

Mean-field theory of two-layers neural networks: dimension-free bounds and kernel limit

Song Mei, Theodor Misiakiewicz, and Andrea Montanari

Stanford University

June 26, 2019

COLT 2019

Gradient dynamics of two-layers neural network

- ▶ Two layers neural network:

$$\Theta = (\theta_1, \dots, \theta_N), \quad \theta_i = (a_i, \mathbf{w}_i) \in \mathbb{R}^D.$$

$$\hat{y}(\mathbf{x}; \Theta) = \frac{1}{N} \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle).$$

- ▶ Risk function:

$$R_N(\Theta) = \mathbb{E}_{\mathbf{x}, y} \left[\left(y - \frac{1}{N} \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, \mathbf{x} \rangle) \right)^2 \right].$$

- ▶ SGD/gradient flow:

$$\Theta^{k+1} = \Theta^k - \eta_k \nabla \ell_N(\Theta^k; \mathbf{x}_k, y_k),$$

$$\frac{d}{dt} \Theta^t = - \nabla R_N(\Theta^t).$$

Gradient dynamics of two-layers neural network

- ▶ Two layers neural network:

$$\Theta = (\theta_1, \dots, \theta_N), \quad \theta_i = (a_i, \mathbf{w}_i) \in \mathbb{R}^D.$$

$$\hat{y}(x; \Theta) = \frac{1}{N} \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, x \rangle).$$

- ▶ Risk function:

$$R_N(\Theta) = \mathbb{E}_{x,y} \left[\left(y - \frac{1}{N} \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, x \rangle) \right)^2 \right].$$

- ▶ SGD/gradient flow:

$$\Theta^{k+1} = \Theta^k - \eta_k \nabla \ell_N(\Theta^k; x_k, y_k),$$

$$\frac{d}{dt} \Theta^t = - \nabla R_N(\Theta^t).$$

Gradient dynamics of two-layers neural network

- ▶ Two layers neural network:

$$\Theta = (\theta_1, \dots, \theta_N), \quad \theta_i = (a_i, \mathbf{w}_i) \in \mathbb{R}^D.$$

$$\hat{y}(x; \Theta) = \frac{1}{N} \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, x \rangle).$$

- ▶ Risk function:

$$R_N(\Theta) = \mathbb{E}_{x,y} \left[\left(y - \frac{1}{N} \sum_{i=1}^N a_i \sigma(\langle \mathbf{w}_i, x \rangle) \right)^2 \right].$$

- ▶ SGD/gradient flow:

$$\Theta^{k+1} = \Theta^k - \eta_k \nabla \ell_N(\Theta^k; \mathbf{x}_k, y_k),$$

$$\frac{d}{dt} \Theta^t = - \nabla R_N(\Theta^t).$$

Two-layers neural networks

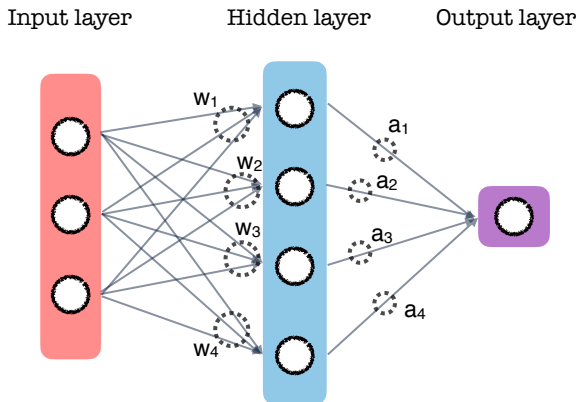


Figure: Architecture for $N = 4$. $\theta_i = (a_i, w_i)$

Related literatures

- ▶ Mean field distributional dynamics:

$$\partial_t \rho_t(\boldsymbol{\theta}) = \nabla \cdot (\nabla \Psi(\boldsymbol{\theta}; \rho_t) \rho_t).$$

- ▶ Non-linear dynamics. Converges in some cases.
- ▶ [Mei, Montanari, Nguyen, 2018], [Rotskoff and Vanden-Eijnden, 2018], [Chizat and Bach, 2018a], [Sirignano and Spiliopoulos, 2018].

- ▶ Neural tangent kernel (NTK) dynamics:

$$\partial_t \|u_t\|_2^2 = -\langle u_t, \mathcal{H}u_t \rangle.$$

- ▶ Linear dynamics. Always converges to 0 empirical risk.
- ▶ [Jacot, Gabriel, and Clement, 2018], [Li and Liang, 2018], [Du, Zhai, Póczos, Singh, 2018].

Related literatures

- ▶ Mean field distributional dynamics:

$$\partial_t \rho_t(\boldsymbol{\theta}) = \nabla \cdot (\nabla \Psi(\boldsymbol{\theta}; \rho_t) \rho_t).$$

- ▶ Non-linear dynamics. Converges in some cases.
- ▶ [Mei, Montanari, Nguyen, 2018], [Rotskoff and Vanden-Eijnden, 2018], [Chizat and Bach, 2018a], [Sirignano and Spiliopoulos, 2018].

- ▶ Neural tangent kernel (NTK) dynamics:

$$\partial_t \|u_t\|_2^2 = -\langle u_t, \mathcal{H}u_t \rangle.$$

- ▶ Linear dynamics. Always converges to 0 empirical risk.
- ▶ [Jacot, Gabriel, and Clement, 2018], [Li and Liang, 2018], [Du, Zhai, Póczos, Singh, 2018].

Related literatures

- ▶ Mean field distributional dynamics:

$$\partial_t \rho_t(\boldsymbol{\theta}) = \nabla \cdot (\nabla \Psi(\boldsymbol{\theta}; \rho_t) \rho_t).$$

- ▶ Non-linear dynamics. Converges in some cases.
- ▶ [Mei, Montanari, Nguyen, 2018], [Rotskoff and Vanden-Eijnden, 2018], [Chizat and Bach, 2018a], [Sirignano and Spiliopoulos, 2018].

- ▶ Neural tangent kernel (NTK) dynamics:

$$\partial_t \|u_t\|_2^2 = -\langle u_t, \mathcal{H}u_t \rangle.$$

- ▶ Linear dynamics. Always converges to 0 empirical risk.
- ▶ [Jacot, Gabriel, and Clement, 2018], [Li and Liang, 2018], [Du, Zhai, Póczos, Singh, 2018].

Related literatures

- ▶ Mean field distributional dynamics:

$$\partial_t \rho_t(\boldsymbol{\theta}) = \nabla \cdot (\nabla \Psi(\boldsymbol{\theta}; \rho_t) \rho_t).$$

- ▶ Non-linear dynamics. Converges in some cases.
- ▶ [Mei, Montanari, Nguyen, 2018], [Rotskoff and Vanden-Eijnden, 2018], [Chizat and Bach, 2018a], [Sirignano and Spiliopoulos, 2018].

- ▶ Neural tangent kernel (NTK) dynamics:

$$\partial_t \| \mathbf{u}_t \|_2^2 = - \langle \mathbf{u}_t, \mathcal{H} \mathbf{u}_t \rangle.$$

- ▶ Linear dynamics. Always converges to 0 empirical risk.
- ▶ [Jacot, Gabriel, and Clement, 2018], [Li and Liang, 2018], [Du, Zhai, Póczos, Singh, 2018].

Related literatures

- ▶ Mean field distributional dynamics:

$$\partial_t \rho_t(\boldsymbol{\theta}) = \nabla \cdot (\nabla \Psi(\boldsymbol{\theta}; \rho_t) \rho_t).$$

- ▶ Non-linear dynamics. Converges in some cases.
- ▶ [Mei, Montanari, Nguyen, 2018], [Rotskoff and Vanden-Eijnden, 2018], [Chizat and Bach, 2018a], [Sirignano and Spiliopoulos, 2018].

- ▶ Neural tangent kernel (NTK) dynamics:

$$\partial_t \| \mathbf{u}_t \|_2^2 = - \langle \mathbf{u}_t, \mathcal{H} \mathbf{u}_t \rangle.$$

- ▶ Linear dynamics. Always converges to 0 empirical risk.
- ▶ [Jacot, Gabriel, and Clement, 2018], [Li and Liang, 2018], [Du, Zhai, Póczos, Singh, 2018].

Related literatures

- ▶ Mean field distributional dynamics:

$$\partial_t \rho_t(\boldsymbol{\theta}) = \nabla \cdot (\nabla \Psi(\boldsymbol{\theta}; \rho_t) \rho_t).$$

- ▶ Non-linear dynamics. Converges in some cases.
- ▶ [Mei, Montanari, Nguyen, 2018], [Rotskoff and Vanden-Eijnden, 2018], [Chizat and Bach, 2018a], [Sirignano and Spiliopoulos, 2018].

- ▶ Neural tangent kernel (NTK) dynamics:

$$\partial_t \| \mathbf{u}_t \|_2^2 = - \langle \mathbf{u}_t, \mathcal{H} \mathbf{u}_t \rangle.$$

- ▶ Linear dynamics. Always converges to 0 empirical risk.
- ▶ [Jacot, Gabriel, and Clement, 2018], [Li and Liang, 2018], [Du, Zhai, Póczos, Singh, 2018].

This work

- (a) Improved bound for SGD - PDE interpolation.
- (b) Relationship of the mean field limit and the kernel limit.

SGD and distributional dynamics (DD)

- ▶ SGD for Θ^k , with $(\mathbf{x}_k, \mathbf{y}_k) \sim \mathbb{P}_{\mathbf{x}, \mathbf{y}}$, $i \in [N]$,

$$\theta_i^{k+1} = \theta_i^k - 2s_k N \nabla_{\theta_i} \ell_N(\Theta^k; \mathbf{x}_k, \mathbf{y}_k). \quad (\text{SGD})$$

- ▶ [MMN18]: $s_k = \varepsilon \hat{\xi}(k\varepsilon)$, $k = t/\varepsilon$, $N \rightarrow \infty$, $\varepsilon \rightarrow 0$:

$$\hat{\rho}_k^{(N)} \equiv \frac{1}{N} \sum_{i=1}^N \delta_{\theta_i^k} \Rightarrow \rho_t \in \mathcal{P}(\mathbb{R}^D) \times [0, \infty).$$

- ▶ Distributional dynamics (DD) for ρ_t ,

$$\partial_t \rho_t(\theta) = 2\xi(t) \nabla_{\theta} \cdot (\rho_t(\theta) \nabla_{\theta} \Psi(\theta; \rho_t)), \quad (\text{DD})$$

where

$$\Psi(\theta; \rho) = \frac{\delta R(\rho)}{\delta \rho(\theta)} = V(\theta) + \int U(\theta, \theta') \rho(d\theta').$$

SGD and distributional dynamics (DD)

- ▶ SGD for Θ^k , with $(\mathbf{x}_k, \mathbf{y}_k) \sim \mathbb{P}_{\mathbf{x}, \mathbf{y}}$, $i \in [N]$,

$$\theta_i^{k+1} = \theta_i^k - 2s_k N \nabla_{\theta_i} \ell_N(\Theta^k; \mathbf{x}_k, \mathbf{y}_k). \quad (\text{SGD})$$

- ▶ [MMN18]: $s_k = \varepsilon \xi(k\varepsilon)$, $k = t/\varepsilon$, $N \rightarrow \infty$, $\varepsilon \rightarrow 0$:

$$\hat{\rho}_k^{(N)} \equiv \frac{1}{N} \sum_{i=1}^N \delta_{\theta_i^k} \Rightarrow \rho_t \in \mathcal{P}(\mathbb{R}^D) \times [0, \infty).$$

- ▶ Distributional dynamics (DD) for ρ_t ,

$$\partial_t \rho_t(\theta) = 2\xi(t) \nabla_{\theta} \cdot (\rho_t(\theta) \nabla_{\theta} \Psi(\theta; \rho_t)), \quad (\text{DD})$$

where

$$\Psi(\theta; \rho) = \frac{\delta R(\rho)}{\delta \rho(\theta)} = V(\theta) + \int U(\theta, \theta') \rho(d\theta').$$

SGD and distributional dynamics (DD)

- ▶ SGD for Θ^k , with $(\mathbf{x}_k, \mathbf{y}_k) \sim \mathbb{P}_{\mathbf{x}, \mathbf{y}}$, $i \in [N]$,

$$\theta_i^{k+1} = \theta_i^k - 2s_k N \nabla_{\theta_i} \ell_N(\Theta^k; \mathbf{x}_k, \mathbf{y}_k). \quad (\text{SGD})$$

- ▶ [MMN18]: $s_k = \varepsilon \xi(k\varepsilon)$, $k = t/\varepsilon$, $N \rightarrow \infty$, $\varepsilon \rightarrow 0$:

$$\hat{\rho}_k^{(N)} \equiv \frac{1}{N} \sum_{i=1}^N \delta_{\theta_i^k} \Rightarrow \rho_t \in \mathcal{P}(\mathbb{R}^D) \times [0, \infty).$$

- ▶ Distributional dynamics (DD) for ρ_t ,

$$\partial_t \rho_t(\boldsymbol{\theta}) = 2\xi(t) \nabla_{\boldsymbol{\theta}} \cdot (\rho_t(\boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}} \Psi(\boldsymbol{\theta}; \rho_t)), \quad (\text{DD})$$

where

$$\Psi(\boldsymbol{\theta}; \rho) = \frac{\delta R(\rho)}{\delta \rho(\boldsymbol{\theta})} = V(\boldsymbol{\theta}) + \int U(\boldsymbol{\theta}, \boldsymbol{\theta}') \rho(d\boldsymbol{\theta}').$$

An improved bound

Assumption

(i) σ bounded; (ii) $\nabla_{\mathbf{w}}\sigma(\langle \mathbf{x}, \mathbf{w} \rangle)$ sub-Gaussian; (iii) $\nabla\Psi$ bdd. Lipschitz.

Theorem (M., Misiakiwicz, Montanari, 2019)

Let $(\theta_i^0)_{i \leq N} \sim_{iid} \rho_0$. Then, $\forall f$ bounded Lipschitz, w.h.p,

$$\sup_{t \leq T} \left| \frac{1}{N} \sum_{i=1}^N f(\theta_i^{\lfloor t/\varepsilon \rfloor}) - \int f(\theta) \rho_t(\theta) \right| \leq \text{Func}(T) \cdot \sqrt{\frac{1}{N}} \vee D\varepsilon.$$

An example: learning a spherically symmetric Lipschitz function using $N = O_d(1)$ neurons and $n = O_d(d)$ samples.

Caveat: this improved bound is not strong. In other cases the factor $\text{Func}(T)$ could potentially be huge.

An improved bound

Assumption

(i) σ bounded; (ii) $\nabla_{\mathbf{w}}\sigma(\langle \mathbf{x}, \mathbf{w} \rangle)$ sub-Gaussian; (iii) $\nabla\Psi$ bdd. Lipschitz.

Theorem (M., Misiakiwicz, Montanari, 2019)

Let $(\boldsymbol{\theta}_i^0)_{i \leq N} \sim_{iid} \rho_0$. Then, $\forall f$ bounded Lipschitz, w.h.p,

$$\sup_{t \leq T} \left| \frac{1}{N} \sum_{i=1}^N f(\boldsymbol{\theta}_i^{\lfloor t/\varepsilon \rfloor}) - \int f(\boldsymbol{\theta}) \rho_t(\boldsymbol{\theta}) \right| \leq \text{Func}(T) \cdot \sqrt{\frac{1}{N}} \vee D\varepsilon.$$

An example: learning a spherically symmetric Lipschitz function using $N = O_d(1)$ neurons and $n = O_d(d)$ samples.

Caveat: this improved bound is not strong. In other cases the factor $\text{Func}(T)$ could potentially be huge.

An improved bound

Assumption

(i) σ bounded; (ii) $\nabla_{\mathbf{w}}\sigma(\langle \mathbf{x}, \mathbf{w} \rangle)$ sub-Gaussian; (iii) $\nabla\Psi$ bdd. Lipschitz.

Theorem (M., Misiakiwicz, Montanari, 2019)

Let $(\boldsymbol{\theta}_i^0)_{i \leq N} \sim_{iid} \rho_0$. Then, $\forall f$ bounded Lipschitz, w.h.p,

$$\sup_{t \leq T} \left| \frac{1}{N} \sum_{i=1}^N f(\boldsymbol{\theta}_i^{\lfloor t/\varepsilon \rfloor}) - \int f(\boldsymbol{\theta}) \rho_t(\boldsymbol{\theta}) \right| \leq \text{Func}(T) \cdot \sqrt{\frac{1}{N}} \vee D\varepsilon.$$

An example: learning a spherically symmetric Lipschitz function using $N = O_d(1)$ neurons and $n = O_d(d)$ samples.

Caveat: this improved bound is not strong. In other cases the factor $\text{Func}(T)$ could potentially be huge.

An improved bound

Assumption

(i) σ bounded; (ii) $\nabla_{\mathbf{w}}\sigma(\langle \mathbf{x}, \mathbf{w} \rangle)$ sub-Gaussian; (iii) $\nabla\Psi$ bdd. Lipschitz.

Theorem (M., Misiakiwicz, Montanari, 2019)

Let $(\boldsymbol{\theta}_i^0)_{i \leq N} \sim_{iid} \rho_0$. Then, $\forall f$ bounded Lipschitz, w.h.p,

$$\sup_{t \leq T} \left| \frac{1}{N} \sum_{i=1}^N f(\boldsymbol{\theta}_i^{\lfloor t/\varepsilon \rfloor}) - \int f(\boldsymbol{\theta}) \rho_t(\boldsymbol{\theta}) \right| \leq \text{Func}(T) \cdot \sqrt{\frac{1}{N}} \vee D\varepsilon.$$

An example: learning a spherically symmetric Lipschitz function using $N = O_d(1)$ neurons and $n = O_d(d)$ samples.

Caveat: this improved bound is not strong. In other cases the factor $\text{Func}(T)$ could potentially be huge.

This work

- (a) Improved bound for SGD - PDE interpolation.
- (b) Relationship of the mean field limit and the kernel limit.

Recovering the kernel limit

Same idea appeared in [Chizat and Bach, 2018b], where the kernel limit was called “lazy training”.

Setup:

Prediction function:
$$\hat{f}_{\alpha,N}(\mathbf{x}; \boldsymbol{\theta}) = \frac{\alpha}{N} \sum_{j=1}^N \sigma_{\star}(\mathbf{x}; \boldsymbol{\theta}_j),$$

Risk function:
$$R_{\alpha,N}(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x}} \left[\left(f(\mathbf{x}) - \hat{f}_{\alpha,N}(\mathbf{x}; \boldsymbol{\theta}) \right)^2 \right],$$

Gradient flow:
$$\frac{d\boldsymbol{\theta}_j^t}{dt} = - \frac{N}{2\alpha^2} \nabla_{\boldsymbol{\theta}_j} R_{\alpha,N}(\boldsymbol{\theta}^t).$$

The coupled dynamics

Denote $\rho_t^{\alpha, N} = (1/N) \sum_{j=1}^N \delta_{\theta_j^t}$. Distributional dynamics:

$$\partial_t \rho_t^{\alpha, N} = (1/\alpha) \nabla_{\theta} \cdot (\rho_t^{\alpha, N} \nabla_{\theta} \Psi(\theta; \rho_t^{\alpha, N})).$$

Denote $u_t^{\alpha, N}(z) = f(z) - \hat{f}_{\alpha, N}(z; \theta^t)$. Residual dynamics:

$$\partial_t \|u_t^{\alpha, N}\|_{L^2}^2 = -\langle u_t^{\alpha, N}, \mathcal{H}_{\rho_t^{\alpha, N}} u_t^{\alpha, N} \rangle.$$

Here

$$\begin{aligned} \mathcal{H}_{\rho}(x, z) &\equiv \int \langle \nabla_{\theta} \sigma_{*}(x; \theta), \nabla_{\theta} \sigma_{*}(z; \theta) \rangle \rho(d\theta), \\ \Psi_{\alpha}(\theta; \rho^{\alpha, N}) &= -\mathbb{E}_x[u_t^{\alpha, N}(x) \sigma_{*}(x; \theta)]. \end{aligned}$$

The coupled dynamics

Denote $\rho_t^{\alpha, N} = (1/N) \sum_{j=1}^N \delta_{\theta_j^t}$. Distributional dynamics:

$$\partial_t \rho_t^{\alpha, N} = (1/\alpha) \nabla_{\theta} \cdot (\rho_t^{\alpha, N} \nabla_{\theta} \Psi(\theta; \rho_t^{\alpha, N})).$$

Denote $u_t^{\alpha, N}(z) = f(z) - \hat{f}_{\alpha, N}(z; \theta^t)$. Residual dynamics:

$$\partial_t \|u_t^{\alpha, N}\|_{L^2}^2 = -\langle u_t^{\alpha, N}, \mathcal{H}_{\rho_t^{\alpha, N}} u_t^{\alpha, N} \rangle.$$

Here

$$\mathcal{H}_{\rho}(x, z) \equiv \int \langle \nabla_{\theta} \sigma_{*}(x; \theta), \nabla_{\theta} \sigma_{*}(z; \theta) \rangle \rho(d\theta),$$
$$\Psi_{\alpha}(\theta; \rho^{\alpha, N}) = -\mathbb{E}_x[u_t^{\alpha, N}(x) \sigma_{*}(x; \theta)].$$

The coupled dynamics

Denote $\rho_t^{\alpha, N} = (1/N) \sum_{j=1}^N \delta_{\theta_j^t}$. Distributional dynamics:

$$\partial_t \rho_t^{\alpha, N} = (1/\alpha) \nabla_{\theta} \cdot (\rho_t^{\alpha, N} \nabla_{\theta} \Psi(\theta; \rho_t^{\alpha, N})).$$

Denote $u_t^{\alpha, N}(z) = f(z) - \hat{f}_{\alpha, N}(z; \theta^t)$. Residual dynamics:

$$\partial_t \|u_t^{\alpha, N}\|_{L^2}^2 = -\langle u_t^{\alpha, N}, \mathcal{H}_{\rho_t^{\alpha, N}} u_t^{\alpha, N} \rangle.$$

Here

$$\begin{aligned} \mathcal{H}_{\rho}(x, z) &\equiv \int \langle \nabla_{\theta} \sigma_{*}(x; \theta), \nabla_{\theta} \sigma_{*}(z; \theta) \rangle \rho(d\theta), \\ \Psi_{\alpha}(\theta; \rho^{\alpha, N}) &= -\mathbb{E}_x [u_t^{\alpha, N}(x) \sigma_{*}(x; \theta)]. \end{aligned}$$

The mean field limit and kernel limit

$$\begin{aligned}\partial_t \rho_t^{\alpha, N} &= (1/\alpha) \nabla_{\theta} \cdot (\rho_t^{\alpha, N} [\nabla_{\theta} \Psi(\theta; \rho_t^{\alpha, N})]), \\ \partial_t \|\mathbf{u}_t^{\alpha, N}\|_{L^2}^2 &= - \langle \mathbf{u}_t^{\alpha, N}, \mathcal{H}_{\rho_t^{\alpha, N}} \mathbf{u}_t^{\alpha, N} \rangle.\end{aligned}$$

- ▶ The mean field limit: fix $\alpha = O(1)$ and let $N \rightarrow \infty$.
- ▶ The kernel limit: let $\alpha \rightarrow \infty$ after $N \rightarrow \infty$.
- ▶ The benefit of kernel limit: the kernel will not change, and the residual dynamics becomes self contained. The empirical risk will converge to 0. Full derivation see appendix H of [Mei, Misiakiewicz, Montanari, 2019].

The mean field limit and kernel limit

$$\begin{aligned}\partial_t \rho_t^{\alpha, N} &= (1/\alpha) \nabla_{\theta} \cdot (\rho_t^{\alpha, N} [\nabla_{\theta} \Psi(\theta; \rho_t^{\alpha, N})]), \\ \partial_t \|\mathbf{u}_t^{\alpha, N}\|_{L^2}^2 &= - \langle \mathbf{u}_t^{\alpha, N}, \mathcal{H}_{\rho_t^{\alpha, N}} \mathbf{u}_t^{\alpha, N} \rangle.\end{aligned}$$

- ▶ The mean field limit: fix $\alpha = O(1)$ and let $N \rightarrow \infty$.
- ▶ The kernel limit: let $\alpha \rightarrow \infty$ after $N \rightarrow \infty$.
- ▶ The benefit of kernel limit: the kernel will not change, and the residual dynamics becomes self contained. The empirical risk will converge to 0. Full derivation see appendix H of [Mei, Misiakiewicz, Montanari, 2019].

The mean field limit and kernel limit

$$\begin{aligned}\partial_t \rho_t^{\alpha, N} &= (1/\alpha) \nabla_{\theta} \cdot (\rho_t^{\alpha, N} [\nabla_{\theta} \Psi(\theta; \rho_t^{\alpha, N})]), \\ \partial_t \|u_t^{\alpha, N}\|_{L^2}^2 &= - \langle u_t^{\alpha, N}, \mathcal{H}_{\rho_t^{\alpha, N}} u_t^{\alpha, N} \rangle.\end{aligned}$$

- ▶ The mean field limit: fix $\alpha = O(1)$ and let $N \rightarrow \infty$.
- ▶ The kernel limit: let $\alpha \rightarrow \infty$ after $N \rightarrow \infty$.
- ▶ The benefit of kernel limit: the kernel will not change, and the residual dynamics becomes self contained. The empirical risk will converge to 0. Full derivation see appendix H of [Mei, Misiakiewicz, Montanari, 2019].

Summary

- ▶ Gave an interpretation of neural tangent kernel in the mean field point of view. (Also in [Chizat and Bach, 2018b])
- ▶ Whether NTK can explain the success of neural network is still an open problem. [Arora, Du, Hu, Li, Salakhutdinov, Wang, 2019], [Ghorbani, Mei, Misiakiwics, and Montanari, 2019a, 2019b], [Allen-Zhu and Li, 2019].

Summary

- ▶ Gave an interpretation of neural tangent kernel in the mean field point of view. (Also in [Chizat and Bach, 2018b])
- ▶ Whether NTK can explain the success of neural network is still an open problem. [Arora, Du, Hu, Li, Salakhutdinov, Wang, 2019], [Ghorbani, Mei, Misiakiwics, and Montanari, 2019a, 2019b], [Allen-Zhu and Li, 2019].