

Scidolyse
Groupe de travail
Mathematics of Deep Learning

On the convergence of gradient descent
for 2-layer neural networks:
A mean-field view

Src: Andrea Montanari (course, '18, '19)
Chizat and Bach '19 (Neurips)

3 challenges:

- Approximation

- Convergence of Gradient Descent

- Generalization

(no over-fitting)

Supervised Learning

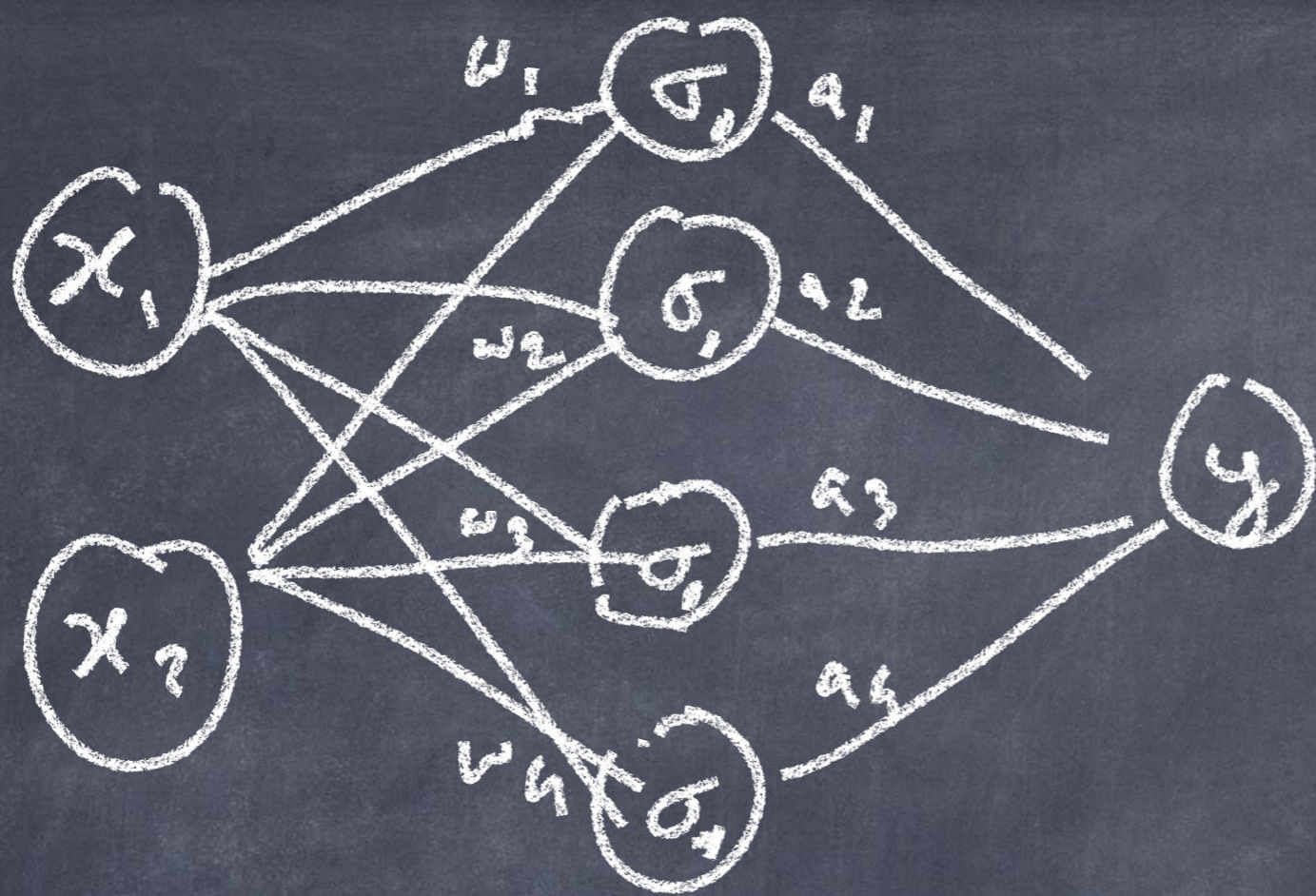
$$\cdot (x, y) \sim \mathcal{P}$$

sample $(x_1, y_1) \dots (x_n, y_n)$

$$\cdot \mathcal{R}(\theta) = \mathbb{E} \left[(f(x; \theta) - y)^2 \right]$$

$$\hat{\mathcal{R}}(\theta) = \mathbb{E}_n \left[(f(x; \theta) - y)^2 \right]$$

$$= \frac{1}{n} \sum_{i=1}^n (f(x_i; \theta) - y_i)^2$$



Neural Nets

1 hidden layer

$$y = f(x; \theta) = \frac{1}{N} \sum_{i=1}^N \sigma(x; \theta_i)$$

$$\theta_i = (a_i, b_i, w_i)$$

$$\sigma(x, \theta) = a \sigma_0(\langle w, x \rangle + b)$$

$$\sigma_0(x) = x + \quad \text{or} \quad \sigma_0(x) = \frac{1}{1 + e^{-2x}} \dots$$

Th [Cybenko '89]

$$\int \overline{\mathbb{E}}(|f(x)|^2) < \infty$$

$$\text{if } \sigma: \mathbb{R} \rightarrow \mathbb{R} \text{ is continuous, } \begin{cases} \sigma(x) \rightarrow 1 \\ \sigma(x) \rightarrow 0 \end{cases} \begin{matrix} x \rightarrow +\infty \\ x \rightarrow -\infty \end{matrix}$$

then $\forall \epsilon > 0, \exists N(\epsilon) \text{ s.t.}$

$$\inf_{(a_i, b_i, w_i)} \overline{\mathbb{E}} \left[f(x) - \frac{1}{N} \sum_{i=1}^N a_i \sigma(w_i x + b_i) \right] \leq \epsilon$$

PL: how to find
the good approximator?

-> Gradient Descent: stepsize

$$\theta^{k+1} = \theta^k + \epsilon^k \nu_k$$

$$\nu_k = -\nabla R(\theta^k)$$

SGD:

$$\nu_k = -\nabla R(\theta^k) + \epsilon^k \underbrace{\nu_k}_{\text{noise}}$$

PF: $R(\theta)$ is not convex in θ
it has local minima, etc...

BUT

still, (S)GD works
especially when the network is
over-parameterized \rightarrow WHY?

Here: $R(\theta) = \mathbb{E} \left[\left(\frac{1}{N} \sum_{i=1}^N \sigma(x_i; \theta) - y \right)^2 \right]$

$$\rightarrow \theta_i^{k+1} = \theta_i^k + 2 \eta \mathbb{E} \left[\sigma(x, \theta_i) \left(y - \frac{1}{N} \sum_{j=1}^N \sigma(x_j, \theta_j) \right) \right]$$

$$R_N(\theta) = \mathbb{E}[y^2] + \frac{2}{N} \sum_{i=1}^N \underbrace{\mathbb{E}[y \sigma(x, \theta_i)]}_{\text{external potential}} + \frac{1}{N^2} \sum_{i,j=1}^N \mathbb{E}[\sigma(x, \theta_i) \sigma(x, \theta_j)]$$

$$R_N(\theta) = R_{\#} + \frac{2}{N} \sum_{i=1}^N V(\theta_i) + \frac{1}{N^2} \sum_{i,j=1}^N U(\theta_i, \theta_j)$$

↑ energy
 ↑ external potential
 ↑ pairwise potential

Important: the kernel U is

semi-definite:

A h bounded, compactly supported

$$\iint U(\theta_1, \theta_2) h(\theta_1) h(\theta_2) d\theta_1 d\theta_2 \geq 0$$

repulsive interaction (in average sense)

Convexification

$$f(x, \theta) = \frac{1}{N} \sum_{i=1}^N \sigma(x, \theta_i)$$

$$= \int \sigma(x, \theta) p(d\theta)$$

when N is large, can be modelled by a density p

$$\rightarrow R_N(\theta)$$

$$\approx R(e) = R_{\#} + 2 \int v(\sigma) e(d\sigma)$$

$$+ \int v(\sigma_1, \sigma_2) p(d\sigma_1) e(d\sigma_2)$$

Prop: $\exists \epsilon_0: \forall e, R(e) \leq \inf_e R(e) + \epsilon_0$

$$\Rightarrow \int v(\sigma, \sigma) p(d\sigma) \leq \kappa$$

then $\lim_{N \rightarrow \infty} |R_N(\theta) - \inf_e R(e)| \leq \frac{\kappa}{N}$

Remark:

$$R(\rho) = R_{\#} + 2 \int v(\theta) p(d\theta) + \int v(\theta_1, \theta_2) (p(\theta_1), p(\theta_2))$$

is convex in ρ !

conv optimization in ∞ -dim $\mathcal{M}_+(\Theta)$

2-dim approach

functional derivative

$$\Psi(\theta, e) = \frac{1}{2} \frac{\delta R(e)}{\delta \rho(\theta)} \approx V(\theta) + \int U(\theta, \theta') \rho(d\theta')$$

= variation of energy when adding 1 particle at θ

ρ_* is a minimum if

$$\text{supp}(\rho_*) \subset \underset{\theta \in \mathbb{R}^D}{\text{argmin}} \Psi(\theta, \rho_*)$$

→ Idea 1:

- Discretize $\Theta \rightarrow \Theta_N$
- Minimize $R(p_N)$ on $\mathcal{M}_1(\Theta_N) \subseteq \mathbb{R}^N$

$$R(p_N) = (R_{\neq} +) + p_N V_N + p_N U_N p_N$$

where $V_N(i) = v(\sigma_i)$ and $U_N(i,j) = v(\sigma_i, \sigma_j)$

Pf: Curse of dimensionality:

requires gigantic N .

→ what follows proves that a good N
does not need to be exponential in D .

Idea 2: Particular approach

$$\rho_N = \frac{1}{N} \sum_{i=1}^N \delta_{\mathcal{O}_i}$$

each particle \mathcal{O}_i moves according to the force of the system: at time $t = h\varepsilon$

$$\text{speed } \mathbb{E}(\dot{\sigma}_i^k | \mathcal{F}_k) = -\nabla_{\sigma_i} V(\sigma_i^k) - \frac{1}{N} \sum_{j=1}^N \nabla_{\sigma_i} V(\sigma_i^k, \sigma_j^k)$$

$$= \mathbb{E}[Y_{\sigma}(x, \sigma_i)] - \frac{1}{N} \sum_{j=1}^N \nabla_{\sigma_i} \mathbb{E}[\sigma(x, \sigma_i) | \sigma(x, \sigma_j)]$$

→ this is exactly (S)FD!

→ prove that the particle system behaves like its continuous equivalent (statistical mechanics)

$$\mathbb{E}[\sigma_i^k | \mathcal{F}_k] = -\nabla \psi(\sigma_i^k, \ell_{\text{loc}})$$

$(P_t)_{t \geq 0}$ continuous time limit.

Continuity equation

$$\frac{\partial \rho}{\partial t} = - \overset{\text{div}}{\nabla} \cdot (\rho(\theta) v(\theta, t))$$



enters: $\rho(\theta) v(\theta, t)$

leaves: $\rho(\theta + d\theta) v(\theta + d\theta, t)$

→ variation of ρ : $\frac{\partial}{\partial \theta} (\rho(\theta) v(\theta, t))$

$$\rightarrow \partial_t \rho_t(\theta) = \nabla_{\theta} \cdot (\rho_t(\theta) \nabla_{\theta} \psi(\theta, \rho_t))$$

Fixed points = densities ρ_* s.t.
all mass sits on zero velocity positions:

$$\text{supp}(\rho_*) \subset \partial \mathbb{R}^D : \nabla \psi(\theta, \rho_*) = 0$$

Thm: if $\|\sigma_0\|_\infty \leq \kappa_2$, $\|\nabla_\theta \sigma_*(x_0)\|_2 \leq \kappa_2$

$$|y_{x_0}| \leq \kappa_2$$

if $\|\nabla_\theta v(0)\|_2 \leq \kappa_3$, $\|\nabla_{\theta_1, \theta_2} v(\theta_1, \theta_2)\|_2 \leq \kappa_3$

$$\|\nabla_\theta v(\theta) - \nabla_\theta v(\theta')\|_2 \leq \kappa_3 \|\theta - \theta'\|_2$$

$$\|\nabla_{\theta_1, \theta_2} v(\theta_1, \theta_2) - \nabla_{\theta_1, \theta_2} v(\theta'_1, \theta'_2)\|_2 \leq \kappa_3 \left(\|\theta_1 - \theta'_1\|_2 + \|\theta_2 - \theta'_2\|_2 \right)$$

$\rho_0 \in \mathcal{S}(\mathbb{R}^D)$. SGD with initialization $(\theta^i)_{i=1}^N = \rho_0$
 step size $\gamma_k = \frac{\epsilon}{2}$

ρ_ϵ : solution of the PDE

Then $\exists C = C(\kappa_i)$, $\forall f: \mathbb{R}^D \times \mathbb{R} \rightarrow \mathbb{R}$ s.t. $\|f\|_\infty \leq 1$, $\|f\|_{\text{Lip}} \leq 1$

$$\sup_{t=0}^{\frac{1}{\epsilon}} \left| R_N(\theta^k) - R(\rho_{k\epsilon}) \right| \leq C e^{-ct} \sqrt{\frac{1}{N}} \sqrt{\frac{1}{\epsilon}} \left[\sqrt{D} + \frac{1}{\epsilon} + 3 \right]$$

with proba $\geq 1 - e^{-3t}$

→ IDE accurate as long as $N \gg D$
 $\epsilon \ll \frac{1}{N}$

→ no curse of dimensionality.

numerical experiments → PDE approx
very accurate in practice.

→ global convergence can be proved
in some cases.

Gradient Flows

$$\dot{x}(t) = -\nabla F(x(t))$$

"continuous time gradient descent"

$$x(t+\varepsilon) = \underset{z \in \mathbb{R}^d}{\operatorname{arg\,min}} \left\{ F(z) + \frac{1}{2\varepsilon} \|z - x(t)\|_2^2 \right\}$$

$$\text{of GD: } x_{k+1} = \underset{z \in \mathbb{R}^d}{\operatorname{arg\,min}} \left\{ F(z) + \frac{1}{2\alpha_k} \|z - x_k\|_2^2 \right\}$$

→ general definition: for a distance
 $d(\cdot, \cdot)$

$$x_\varepsilon((k+1)\varepsilon) = \underset{z \in \mathbb{R}^d}{\operatorname{arg\,min}} \left\{ F(z) + \frac{1}{2\varepsilon} d(z, x_\varepsilon(k)) \right\}$$

$$x(k) = \lim_{\varepsilon \rightarrow 0} x_\varepsilon(k)$$

is the gradient flow of the cost
function F on X for the metric d

Prop:

$$\partial_t \rho(t) = \nabla_{\partial} \cdot (\rho_{+}(t) \nabla_{\partial} \psi(t, \rho))$$

is the gradient flow for the cost $R(\rho)$ in Wasserstein metric

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

Application of μ and ν = probability on $\mathbb{R}^d \times \mathbb{R}^d$
with marginals μ and ν

Noisy GD

$$\theta_i^{k+1} = \theta_i^k + \eta_k \nabla_{\theta_i} \sigma(x_k, \theta_i^k) \left(y_k - \frac{1}{N} \sum_{i=1}^N \sigma(x_k, \theta_i^k) \right) + \sqrt{\eta_k} g_i^k \quad g_i^k \sim \mathcal{N}(0, I_D)$$

$$\rightarrow \partial_t p_t(\theta) = \nabla_{\theta} \cdot (p_t(\theta) \nabla_{\theta} \Psi(\theta, p_t)) + T \Delta p_t(\theta)$$

$$F(p) = \frac{1}{2} R(p) - TS(p) \quad \text{free energy}$$

$$\text{where } S(p) = - \int p(\theta) \log p(\theta) d\theta \quad \text{entropy}$$

$$p_t(\theta) = \frac{1}{Z(\beta)} \exp(-\beta \Psi(\theta, p_t)) \quad \text{Boltzmann equation}$$

One can prove convergence in a time
that depends on D but not on N .

→ SGD reaches a near-optimum
in time independent of the number of neurons