# A Mean Field View of the Landscape of Two-Layer Neural Network (1)
## Song Mei, Andrea Montanari, Phan-Minh Nguyen

## I - Learning with a shallow NN

$$y_i = f(x_i) + \varepsilon_i$$

depth 2 width N neural net $\quad \hat{f}(x, \vec{\theta}) := \frac{1}{N} \sum_{i=1}^{N} \sigma_*(x, \theta_i) \qquad \vec{\theta} \in (\mathbb{R}^d)^N = \omega^N \qquad \omega \subset \mathbb{R}^D$

typ. $d=3$ $\quad \theta_i = (a_i, b_i, \omega_i) \quad \sigma_*(x, \theta_i) = a_i \sigma(\langle \omega_i, x \rangle + b_i) \qquad \sigma = $ rel. nonlip.

loss $R_N(\vec{\theta}, x, y) = \mathbb{E}\left[ (y - \hat{f}(x\vec{\theta}))^2 \right] \qquad$ avg risk: $R_N(\vec{\theta}) = \mathbb{E}[R_N(\theta, X, Y)]$

minimized by SGD: $(x_k, y_k)$ iid $\quad$ Robbins Monroe

$$\vec{\theta}^{k+1} = \vec{\theta}^k - s_k \nabla_\theta R_N(\vec{\theta}, X_k, Y_k) \qquad \mathbb{E}[\nabla_\theta R_N(\vec{\theta}, X_k, Y_k)] = \nabla R_n(\vec{\theta})$$

here: $\vec{\theta}_i^{n+1} = \vec{\theta}_i^n + 2 s_k \nabla_{\theta_i} \sigma_*(X_k, \theta_i^k)\left( Y_k - \frac{1}{N} \sum_{j=1}^{N} \sigma_*(X_k, \theta_j^k) \right)$

$\underline{Pb}$: does it converge to $\min R_N$ ? $\to$ predict dynamics

ex: $Y_i \sim U(f_1)$ $X_i \sim N(0, (1+4)I_i)$ $Z_i \sim N(0, 1)$ Montanari

or $Y = \hat{f}(x, \vec{\theta}_0)$ and try to fit $\hat{f}$ $\quad$ (digest/Bach)

## II Mean Field View

2 "idealizations": $\quad$ (gradient) $\to$ expectation

$\qquad\qquad\qquad$ mean-field : $\infty$ of neurons (particles)

$$R_N(\vec{\theta}, x, y) = y^2 + \frac{2}{N} \sum_{i=1}^{N} -y \sigma_*(x, \theta_i) + \frac{1}{N^2} \sum_{i,j} \sigma_*(x, \theta_i) \sigma_*(x, \theta_j)$$

$\rho_N = \frac{1}{N} \sum \delta_{\theta_i}$

$\rho_N^2 = \frac{1}{N^2} \sum_{i,j} \delta_{\theta_i, \theta_j}$

$$= y^2 + 2 \int \underbrace{-y \sigma_*(x, \theta) d\rho_N(\theta)}_{V(\theta, x, y)} + \iint \underbrace{\sigma_*(x, \theta) \sigma_*(x, \theta') d\rho_N^2(\theta)}_{U(\theta, \theta', x, y)}$$

$V(\vec{\theta}) = \mathbb{E}[V(\vec{\theta}, x, y)] = $ external potential $\qquad U(\theta, \theta') = \mathbb{E}[U(\theta, \theta', X, Y)] = $ pairwise potential

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \geq 0$, s. definite $\to$ repulsive

$R_N(\vec{\theta}) = R_\# + 2 \int V(\theta) d\rho_N(\theta) + \iint U(\theta, \theta') d\rho_N^2(\theta) \qquad \int V(\theta, \theta') 2(\theta) d\theta d\theta'$

$\mathbb{E}(y^2)$ $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \downarrow$ idealize $\qquad \mathbb{E}[\sigma_*(x, \theta) \rho(\theta)]^2$

$$R(\rho) = R_\# + 2 \int V(\theta) d\rho(\theta) + \iint U(\theta, \theta') d\rho(\theta) d\rho(\theta') \qquad \to R_N(\theta) = R(\rho_N)$$

idea: $\quad \left| \inf_{\vec{\theta} \in \omega^N} R_N(\vec{\theta}) - \inf_\rho R(\rho) \right| \leq \frac{K}{N}$

$\qquad\qquad$ if $\exists K > 0, \varepsilon_0 > 0$ st $\forall \rho = R(\rho) \leq \inf_\rho R(\rho) + \varepsilon_0, \int U(\theta, \theta') d\theta d\theta' \leq K$

Proof: let $\rho^*$ st $R(\rho^*) \leq R(\rho) + \varepsilon$ and let $\theta_1, ..., \theta_N \sim \vec{\theta}^*$

then $\mathbb{E}[R_N(\vec{\theta}^*)] = \int \mathbb{E}[R_N(\vec{\theta}, x, y)] d\rho^*(\theta) = R_\# + \int V(\theta) d\rho^*(\theta) + \int \frac{1}{N^2} \sum_{i,j} U(\theta_i, \theta_j) d\rho^*(\theta) d\rho^*(\theta)$

$= R_\# + \int V(\theta) d\rho^*(\theta) + \int U(\theta, \theta') d\rho^*(\theta) d\rho^*(\theta') + \frac{1}{N}\left( \int U(\theta, \theta) d\rho^*(\theta) - \iint U(\theta, \theta') d\rho^*(\theta) d\rho^*(\theta') \right)$

$\leq R(\rho) + \frac{K}{N}$

## III. Mean Field Dynamics

What does SGD correspond to in the continuous world?

answer: Gradient flow for $R(\rho)$ in Wasserstein metric

SGD: $\quad \theta_i^{h+1} = \theta_i^h + s_h \underbrace{-\nabla_{\theta_i} R_N(\vec{\theta}^h, x_h, y_h)}_{} \qquad \rho_N^h = \frac{1}{N}\sum_i^N \delta_{\theta_i^h}$

$$v(\theta_i^h, \rho_N^h, x_h, y_h) = 2\nabla V(\theta_i^h, x_h, y_h) - 2\int \nabla U(\theta_i^h, \theta', x_h, y_h)\, d\rho_N^h(\theta')$$

$$= -\nabla \psi(\theta_i^h, \rho_N^h, x_h, y_h) = -\nabla \psi(\theta_i^h, \rho_N^h) + \text{noise}$$

where $\quad \psi(\vec{\theta}, \rho, x, y) = 2V(\theta, x, y) + 2\int U(\theta, \theta', x, y)\, d\rho(\theta')$

$\quad\quad \psi(\vec{\theta}, \rho) = \mathbb{E}\left[\psi(\theta, \rho, x_h, y_h)\right] = 2V(\theta) + 2\int U(\theta, \theta')\, d\rho(\theta')$

$$= \frac{\delta R(\rho)}{\delta \rho(\theta)} = \text{"additional energy when adding one particle at } \theta \text{"}$$
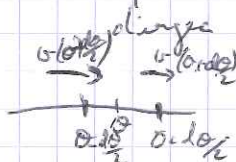
$\quad \rightarrow$ idealized particle speed: $v(\theta, \rho) = -\nabla \psi(\theta, \rho)$

$\quad$ such that $\mathbb{E}\left[v(\theta_i^h, \rho_h, x_h, y_h)\right] = v(\theta_i^h, \rho_h)$

Continuous dynamic: Continuity equation

$$(\ast)\qquad \frac{\partial \rho_t(\theta)}{\partial t} = -\nabla_\theta \cdot \left[\rho_t(\theta)\, v(\theta, \rho_t)\right] = -\nabla_\theta \cdot \left[\rho_t(\theta)\, \nabla\psi(\theta, \rho_t)\right]$$

intuition in 1D:

$$\rho_{t+dt}(\theta) - \rho_t(\theta) = \left[-\rho(\theta + \frac{d\theta}{2})\, v(\theta + \frac{d\theta}{2}) + v(\theta - \frac{d\theta}{2})\, \rho(\theta - \frac{d\theta}{2})\right] dt$$

• Gradient flow: $\dot{x}(t) = -\nabla F(x(t)) \qquad x(t+\varepsilon) = \operatorname{argmin}_z \left\{ F(z) + \frac{1}{2\varepsilon}\|z - x(t)\|^2 \right\}$

• Wasserstein metric: $\rho_{t+\varepsilon} = \operatorname{argmin}_{\rho \in \mathcal{P}_2(\theta)} \left\{ R(\rho) + \frac{1}{2\varepsilon} W_2(\rho, \rho_t)^2 \right\}$

$\quad$ where $W_2(\rho, \rho')^2 = \inf_{\gamma \in C(\rho, \rho')} \int_{\theta} \|\theta - \theta'\|^2 \gamma(d\theta, d\theta')$

## IV   Approximation with Particles

**Thm 5.1:** Assumptions: $\rho(x,0) \mapsto \sigma_*(x,0)$ bounded with sub-Gaussian gradient

$$\|\sigma_*\|_{\infty} \leq K_2 \qquad \|\nabla_x \sigma_*(x,0)\|_{\psi_2} \leq K_2$$

· Obs $|y_n| \leq K_2$

· gradient of $v$ and $U$ are bounded, Lipschitz continuous:

$$\|\nabla_\theta V(\theta)\|_2 \qquad \|\nabla_{\theta_1} U(\theta_1,\theta_2)\| \leq K_3$$

$$\|\nabla V(\theta) - \nabla V(\theta')\|_2 \leq K_3 \|\theta - \theta'\|_2$$

$$\|\nabla_{\theta_1} U(\theta_1,\theta_2) - \nabla_{\theta_1} U(\theta_1',\theta_2')\| \leq K_3 \left\| \binom{\theta_1}{\theta_2} - \binom{\theta_1'}{\theta_2'} \right\|_2$$

For $\rho \in M_1(\mathbb{R})$, consider SGD with initialization $(\theta_i^0)_{1 \leq i \leq N} \sim \rho_0$
and step size $s_n = \frac{s}{N}$. For $t \geq 0$, let $\rho_t$ be the solution of ⊛

Then there exists $C = C(K_i)$ s.t. $\forall \, \rho \in \mathcal{C}(\Theta) \times \mathbb{R} \to \mathbb{R}$ , $\|f\|_\infty \leq 1$ , $\|f\|_{\theta_i} \leq 1$

$\varepsilon \leq 1$ : 
$$\sup_{k \in [0,\frac{T}{s}]} \left| \int \frac{1}{N} \sum_i f(\theta_i^k) - \int\int f \, d\rho_{ks} \right|$$

and in particular:
$$\sup_{k \in [0,\frac{T}{s}] \cap \mathbb{N}} \left| R_N(\vec{\theta}^k) - R(\rho_{ks}) \right| \leq C e^{CT} \sqrt{\frac{1}{N} \vee \varepsilon} \left( \sqrt{D + \log N_2} \, z + z \right)$$

with proba $\geq 1 - e^{-z^2}$

$\Rightarrow$ PDE approx accurate as soon as $N \gg D$ , $\varepsilon \ll \frac{1}{D}$

speed is indep of $N$ !

example: $y_n = (2 \pm 1, y) \qquad X_n = (1 + y_n \Delta) Z_n$ when $Z_n \sim N(0,1)$
(= radial optimal classifier)

$\theta_i = w_i \binom{no \; a_i}{no \; b_i}$ then $\rho_t$ is radial, $\rho_t(|w|)$ prediction will $\frac{1}{N}\sum_{i=1}^N \delta_{|w_i|}$

can see learning fail in some examples

Analysis "Propagation of chaos" inspired of Sznitzman 91 "topics propagation of chaos.

$\Theta_i^0 \stackrel{iid}{\sim} \rho_0$

$\Theta_i^{n+1} = \Theta_i^n + \delta_n v(\Theta_i^n, \rho_N^n, x_n, y_n)$

$t = n\varepsilon$
$\delta_n = \varepsilon/2$

$\rho_0$

$\dfrac{\partial \rho_t(\theta)}{\partial t} = -\nabla_\theta \cdot (\rho_t(\theta) v(\theta, \rho_t))$

squeeze (hybrid called "non linear dynamics")

$\bar{\Theta}_i^0 \stackrel{iid}{\sim} \rho_0$

$\dfrac{d}{dt} \bar{\Theta}_i^t = v(\Theta_i^t, \rho_t) = -\psi(\Theta_i^t, \rho_t)$

$\rho_{i,\rho} : (\bar{\Theta}_i^t) \stackrel{iid}{\sim} \rho_t$

$$\| \Theta_i^{T/\varepsilon} - \bar{\Theta}_i^T \| = \int_0^T \| \psi(\bar{\Theta}_i^t, \rho_t) - \psi(\Theta_i^{\lfloor t/\varepsilon \rfloor}, \rho_N^{\lfloor t/\varepsilon \rfloor}, x_{\lfloor t/\varepsilon \rfloor}, y_{\lfloor t/\varepsilon \rfloor}) \| \, dt$$

$$\leq \int_0^T \| \psi(\bar{\Theta}_i^t, \rho_t) - \psi(\Theta_i^{\lfloor t/\varepsilon \rfloor}, \rho_N^{\lfloor t/\varepsilon \rfloor}) \| \, dt$$

$$+ \left\| \varepsilon \sum_{n=0}^{T/\varepsilon - 1} \psi(\Theta_i^{\lfloor n\varepsilon \rfloor}, \rho_N^{\lfloor n\varepsilon \rfloor}) - \psi(\Theta_i^{\lfloor n\varepsilon \rfloor}, \rho_N^{\lfloor n\varepsilon \rfloor}, x_{\lfloor n\varepsilon \rfloor}, y_{\lfloor n\varepsilon \rfloor}) \right\| \quad (3)$$

$$\leq \int_0^T \| \psi(\bar{\Theta}_i^t, \rho_t) - \psi(\Theta_i^{\lfloor t/\varepsilon \rfloor}, \rho_t) \| \, dt + \boxed{\phantom{xx}}\Theta_i \Big/ \leq k_3 \int \| \bar{\Theta}_i^t - \Theta_i^{t/\varepsilon} \| \, dt + \text{Gronwall}$$

$$+ \int_0^T \| \psi(\Theta_i^{\lfloor t/\varepsilon \rfloor}, \rho_t) - \psi(\Theta_i^{\lfloor t/\varepsilon \rfloor}, \rho_N) \| \, dt \Big/ \leq k_3 \int \frac{1}{N} \sum_p \| \bar{\Theta}_i^t - \Theta_i^{t/\varepsilon} \| \, dt \to \text{Gronwall}$$

$$+ \quad (3) \qquad \| \text{sum of \_\_\_ iords} \to \text{by Azuma-Hoeffding}$$

$$\leq \boxed{\phantom{xxx}}$$