

MEAN-FIELD THEORY OF NEURAL NETWORKS

Today: quick overview + planning rest of semester

SETTING:

- Given iid data $\{(\alpha_i, y_i)\}_{1 \leq i \leq m}$

$$\text{e.g. } y_i = f_*(\alpha_i) + \varepsilon_i$$

- covariates $\alpha_i \in \mathbb{R}^d$ $\alpha_i \stackrel{\text{iid}}{\sim} P$

- noise ε_i independent and $E[\varepsilon_i] = 0$

- Parametric class of functions $f: (\alpha, \Theta) \mapsto f(\alpha; \Theta)$

- Test error with squared loss:

$$R(\Theta) = E_{y|\alpha} \{ (y - f(\alpha; \Theta))^2 \}$$

E.g., fit $\hat{f}(\alpha, \hat{\Theta})$ by minimizing

$$\hat{R}_m(\Theta) = \frac{1}{m} \sum_{i=1}^m (y_i - f(\alpha_i; \Theta))^2$$

2-LAYERS NEURAL NETWORKS:

$$f(x; \Theta) = \frac{1}{N} \sum_{i=1}^N \sigma_*(x; \Theta_i)$$

$N = \#$ of hidden units (neurons)

$\sigma_* : \mathbb{R}^d \times \mathbb{R}^D \rightarrow \mathbb{R}$ activation function

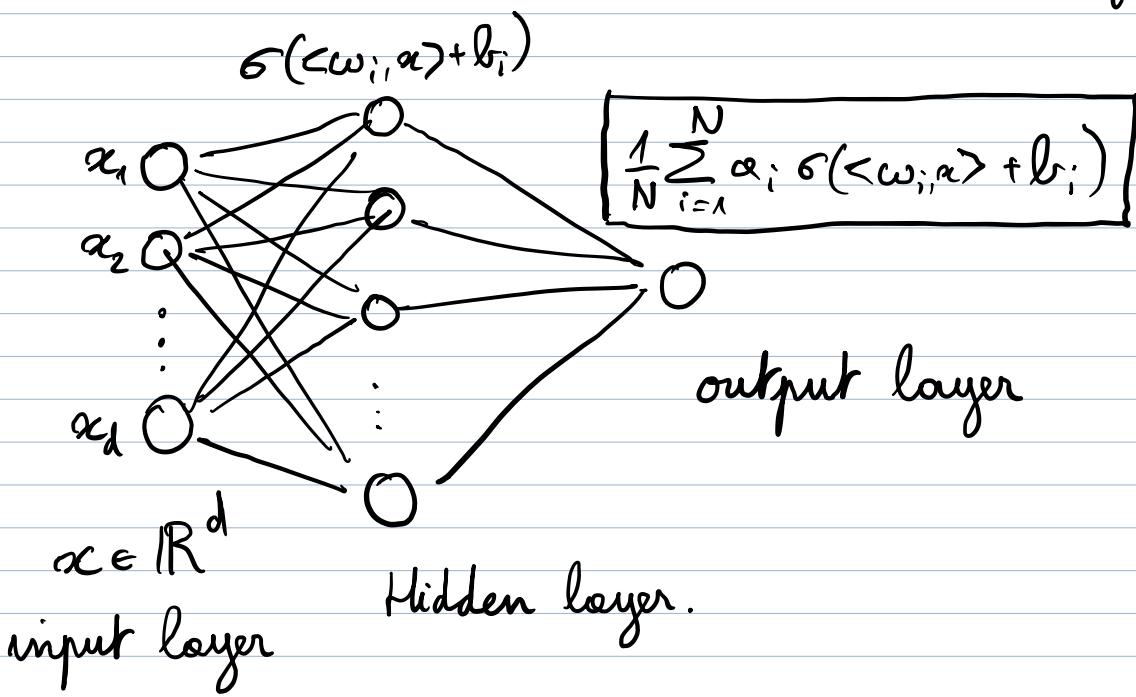
$\Theta_i \in \mathbb{R}^D$ parameters

$$\Theta = (\Theta_1, \dots, \Theta_N) \in \mathbb{R}^{ND}$$

Standard choice : $\sigma_*(x; \Theta_i) = a_i \sigma(\langle w_i, x \rangle + b_i)$

$$\Theta_i = (a_i, b_i, w_i) \in \mathbb{R}^{d+2}$$

e.g. $\sigma(x) = \max(x, 0)$, (ReLU) $\sigma(x) = \frac{1}{1 + e^{-2x}}$ (sigmoid)



Why 2-layers NNs? \Rightarrow rich enough class of functions

$$\rightarrow \mathbb{E}[f_*(x)^2] < \infty$$

Theorem [Cybenko, 1989] Take $x \sim P$, $f_* \in L^2(P)$ and

$$\sigma: \mathbb{R} \rightarrow \mathbb{R} \text{ continuous with } \begin{cases} \lim_{x \rightarrow \infty} \sigma(x) = 1 \\ \lim_{x \rightarrow -\infty} \sigma(x) = 0 \end{cases}$$

Then for any $\varepsilon > 0$, there exists $N = N(\varepsilon)$ such that

$$\inf_{\{(a_i, b_i, \omega_i)\}} \mathbb{E} \left[\left(f_*(x) - \frac{1}{N} \sum_{i=1}^N a_i \sigma(\langle \omega_i, x \rangle + b_i) \right)^2 \right] \leq \varepsilon$$

\Rightarrow 2-layers NN can approximate any reasonable functions

\Rightarrow From classical theory of universal approximation

it is often more insightful to think in terms of empirical distribution:

$$f_N(x; \Theta) = \frac{1}{N} \sum_{i=1}^N \sigma_*(x; \Theta_i) = \int \sigma_*(x; \Theta) \hat{\rho}_N(d\Theta)$$

$$=: f(x; \hat{\rho}_N) \text{ with } \hat{\rho}_N = \frac{1}{N} \sum_{i=1}^N \delta_{\Theta_i}$$

$$f(x; \rho) = \int \sigma_*(x; \theta) \rho(d\theta) \quad \text{for general distribution} \\ \rho \in \mathcal{P}(\mathbb{R}^D)$$

MEAN-FIELD LIMIT : STATICS

Preliminary: look at the population risk

$$R_N(\theta) = \mathbb{E}_{y, \alpha} \left[\left(y - \frac{1}{N} \sum_{i=1}^N \sigma_*(\alpha; \theta_i) \right)^2 \right]$$

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

where $\bullet R_\# = \mathbb{E}[y^2]$

$\bullet V(\theta) = -\mathbb{E}[y \sigma_*(\alpha; \theta)]$ ✓

$\bullet U(\theta_1, \theta_2) = \mathbb{E}[\sigma_*(\alpha; \theta_1) \sigma_*(\alpha; \theta_2)]$

Rmk: Can think about $R_N(\theta)$ as the energy of N particles in D dimensions moving in an external potential $V(\theta)$ and with pairwise potentials $U(\theta_i, \theta_j)$

→ kernel U is PSD (U is a repulsive interaction)

Again can replace $\hat{\rho}_N \rightarrow \rho \in \mathcal{P}(\mathbb{R}^D)$

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

$$U(\theta_1, \theta_2) = \int \sigma_x(x; \theta_1) \sigma_x(x; \theta_2) P(dx)$$

Properties of this risk:

① $\rho \mapsto R(\rho)$ convex function on $\mathcal{P}(\mathbb{R}^D)$
 $\frac{1}{2}(\rho_1 + \rho_2)$

② Minimizing $R_N(\theta)$ and $R(\rho)$ is not much different

Prop. Assume $\exists \varepsilon > 0$ such that for any $\rho \in \mathcal{P}(\mathbb{R}^D)$
if $R(\rho) \leq \inf_{\rho} R(\rho) + \varepsilon \Rightarrow \int U(\theta, \theta) \rho(d\theta) \leq K$

Then $|\inf_{\theta} R_N(\theta) - \inf_{\rho} R(\rho)| \leq \frac{K}{N}$

$$\inf_{\theta} R_N(\theta) \geq \inf_{\rho} R(\rho) \quad \theta_i \sim_{\text{iid}} \rho \quad |R_N(\theta) - \mathbb{E}R_N(\theta)| \leq \frac{1}{N} \int U(\theta, \theta) \rho(d\theta)$$

③ Define functional derivative $\mathbb{E}[R_N(\theta)] = R(\rho) +$

$$\Psi(\theta; \rho) \equiv \frac{1}{2} \frac{\partial R}{\partial \rho(\theta)} = V(\theta) + \int U(\theta, \theta') \rho(d\theta')$$

("additional energy of adding a single particle at $\theta \in \mathbb{R}^D$)

Global minima: distributions ρ^* such that

$$\text{support}(\rho^*) \subseteq \operatorname{argmin}_{\theta \in \mathbb{R}^D} \Psi(\theta; \rho^*)$$

(\rightarrow "Energy cannot decrease by moving mass from $\text{supp}(\rho^*)$ "
infinitesimal)

STOCHASTIC GRADIENT DESCENT

We want to minimize $R_N(\Theta)$: do SGD

$$\Theta^{k+1} = \Theta^k - \frac{\varepsilon}{2} \nabla_{\Theta} \underset{=} l(y_k, \alpha_k; \Theta_i^k)$$

→ consider one-pass over the data: (y_k, α_k) are iid

$$\text{Get } \Theta_i^{k+1} = \Theta_i^k + \underset{=} \varepsilon \nabla_{\Theta_i} \sigma_*(\alpha_k; \Theta_i^k) \left(y_k - \frac{1}{N} \sum_{i=1}^N \sigma_*(\alpha_k; \Theta_i^k) \right)$$


MEAN-FIELD LIMIT OF SGD :

DISTRIBUTIONAL DYNAMICS

We will take both $\varepsilon \rightarrow 0$ (continuous time limit)
 $N \rightarrow \infty$

SGD dynamics describes a set of N particles, with
velocity of particle i :

$$v_i^k = \frac{\Theta_i^{k+1} - \Theta_i^k}{\varepsilon} = \nabla_{\Theta_i} (y_k \sigma_*(\alpha_k; \Theta_i^k)) - \frac{1}{N} \sum_{j=1}^N (\nabla_{\Theta_i} \sigma_*(\alpha_k; \Theta_j^k)) \sigma_*(\alpha_k; \Theta_j^k)$$


→ taking expectation wrt y_k, α_k (denote F_k)

$$\mathbb{E}[v_i^k | \mathcal{F}_k] = -\underbrace{\nabla_{\theta_i} V(\theta_i^k)}_{\text{red}} - \frac{1}{N} \sum_{j=1}^N \underbrace{\nabla_{\theta_i} U(\theta_i^k, \theta_j^k)}_{\text{red}}$$

→ if at time k , density of particles $\approx \rho_{k\varepsilon}$

$$\mathbb{E}[v_i^k | \mathcal{F}_k] \approx \underbrace{v(\theta_i^k; \rho_t)}_{\text{red}} = -\underbrace{\nabla_{\theta_i} \Psi(\theta_i^k; \rho_{k\varepsilon})}_{\text{red}}$$

→ Denote $t = k\varepsilon$

→ density ρ_t should satisfy the continuity equation

$$\partial_t \rho_t(\theta) + \nabla_{\theta} \cdot (\rho_t(\theta) v(\theta; \rho_t)) = 0$$

i.e.

$$\partial_t \rho_t = \nabla_{\theta} \cdot (\rho_t \nabla_{\theta} \Psi(\theta; \rho_t))$$

DISTRIBUTIONAL DYNAMICS

Several groups in 2018 showed

[Mei, Montanari, Nguyen], [Chizat, Bach]

[Rotskoff, Vanden-Eijnden], [Sirignano, Spiliopoulos]

$$\lim_{N \rightarrow \infty} \lim_{\varepsilon \rightarrow 0} \| f(\cdot; \rho_{(N)}^{L^t/\varepsilon}) - f(\cdot; \rho_t) \|_{L^2} = 0$$

$$\lim_{N \rightarrow \infty} \lim_{\varepsilon \rightarrow 0} | R_N(\theta^{L^t/\varepsilon}) - R(\rho_t) | = 0$$

NOISY SGD

$$\theta_i^{k+1} = \theta_i^k - \frac{\varepsilon}{2} \nabla_{\theta_i} l(y_k, \alpha_k; \theta^k) + \sqrt{\varepsilon/B} g_i^k$$

$g_i^k \sim N(0, \text{Id}_D)$

Get: $\partial_t \rho_t = \nabla_{\theta} \cdot (\rho_t \nabla_{\theta} \Psi(\theta; \rho_t)) + \frac{1}{B} \Delta \rho_t$

+ diffusion term

GRADIENT FLOW

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

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

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

→ if diffusion term, GF of free energy

$$F(\rho) = \frac{1}{2} R(\rho) - \frac{1}{B} S(\rho)$$

with $S(\rho) = - \int \rho(\theta) \log \rho(\theta) d\theta$

$$\frac{1}{N} \sum_{i=1}^N \alpha_i^{(3)} \sigma(\langle w_i^{(2)}, x^{(1)} \rangle)$$

$$\alpha^{(1)} = \sigma(\underbrace{\langle w_{ij}^{(1)}, x \rangle}_{\text{ }})$$

TOPICS:

- ① If anyone is interested in presenting some mathematical topics related to MF limit
 (Gradient Flow in Wasserstein spaces etc.)

Ⓐ Convergence of SGD process to PDE

E.g., if we are interested in the risk

$$\sup_{k \in [0, T/\varepsilon] \cap \mathbb{N}} |R_N(\theta^k) - R(\rho^{k\varepsilon})| \leq \text{Err}(T, N, \varepsilon, \dots)$$

$$\text{E.g., } \sigma_*(x; \theta) = a \sigma(x; \omega) \quad v(\omega) = -\mathbb{E}[c y \sigma(x; \omega)] \\ u(\omega_1, \omega_2) = \mathbb{E}[\sigma(x; \omega_1) \sigma(x; \omega_2)]$$

and $|y|, \|\sigma\|_\infty, \|\nabla v(\omega)\|_2, \|\nabla u(\omega_1, \omega_2)\|_2, \|\nabla v(\omega)\|_{\text{Lip}}, \|\nabla u(\omega_1, \omega_2)\|_{\text{Lip}} \leq k$
 $\nabla_\omega \sigma(\cdot, \omega)$ K-subgaussian

Then with proba $\geq 1 - e^{-\beta^2}$, $\tilde{K} = \text{poly}(K)$

$$\text{Err}(T, N, \varepsilon, \dots) \leq \tilde{K} \frac{\tilde{K} T^3}{\sqrt{N}} \left\{ \frac{\sqrt{\log N} + \beta}{\sqrt{N}} + (\sqrt{D} + \log N + \beta) \sqrt{\varepsilon} \right\}$$

[Misiakiewicz, Mei, Montanari, 2019]

Proof: ~~couplings + propagation of chaos + concentration of measure~~

→ For $T = O(1)$, $N \gg \text{cste}$ that only depend on K

→ still ~~e^{KT^3}~~ (worst case analysis)

→ Can we do better?

(B) Global convergence of the PDE

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

Noisy SGD:

$$\text{unique minimizer } \rho_*(\theta) = \frac{1}{Z(T)} \exp(-\beta \Psi(\theta; \rho_*))$$

Free energy strictly decreasing along PDE path

$$F(\rho_t) \rightarrow \inf_{\rho} F(\rho)$$

$$\text{and for } \beta > \beta_*(\delta), \quad \lim_{t \rightarrow \infty} R(\rho_t) \leq \inf_{\rho} R(\rho) + \delta$$

Noiseless SGD

ρ fixed point of PDE iff $\text{support}(\rho) \subseteq \{\theta : \nabla_\theta \Psi(\theta; \rho) = 0\}$

[Chizat, Bach, 2018]

Morse-Sard type condition
+ good initialization ↪ $\Rightarrow \rho_\infty$ is a global
+ $\rho_t \rightarrow \rho_\infty$ in W_2 optimizer

[Wojtowytch, 2020]

extends these results

[Nguyen, Pham, 2020]

replace Morse-Sard type condition by a topology argument

C Rates of convergence of the PDE to minimizer

E.g. in noisy SGD.

Important to transfer results to practical settings

[Javanmard, Mondelli, Montanari, 2019]

$$\sigma_*(\alpha; \omega) = \underline{K^\delta(\alpha - \omega_i)}$$

$$K^\delta(\alpha) = \underline{s^{-d} K(\alpha/s)}$$

$K: \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$ 1st order kernel
with compact support

Target function: $\underline{f_*}$ α -strongly concave + smooth

Then $R(\rho)$ is ~~α~~ -displacement convex

$$R(\rho_t) \leq R(\rho_0) e^{-2\alpha t}$$



\Rightarrow In general, $R(\rho)$ not displacement convex.

① Representation vs Optimization

① Representation: $\frac{1}{N} \sum_{i=1}^N a_i \sigma(x; \omega_i) \rightarrow f(x)$

② Optimization: $f(x; \rho_t) \rightarrow f(x)$

Is it easier to optimize for some class of functions?

E.g. What are easier to represent?

Barron - norm: $f(x) = \int_{\Omega} \alpha \sigma(\omega^T \alpha) \rho(d\alpha, d\omega)$

$$= \mathbb{E}_{\rho} [\alpha \sigma(\omega^T \alpha)]$$

$$\|f\|_B = \inf_{\rho: f = \mathbb{E}_{\rho}} \mathbb{E}_{\rho} [|\alpha| \| \omega \|_1]$$

Prop: We have

$$\inf_{f_N: N \text{ neurons}} \|f_* - f_N\|_{L^2} \lesssim \frac{\|f\|_B}{\sqrt{N}}$$

\exists 1-lipschitz function such that.

$$(*) \quad \inf_{f_N: N \text{ neurons}} \|f_* - f_N\|_{L^2} \geq \frac{\delta}{N^{1/d}}$$

$$\|f(\cdot, \rho_t)\|_B \leq \int |\theta|^2 \rho_t(d\theta) = o(t) \quad (**)$$

(second moment \nearrow slower than linearly)

$$(*) + (**) \quad R(\rho^t) \geq c t^{-\frac{c}{d+2}}$$

[Wojkowycz, Weinan, 2020]

(E) Multilayer NNs

Several attempts

[Araujo, Oliveira, Yukimura, 2019]

[Nguyen, Pham, 2020] neuronal embedding

→ degeneracy of the middle layers that collapse into one neuron

→ basically can factorize distribution weights

$$P^t = P_{(0)}^t \otimes P_{(1)}^t \otimes \dots \otimes P_{(L)}^t$$

[Feng, Lee, Yang, Zheng, 2020]

attempt to overcome degeneracy using new representation "neural feature flow"

(F) NTK vs mean-field limit

→ Might be of general interest

Can we compare mean-field and NTK solutions?

Can we study the transition between the two?

[Misiakiewicz, Mei, Montanari, 2019]

NTK appears at beginning of dynamics.

[Chizat, Bach, 2020]

MF find F_1 -max margin solution
 \neq NTK — F_2 —

- How many weeks?
 - Can invite speakers