

Generalization of Scaled Deep ResNets in the Mean-Field Regime

Yihang Chen ¹ Fanghui Liu² Yiping Lu³ Grigorios G Chrysos ⁴ Volkan Cevher ¹

¹ LIONS, EPFL ²University of Warwick

³New York University ⁴University of Wisconsin-Madison

Published as a conference paper at ICLR 2024 (Spotlight).

lions@epfl

EPFL



ICLR

OverView

Question:

- ▶ Can we build a generalization analysis of trained Deep ResNets in the mean-field setting?

Contributions:

- ▶ The first **minimum eigenvalue estimation** (lower bound) of the Gram matrix of the gradients for deep ResNet parameterized by the ResNet encoder's parameters and MLP predictor's parameters in the mean-field regime.
- ▶ The paper proves that the **KL divergence** of feature encoder ν and output layer ν can be bounded by a constant (depending only on network architecture parameters) during the training, which facilitates our generalization analysis.
- ▶ This paper builds the connection between the Rademacher complexity result and KL divergence, and then derive the **convergence rate** $\mathcal{O}(1/\sqrt{n})$ **for generalization**.

Problem Settings

Basic Settings:

- ▶ The training set $\mathcal{D}_n = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ is drawn from an unknown distribution μ on $\mathcal{X} \times \mathcal{Y}$, and μ_X is the marginal distribution of μ over \mathcal{X} .
- ▶ We consider a binary classification task, denoted by minimizing the expected risk, let $\ell_{0-1}(f, y) := \mathbb{1}\{yf < 0\}$.
- ▶ We employ the squared loss in ERM in training, i.e, $\ell(f, y) := \frac{1}{2}(y - f)^2$.
- ▶ The hypothesis f is parameterized by the ResNet feature encoder and a non-linear predictor, $f_{\tau, \nu}$. The empirical loss $\widehat{L}(\tau, \nu) := \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_n} \ell(f_{\tau, \nu}(\mathbf{x}), y(\mathbf{x}))$.

Problem Settings

Network Structure: (α, β will be determined later)

- ▶ Discrete

$$z_{l+1}(\mathbf{x}) = z_l(\mathbf{x}) + \frac{\alpha}{ML} \sum_{m=1}^M \sigma(z_l(\mathbf{x}), \boldsymbol{\theta}_{l,m}) \in \mathbb{R}^d, \quad l \in [L-1], \quad (1)$$

$$f_{\Omega_K, \Theta_{L,M}}(\mathbf{x}) = \frac{\beta}{K} \sum_{k=1}^K h(z_L, \boldsymbol{\omega}_k) \in \mathbb{R},$$

- ▶ The following ODE models the infinite depth infinite width ResNet.

$$\frac{dz(\mathbf{x}, s)}{ds} = \alpha \cdot \int_{\mathbb{R}^{k_\nu}} \sigma(z(\mathbf{x}, s), \boldsymbol{\theta}) d\nu(\boldsymbol{\theta}, s), \quad s \in [0, 1], \quad z(\mathbf{x}, 0) = \mathbf{x}. \quad (2)$$

We denote the solution of Equation (2) as $Z_\nu(\mathbf{x}, s)$.

- ▶ The whole network can be written as

$$f_{\tau, \nu}(\mathbf{x}) := \beta \cdot \int_{\mathbb{R}^{k_\tau}} h(Z_\nu(\mathbf{x}, 1), \boldsymbol{\omega}) d\tau(\boldsymbol{\omega}),$$

Assumptions

Assumption (Assumptions on data)

We assume that for $\mathbf{x}_i \neq \mathbf{x}_j \sim \mu_X$, the following holds with probability 1,

$$\|\mathbf{x}_i\|_2 = 1, |y(\mathbf{x}_i)| \leq 1, \langle \mathbf{x}_i, \mathbf{x}_j \rangle \leq C_{\max} < 1, \forall i, j \in [n].$$

Assumption (Assumption on initialization)

The initial distribution τ_0, ν_0 is standard Gaussian: $(\tau_0, \nu_0)(\boldsymbol{\omega}, \boldsymbol{\theta}, s) \propto \exp\left(-\frac{\|\boldsymbol{\omega}\|_2^2 + \|\boldsymbol{\theta}\|_2^2}{2}\right), \forall s \in [0, 1]$.

Assumption (Assumptions on activation σ, h)

Let $\boldsymbol{\theta} := (\mathbf{u}, \mathbf{w}, b) \in \mathbb{R}^{k_\nu}$, where $\mathbf{u}, \mathbf{w} \in \mathbb{R}^{k_\nu}, b \in \mathbb{R}$, i.e. $k_\nu = 2d + 1$; $\boldsymbol{\omega} := (a, \mathbf{w}, b) \in \mathbb{R}^{k_\tau}$, where $\mathbf{w} \in \mathbb{R}^{k_\nu}, a, b \in \mathbb{R}$, i.e. $k_\tau = d + 2$. For any $\mathbf{z} \in \mathbb{R}^{k_\nu}$, we assume

$$\sigma(\mathbf{z}, \boldsymbol{\theta}) = \mathbf{u}\sigma_0(\mathbf{w}^\top \mathbf{z} + b), \quad h(\mathbf{z}, \boldsymbol{\omega}) = a\sigma_0(\mathbf{w}^\top \mathbf{z} + b), \quad \sigma_0 : \mathbb{R} \rightarrow \mathbb{R}.$$

In addition, we have the following assumption on σ_0 . $|\sigma_0(x)| \leq C_1 \max(|x|, 1), |\sigma_0'(x)| \leq C_1, |\sigma_0''(x)| \leq C_1$, and let $\mu_i(\sigma_0)$ be the i -th Hermite coefficient of σ_0 .

Gradient Evolution

- ▶ The evolution of the ResNet layers $\nu(\boldsymbol{\theta}, s)$ can be characterized as

$$\frac{\partial \nu}{\partial t}(\boldsymbol{\theta}, s, t) = \nabla_{\boldsymbol{\theta}} \cdot \left(\nu(\boldsymbol{\theta}, s, t) \nabla_{\boldsymbol{\theta}} \frac{\delta \widehat{L}(\tau, \nu)}{\delta \nu}(\boldsymbol{\theta}, s, t) \right), \quad t \geq 0, \quad (3)$$

- ▶ The evolution of the final layer distribution $\tau(\boldsymbol{\omega})$ can be characterized as

$$\frac{\partial \tau}{\partial t}(\boldsymbol{\omega}, t) = \nabla_{\boldsymbol{\omega}} \cdot \left(\tau(\boldsymbol{\omega}, t) \nabla_{\boldsymbol{\omega}} \frac{\delta \widehat{L}(\tau, \nu)}{\delta \tau}(\boldsymbol{\omega}, t) \right), \quad t \geq 0, \quad (4)$$

where the functional derivative

$$\frac{\delta \widehat{L}(\tau, \nu)}{\delta \tau}(\boldsymbol{\omega}) = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_n} [\beta \cdot (f_{\tau, \nu}(\mathbf{x}) - y(\mathbf{x})) \cdot h(\mathbf{Z}_{\nu}(\mathbf{x}, 1), \boldsymbol{\omega})].$$

Gram Matrix

- ▶ We define one Gram matrix for the ResNet layers, $\mathbf{G}_1(\tau, \nu)$ by

$$\mathbf{G}_1(\tau, \nu) = \int_0^1 \mathbf{G}_1(\tau, \nu, s) ds$$
$$\mathbf{G}_1(\tau, \nu, s) = \mathbb{E}_{\boldsymbol{\theta} \sim \nu(\cdot, s)} \mathbf{J}_1(\tau, \nu, \boldsymbol{\theta}, s) \mathbf{J}_1(\tau, \nu, \boldsymbol{\theta}, s)^\top .$$

- ▶ We define the Gram matrix for the MLP parameter distribution τ , $\mathbf{G}_2(\tau, \nu)$ by $\mathbf{G}_2(\tau, \nu) = \mathbb{E}_{\boldsymbol{\omega} \sim \tau(\cdot)} \mathbf{J}_2(\nu, \boldsymbol{\omega}) \mathbf{J}_2(\nu, \boldsymbol{\omega})^\top$, where the row vector of \mathbf{J}_2 is defined as

$$(\mathbf{J}_2(\nu, \boldsymbol{\omega}))_{i,\cdot} = \nabla_{\boldsymbol{\omega}} h(\mathbf{Z}_\nu(\mathbf{x}_i, 1), \boldsymbol{\omega}), \quad 1 \leq i \leq n .$$

- ▶ The Gram matrix for the whole network is $\mathbf{G} = \alpha^2 \mathbf{G}_1 + \mathbf{G}_2$.

Minimum Eigenvalue

Lemma

There exist a constant $\Lambda := \Lambda(d)$, only depending on the dimension d , such that $\lambda_{\min}[\mathbf{G}(\tau_0, \nu_0)]$ is lower bounded by

$$\lambda_0 := \lambda_{\min}(\mathbf{G}(\tau_0, \nu_0)) \geq \lambda_{\min}(\mathbf{G}_2(\tau_0, \nu_0)) \geq \Lambda(d).$$

Theorem

Assume the PDE Eqn. 4 has solution $\tau_t \in \mathcal{P}^2$, and the PDE Eqn. 3 has solution $\nu_t \in \mathcal{C}(\mathcal{P}^2; [0, 1])$. Under Assumption 1, 2, 3, for some constant C_{KL} dependent on d, α , taking $\bar{\beta} := \frac{\beta}{n} > \frac{4\sqrt{C_{\text{KL}}(d, \alpha)}}{\Lambda r_{\max}}$, the following results hold for all $t \in [0, \infty)$:

$$\widehat{L}(\tau_t, \nu_t) \leq e^{-\frac{\beta^2 \Lambda}{2n} t} \widehat{L}(\tau_0, \nu_0), \quad \text{KL}(\tau_t \| \tau_0) \leq \frac{C_{\text{KL}}(d, \alpha)}{\Lambda^2 \bar{\beta}^2}, \quad \text{KL}(\nu_t \| \nu_0) \leq \frac{C_{\text{KL}}(d, \alpha)}{\Lambda^2 \bar{\beta}^2}.$$

where the radius r_{\max} is defined such that if $\nu \in \mathcal{C}(\mathcal{P}^2; [0, 1])$, $\tau \in \mathcal{P}^2$, $\max\{\mathcal{W}_2(\nu, \nu_0), \mathcal{W}_2(\tau, \tau_0)\} \leq r_{\max}$, we have $\lambda_{\min}(\mathbf{G}_2(\tau, \nu)) \geq \frac{\lambda_0}{2}$.

Generalization

Theorem (Generalization)

Assume $\tau_y \in \mathcal{C}(\mathcal{P}^2; [0, 1])$ and $\nu_y \in \mathcal{P}^2$ be the ground truth distributions, such that, $y(\mathbf{x}) = \mathbb{E}_{\omega \sim \tau_y} h(\mathbf{Z}_{\nu_y}(\mathbf{x}, 1), \omega)$. Under the Assumption 1, 2 and 3, we set $\beta > \Omega(\sqrt{n})$. For any $\delta > 0$, with probability at least $1 - \delta$, the following bound holds:

$$\mathbb{E}_{\mathbf{x} \sim \mu_X} \ell_{0-1}(f_{\tau_*, \nu_*}(\mathbf{x}), y(\mathbf{x})) \lesssim 1/\sqrt{n} + 6 \sqrt{\log(2/\delta)/2n},$$

where \lesssim hides the constant dependence on d, α .