psi_\alpha(\theta; \rho_t) + \frac{1}{\beta} \log \rho_t(\theta) \right] \right\|_2^2 \rho_t(\theta) d\theta$$

If  $\dot{\rho} = 0$  iff  $\rho^*$  satisfy equation above

In particular this means that

$$\lim_{t \rightarrow \infty} F_{\alpha, \beta}(\rho_t) = \min_{\rho} F_{\alpha, \beta}(\rho)$$

If  $\beta$  taken sufficiently small  $\min_{\rho} F_{\alpha, \beta}(\rho) \approx \min_{\rho} R(\rho)$

[We need  $\beta = \mathcal{R}(\mathcal{D})$ ]

We get:  $\lim_{t \rightarrow \infty} R(\rho_t) \leq \inf_{\rho} R(\rho) + \Theta\left(\frac{\mathcal{D}}{\beta}\right)$

Great!!! Global CV as long as  $\beta < \infty$   
but what about convergence time?

Standard approach: prove Log-Sobolev inequality  
for this setting

→ If LSI with LS constant  $\alpha$  then

$$F_{\lambda, \beta}(\rho_t) - F_{\lambda, \beta}(\rho_*) \leq e^{-C\frac{\alpha}{\beta}t} (F_{\lambda, \beta}(\rho_0) - F_{\lambda, \beta}(\rho_*))$$

However: best we can show  $\alpha \geq \lambda \beta e^{-CB}$

Hence if  $\beta = \Theta(D)$   $\alpha \asymp e^{-\Theta(D)}$

and require  $T = e^{\Theta(D)}$  to converge

→ This cannot be improved in general (i.e., there exist settings where we will need  $e^D$  time to CV).

# Returning to single neuron learning example

(30)

We return to the example of learning a single neuron which we started with

This illustrate how the PDE can be simplified in settings with symmetries.

$$\text{assume } \sigma'(t) > 0 \quad \forall t \in \mathbb{R}$$

Consider:  $\alpha_i \sim N(0, I_d)$   $y_i = \sigma(\langle w_*, \alpha_i \rangle) + \varepsilon_i$

Model:  $f(x; \theta) = \int a \sigma(\langle w, x \rangle) p(dw)$

Training: "uninformative initializer"

$$p_0(dw) = P_A \otimes N(0, \frac{\lambda^2}{d} I_d)$$

$a$  initialized  $\perp$   $w$  initialized

$$\partial_t e_t = \nabla \cdot [e_t \nabla_{\theta} \mathcal{W}(\theta; e_t)]$$

(31)

Exploring symmetries: data distribution  $(y, \omega)$  remain invariant under a rotation that leaves  $\omega_*$  unchanged

Hence: if  $R \in \mathbb{R}^{d \times d}$  s.t  $R\omega_* = \omega_*$

then  $R_{\#} \rho_t$  is a solution of the PDE with initialization  $R_{\#} \rho_0$

Hence  $R_{\#} \rho_0 = \rho_0 \rightarrow$  by uniqueness of the solution

we must have  $R_{\#} \rho_t = \rho_t$  for all  $t \geq 0$ .

Hence  $\rho_t$  is invariant under rotations that leave  $R\omega_* = \omega_*$   
i.e. only depend on  $(\alpha, \langle \omega_*, \omega \rangle, \|P_{\perp} \omega\|_2)$

$$\begin{array}{c} \parallel \\ \alpha \\ \parallel \\ \pi \end{array}$$

$$\rho_t(d\alpha d\omega) = \bar{\rho}_t(d\alpha, ds, dr) \otimes \text{Unif}(\mathbb{S}^{d-2})$$

Denote  $z = (\alpha, s, r)$

$$\Psi(z; \bar{\rho}_t) = \Psi(\theta; \rho_t) = V(\theta) + \int U(\theta, \theta') \rho_t(d\theta')$$

$$V(\theta) = -\mathbb{E}[y \alpha \sigma(\langle \omega, \alpha \rangle)] \quad U(\theta, \theta') = \mathbb{E}[\alpha \alpha' \sigma(\langle \omega, \alpha \rangle) \sigma(\langle \omega, \alpha' \rangle)]$$

$$V(\theta) = -\alpha \mathbb{E}[\sigma(\langle \omega_s, \alpha \rangle) \sigma(\langle \omega, \alpha \rangle)]$$

$$= -\alpha \mathbb{E}_{G_1, G_2} \left[ \sigma(G_1) \sigma(sG_1 + nG_2) \right] \\ \sim N(0, 1)$$

$$U(\theta, \theta') = \alpha \alpha' \mathbb{E}_{\substack{G_1, G_2, G_3 \\ \sim N(0, 1)}} \left[ \sigma(sG_1 + nG_2) \sigma(s'G_1 + \frac{\langle \omega_s, \omega'_s \rangle}{n} G_2 + \sqrt{n'^2 - \frac{\langle \omega_s, \omega'_s \rangle^2}{n^2}} G_3) \right]$$

When taking integral over  $\omega'_s = n' u$   $u \sim \text{Unif}(S^{d-2})$   
 this indeed only depend on  $(\alpha, s, n)$ .

Hence the dynamics only depends on 3 parameters instead of  $d+1$   
 Reduced form  $\bar{\rho}_t \in P(\mathbb{R}^3)$

$$\partial_t \bar{\rho}_t = \partial_a (\bar{\rho}_t \partial_a \Psi(g, \bar{\rho}_t)) + \partial_s (\bar{\rho}_t \partial_s \Psi(g, \bar{\rho}_t)) + \frac{1}{n} \partial_n (n \bar{\rho}_t \partial_n \Psi(g, \bar{\rho}_t))$$

Initialization translates into  $\bar{\rho}_0 = \underbrace{P_A}_{a^0} \otimes \underbrace{N(0, \frac{\gamma^2}{d})}_{s^0} \otimes \underbrace{Q_{d-1, \gamma}}_{n^0}$

$$\text{where } Q_{d-1, \gamma} = \frac{\gamma}{\sqrt{d+1}} \sqrt{\chi^2_{(d)}}$$

The above PDE can be solved numerically efficiently

$$\lim_{t \rightarrow \infty} R(\rho_t) = 0$$

excess test error

For  $\varepsilon > 0$ ,  $\exists T = T(\varepsilon)$  indep of  $d$  s.t.  $R(\rho_T) \leq \varepsilon$

Then taking  $M \gtrsim \frac{e^{KT(\varepsilon)}}{\varepsilon^2} = \mathbb{W}_d(1)$

$$\eta \lesssim \frac{e^{-KT(\varepsilon)}}{\varepsilon^2 d} = \mathbb{W}_d\left(\frac{1}{d}\right)$$

We get  $R_M(\mathbb{W}^{\frac{T}{\eta}}) \leq R(\rho_T) + \varepsilon \leq 2\varepsilon$

Total # samples =  $\mathbb{W}_d(d)$  # neurons =  $\mathbb{W}_d(1)$

(Success?)

## Conclusion

Mean-Field scaling:  $\rho_t \neq \rho_0$  "feature learning"

From single neuron example: can indeed adapt to low-dim structure and vastly outperform NTK.

Goal of  $M \rightarrow \infty$  is to simplify analysis of SGD.

Did we succeed?

Pros:

- Compact description of dynamics in terms of PDE
- When setting has symmetries, can exploit them to reduce the effective dimension of the parameters

Cons:

- These PDE are hard to analyze in general
- Hard to prove qualitative/quantitative CV guarantees
- only good approx for  $O(D)$  steps of SGD

Ongoing line of work!!!