Mean Field Theory and Tangent Kernel Theory of Neural Networks

Song Mei

Stanford University

October 22, 2019

Stanford STATS385 guest lecture

Song Mei (Stanford University)

Mean Field and Tangent Kernel

October 22, 2019 1/55

- 10

Deep learning applications



Applications

- Computer vision (autonomous vehicles).
- ▶ Generative modeling (WaveNet for generating speech).
- Reinforcement learning (Go playing).

Mathematical challenges/mysteries



• Optimization \rightarrow Non-convexity.

• Generalization \rightarrow Overparameterization.

1. Mean field theory

2. Tangent kernel theory

3. Transitions between mean field and tangent kernel regime

- 32

<ロト <回ト < 注ト < 注ト

Agenda

1. Mean field theory

2. Tangent kernel theory

3. Transitions between mean field and tangent kernel regime

- 34

<ロト <回ト < 注ト < 注ト



Figure: $\boldsymbol{\theta}_i = (a_i, \boldsymbol{w}_i)$.

Song Mei (Stanford University)

Mean Field and Tangent Kernei

э

▶ Parameters: $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_N) \in \mathbb{R}^{N \times D}$.

Prediction function:

$$f(\boldsymbol{x}; \boldsymbol{ heta}) = rac{1}{N} \sum_{i=1}^N \sigma_\star(\boldsymbol{x}; \boldsymbol{ heta}_i) = rac{1}{N} \sum_{i=1}^N a_i \sigma(\langle \boldsymbol{w}_i, \boldsymbol{x}
angle).$$

Risk function:

$$R_N(oldsymbol{ heta}) = \mathbb{E}_{oldsymbol{x},y} \Big[\ell\Big(y, rac{1}{N}\sum_{i=1}^N \sigma_\star(oldsymbol{x};oldsymbol{ heta}_i)\Big) \Big].$$

• Gradient flow for R_N :

$$\frac{\mathrm{d}}{\mathrm{d}t}\boldsymbol{\theta}_{i}^{t} = -N\xi(t)\nabla_{\boldsymbol{\theta}_{i}}R_{N}(\boldsymbol{\theta}^{t}). \qquad (\mathrm{GF})$$

(日)、(四)、(日)、(日)、

Difficulty: non-convexity with local minimizers!

Song Mei (Stanford University)

- ▶ Parameters: $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_N) \in \mathbb{R}^{N \times D}$.
- Prediction function:

$$f(\boldsymbol{x};\boldsymbol{\theta}) = rac{1}{N}\sum_{i=1}^{N}\sigma_{\star}(\boldsymbol{x};\boldsymbol{\theta}_{i}) = rac{1}{N}\sum_{i=1}^{N}a_{i}\sigma(\langle \boldsymbol{w}_{i},\boldsymbol{x} \rangle).$$

Risk function:

$$R_N(oldsymbol{ heta}) = \mathbb{E}_{oldsymbol{x}, y} \Big[\ell\Big(y, rac{1}{N} \sum_{i=1}^N \sigma_\star(oldsymbol{x}; oldsymbol{ heta}_i) \Big) \Big].$$

• Gradient flow for R_N :

$$\frac{\mathrm{d}}{\mathrm{d}t}\boldsymbol{\theta}_{i}^{t} = -N\xi(t)\nabla_{\boldsymbol{\theta}_{i}}R_{N}(\boldsymbol{\theta}^{t}). \qquad (\mathrm{GF})$$

イロト イポト イヨト イヨト

Difficulty: non-convexity with local minimizers!

Song Mei (Stanford University)

- ▶ Parameters: $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_N) \in \mathbb{R}^{N \times D}$.
- Prediction function:

$$f(\boldsymbol{x};\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^{N} \sigma_{\star}(\boldsymbol{x};\boldsymbol{\theta}_{i}) = \frac{1}{N} \sum_{i=1}^{N} \boldsymbol{a}_{i} \sigma(\langle \boldsymbol{w}_{i}, \boldsymbol{x} \rangle).$$

Risk function:

$$R_N(oldsymbol{ heta}) = \mathbb{E}_{oldsymbol{x},oldsymbol{y}} \Big[\ell\Big(y, rac{1}{N} \sum_{i=1}^N \sigma_\star(oldsymbol{x};oldsymbol{ heta}_i) \Big) \Big].$$

• Gradient flow for R_N :

$$\frac{\mathrm{d}}{\mathrm{d}t}\boldsymbol{\theta}_{i}^{t} = -N\xi(t)\nabla_{\boldsymbol{\theta}_{i}}R_{N}(\boldsymbol{\theta}^{t}). \qquad (\mathrm{GF})$$

イロト イポト イヨト イヨト

Difficulty: non-convexity with local minimizers!

- ▶ Parameters: $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_N) \in \mathbb{R}^{N \times D}$.
- Prediction function:

$$f(\boldsymbol{x};\boldsymbol{\theta}) = rac{1}{N}\sum_{i=1}^N \sigma_\star(\boldsymbol{x};\boldsymbol{\theta}_i) = rac{1}{N}\sum_{i=1}^N a_i\sigma(\langle \boldsymbol{w}_i, \boldsymbol{x} \rangle).$$

Risk function:

$$R_N(oldsymbol{ heta}) = \mathbb{E}_{oldsymbol{x},oldsymbol{y}} \Big[\ell\Big(y, rac{1}{N}\sum_{i=1}^N \sigma_\star(oldsymbol{x};oldsymbol{ heta}_i)\Big) \Big].$$

• Gradient flow for R_N :

$$\frac{\mathrm{d}}{\mathrm{d}t}\boldsymbol{\theta}_{i}^{t} = -N\xi(t)\nabla_{\boldsymbol{\theta}_{i}}R_{N}(\boldsymbol{\theta}^{t}). \tag{GF}$$

イロト イポト イヨト イヨト

Difficulty: non-convexity with local minimizers!

- ▶ Parameters: $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_N) \in \mathbb{R}^{N \times D}$.
- Prediction function:

$$f(\boldsymbol{x};\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^{N} \sigma_{\star}(\boldsymbol{x};\boldsymbol{\theta}_{i}) = \frac{1}{N} \sum_{i=1}^{N} \boldsymbol{a}_{i} \sigma(\langle \boldsymbol{w}_{i}, \boldsymbol{x} \rangle).$$

Risk function:

$$R_N(oldsymbol{ heta}) = \mathbb{E}_{oldsymbol{x},oldsymbol{y}} \Big[\ell\Big(y, rac{1}{N}\sum_{i=1}^N \sigma_\star(oldsymbol{x};oldsymbol{ heta}_i)\Big) \Big].$$

• Gradient flow for R_N :

$$\frac{\mathrm{d}}{\mathrm{d}t}\boldsymbol{\theta}_{i}^{t} = -N\xi(t)\nabla_{\boldsymbol{\theta}_{i}}R_{N}(\boldsymbol{\theta}^{t}). \tag{GF}$$

<ロト <問ト < 注ト < 注ト

Difficulty: non-convexity with local minimizers!

Song Mei (Stanford University)

3

Landscape analysis?

$$R_N(oldsymbol{ heta}) = \mathbb{E}_{oldsymbol{x},oldsymbol{y}} \Big[\Big(y - rac{1}{N} \sum_{j=1}^N \sigma_\star(oldsymbol{x};oldsymbol{ heta}_j) \Big)^2 \Big].$$

- [Kawaguchi, 2016], [Freeman, Bruna, 2016]: linear network has no spurious local min.
- [Soltanolkotabi, Javanmard, Lee, 2017]: Quadratic two-layers NN has no spurious local min.
- [Zhong, Song, Jain, Bartlett, Dhillon, 2017]: Local strong convexity of two layers NN.
- [Soudry, Carmon, 2016], [Ge, Lee, Ma, 2017], [Tian, 2017], [Soltanolkotabi, 2017], [Li, Yuan, 2017] ...

<ロト <問ト < 注ト < 注ト

Mean field perspective: Emp. dist. of weights

Prediction function

$$f(\boldsymbol{x}; \boldsymbol{\theta}) = rac{1}{N} \sum_{i=1}^{N} \sigma_{\star}(\boldsymbol{x}; \boldsymbol{\theta}_{i}) = \int \sigma_{\star}(\boldsymbol{x}; \boldsymbol{\theta}) \hat{
ho}_{N}(\mathrm{d}\boldsymbol{\theta}).$$

Empirical distribution of the weights:

$$\hat{\rho}_N = rac{1}{N}\sum_{i=1}^N \delta_{oldsymbol{ heta}_i} \in \mathcal{P}(\mathbb{R}^D).$$

▶ Risk functional $R : \mathcal{P}(\mathbb{R}^D) \to \mathbb{R}$

$$R(
ho) = \mathbb{E}_{x,y} \Big[\ell \Big(y, \int \sigma_{\star}(x; \theta)
ho(\mathrm{d}\theta) \Big) \Big].$$

- 32

<ロト <回ト < 注ト < 注ト

Induced dynamics on empirical distribution

• Gradient flow on particles $\{\boldsymbol{\theta}_i^t\}_{i \in [N]}, \, \boldsymbol{\theta}_i^t \in \mathbb{R}^D,$

$$\frac{\mathrm{d}}{\mathrm{d}t}\boldsymbol{\theta}_{i}^{t} = -N\xi(t)\nabla_{\boldsymbol{\theta}_{i}}R_{N}(\boldsymbol{\theta}^{t}), \quad \boldsymbol{\theta}_{i}^{0} = \boldsymbol{\theta}_{i}^{0}. \tag{GF}$$

▶ Define dynamics on distribution $ho_{N,t} \in \mathcal{P}(\mathbb{R}^D),$

$$\partial_t \rho_{N,t}(\boldsymbol{\theta}) = 2\xi(t) \nabla_{\boldsymbol{\theta}} \cdot (\rho_{N,t}(\boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}} \Psi(\boldsymbol{\theta}; \rho_{N,t})),$$

with

$$ho_{N,0} = rac{1}{N}\sum_{i=1}^N \delta_{oldsymbol{ heta}_i^0}, \qquad \Psi(oldsymbol{ heta};
ho) = rac{\delta R}{\delta
ho}(oldsymbol{ heta};
ho).$$

• Claim: $\rho_{N,t} = (1/N) \sum_{i=1}^{N} \delta_{\boldsymbol{\theta}_{i}^{t}}$.

- 34

<ロ> (日) (日) (日) (日) (日)

Induced dynamics on empirical distribution

• Gradient flow on particles $\{\boldsymbol{\theta}_i^t\}_{i \in [N]}, \, \boldsymbol{\theta}_i^t \in \mathbb{R}^D,$

$$\frac{\mathrm{d}}{\mathrm{d}t}\boldsymbol{\theta}_{i}^{t} = -N\xi(t)\nabla_{\boldsymbol{\theta}_{i}}R_{N}(\boldsymbol{\theta}^{t}), \quad \boldsymbol{\theta}_{i}^{0} = \boldsymbol{\theta}_{i}^{0}. \tag{GF}$$

▶ Define dynamics on distribution $\rho_{N,t} \in \mathcal{P}(\mathbb{R}^D)$,

. .

$$\partial_t \rho_{N,t}(\boldsymbol{\theta}) = 2\xi(t) \nabla_{\boldsymbol{\theta}} \cdot (\rho_{N,t}(\boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}} \Psi(\boldsymbol{\theta}; \rho_{N,t})), \qquad (\text{PDE})$$

with

$$ho_{N,0} = rac{1}{N}\sum_{i=1}^N \delta_{oldsymbol{ heta}_i}^{
ho}, \qquad \Psi(oldsymbol{ heta};oldsymbol{
ho}) = rac{\delta R}{\delta
ho}(oldsymbol{ heta};oldsymbol{
ho}).$$

• Claim: $\rho_{N,t} = (1/N) \sum_{i=1}^{N} \delta_{\boldsymbol{\theta}_{i}^{t}}$.

- 34

<ロト <問ト < 注ト < 注ト

Induced dynamics on empirical distribution

• Gradient flow on particles $\{\boldsymbol{\theta}_i^t\}_{i \in [N]}, \, \boldsymbol{\theta}_i^t \in \mathbb{R}^D,$

$$\frac{\mathrm{d}}{\mathrm{d}t}\boldsymbol{\theta}_{i}^{t} = -N\xi(t)\nabla_{\boldsymbol{\theta}_{i}}R_{N}(\boldsymbol{\theta}^{t}), \quad \boldsymbol{\theta}_{i}^{0} = \boldsymbol{\theta}_{i}^{0}. \tag{GF}$$

▶ Define dynamics on distribution $\rho_{N,t} \in \mathcal{P}(\mathbb{R}^D)$,

. .

$$\partial_t \rho_{N,t}(\boldsymbol{\theta}) = 2\xi(t) \nabla_{\boldsymbol{\theta}} \cdot (\rho_{N,t}(\boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}} \Psi(\boldsymbol{\theta}; \rho_{N,t})), \qquad (\text{PDE}$$

with

$$ho_{N,0} = rac{1}{N} \sum_{i=1}^N \delta_{oldsymbol{ heta}_i}^0, \qquad \Psi(oldsymbol{ heta}; oldsymbol{
ho}) = rac{\delta R}{\delta
ho}(oldsymbol{ heta}; oldsymbol{
ho}).$$

• Claim: $\rho_{N,t} = (1/N) \sum_{i=1}^{N} \delta_{\boldsymbol{\theta}_{i}^{t}}$.

- 32

<ロト <回ト < 注ト < 注ト



A short proof

Test function + Chain rule + Integration by part.

$$\begin{split} \frac{\mathrm{d}}{\mathrm{d}t} \int f(\boldsymbol{\theta}) \rho_{N,t}(\mathrm{d}\boldsymbol{\theta}) &= \frac{\mathrm{d}}{\mathrm{d}t} \Big[\frac{1}{N} \sum_{i=1}^{N} f(\boldsymbol{\theta}_{i}^{t}) \Big] = -\frac{1}{N} \sum_{i=1}^{N} \langle \nabla f(\boldsymbol{\theta}_{i}^{t}), N \nabla_{\boldsymbol{\theta}_{i}} R_{N}(\boldsymbol{\theta}^{t}) \rangle \\ &= -\frac{1}{N} \sum_{i=1}^{N} \left\langle \nabla f(\boldsymbol{\theta}_{i}^{t}), \nabla [\delta R/\delta \rho](\boldsymbol{\theta}_{i}^{t}; \rho_{N,t}) \right\rangle \\ &= -\int \left\langle \nabla f, \nabla [\delta R/\delta \rho](\boldsymbol{\theta}; \rho_{N,t}) \right\rangle \rho_{N,t}(\mathrm{d}\boldsymbol{\theta}). \end{split}$$

- 2

・ロト ・四ト ・ヨト ・ヨト

What is this PDE?

$$\partial_t \rho_t = \nabla_{\boldsymbol{\theta}} \cdot \left(\rho_t \nabla_{\boldsymbol{\theta}} \Psi(\boldsymbol{\theta}; \rho_t) \right).$$

Existence and uniqueness: [Sznitman, 1991].

- Physics: nonlinear transport equation describing motions of particles with pairwise interaction (mean field approach).
- Math: Gradient flow of R(ρ) in the metric space (P(R^D), W₂). [Jordan, Kinderlehrer, Otto, 1998], [Ambrosio, Gigli, Savaré, 2006], [Carrillo, McCann, Villani, 2013]

What is this PDE?

$$\partial_t \rho_t = \nabla_{\boldsymbol{\theta}} \cdot \left(\rho_t \nabla_{\boldsymbol{\theta}} \Psi(\boldsymbol{\theta}; \rho_t) \right).$$

Existence and uniqueness: [Sznitman, 1991].

- Physics: nonlinear transport equation describing motions of particles with pairwise interaction (mean field approach).
- Math: Gradient flow of R(ρ) in the metric space (P(R^D), W₂). [Jordan, Kinderlehrer, Otto, 1998], [Ambrosio, Gigli, Savaré, 2006], [Carrillo, McCann, Villani, 2013]

- 20

Converge of $ho_{N,t}$ to ho_t as $N o \infty$

Let $\rho_{N,t}$ be the solution of $(\boldsymbol{\theta}_i^0 \sim_{iid} \rho_0 \in \mathcal{P}(\mathbb{R}^D))$

$$\partial_t \rho_{N,t} = \nabla_{\boldsymbol{\theta}} \cdot \Big(\rho_{N,t} \nabla_{\boldsymbol{\theta}} \Psi(\boldsymbol{\theta}; \rho_{N,t}) \Big), \quad \rho_{N,0} = \frac{1}{N} \sum_{i=1}^N \delta_{\boldsymbol{\theta}_i^0}.$$

Let ρ_t be the solution of

$$\partial_t
ho_t =
abla_{oldsymbol{ heta}} \cdot \Big(
ho_t
abla_{oldsymbol{ heta}} \Psi(oldsymbol{ heta};
ho_t) \Big), \quad
ho_t|_{t=0} =
ho_0.$$

- 34

ъ*т*

Converge of $ho_{N,t}$ to ho_t as $N o \infty$

Let $\rho_{N,t}$ be the solution of $(\boldsymbol{\theta}_i^0 \sim_{iid} \rho_0 \in \mathcal{P}(\mathbb{R}^D))$

$$\partial_t \rho_{N,t} =
abla_{m{ heta}} \cdot \Big(
ho_{N,t}
abla_{m{ heta}} \Psi(m{ heta};
ho_{N,t}) \Big), \quad
ho_{N,0} = rac{1}{N} \sum_{i=1}^N \delta_{m{ heta}_i^0}.$$

Let ρ_t be the solution of

$$\partial_t \rho_t = \nabla_{\boldsymbol{\theta}} \cdot \left(\rho_t \nabla_{\boldsymbol{\theta}} \Psi(\boldsymbol{\theta}; \rho_t) \right), \quad \rho_t|_{t=0} = \rho_0.$$

3

3.7

Converge of $ho_{N,t}$ to ho_t as $N o \infty$

Theorem (M., Montanari, and Nguyen, 2018)

Under some assumptions. Let $(\theta_i^0)_{i \leq N} \sim_{iid} \rho_0$. Denote $\rho_{N,t}$: emp. dist. of sol. of grad. flow with init. θ^0 . Denote ρ_t : sol. of PDE with init. ρ_0 . Then, $\forall f$ bounded Lipschitz,

$$\sup_{t \leq T} \left| \int f(\boldsymbol{\theta}) \rho_{N,t}(\mathrm{d}\boldsymbol{\theta}) - \int f(\boldsymbol{\theta}) \rho_t(\mathrm{d}\boldsymbol{\theta}) \right| \leq K e^{KT} \sqrt{\frac{\log(N/\delta)}{N}}$$

with probability at least $1 - \delta$.

See also [Rotskoff, Vanden-Eijnden, 2018], [Sirignano, Spiliopoulos, 2018].

・ロト ・ 一下・ ・ ヨト ・ ヨー・

Immediate implication

$$egin{aligned} &
ho_{N,t} o
ho_t, \quad N o \infty \ &\partial
ho_t =
abla \cdot ig(
ho_t
abla \Psi(oldsymbol{ heta};
ho_t) ig), \quad eta_t ert_{t=0} =
ho_0. \end{aligned}$$

Convergence speed of N-neuron gradient flow is independent of N!

• Effective dimension from $N \times D$ to D.

- 31

SGD v.s. PDE



æ

(a)

SGD v.s. PDE



Song Mei (Stanford University)

- 2

<ロト <回ト < 注ト < 注ト

Does $R(\rho_t) \to \min_{\rho} R(\rho)$ as $t \to \infty$?

In general, no convergence guarantees.

But sometimes, yes.

[M., Montanari, Nguyen, 2018]

Case by case: a special mixture of two Gaussians.

▶ Noisy SGD (PDE with diffusion term).

◆□ ▶ ◆□ ▶ ▲ ■ ▶ ▲ ■ ▶ ● ● ●

PDE with diffusion term

Noisy gradient flow

$$\mathrm{d} oldsymbol{ heta}_i^t = -N
abla_{oldsymbol{ heta}_i} R_N(oldsymbol{ heta}^t) \mathrm{d} t + rac{1}{\sqrt{eta}} \mathrm{d} oldsymbol{W}_i^t.$$

PDE with diffusion term

$$\partial_t \rho_t = \nabla_{\boldsymbol{\theta}} \cdot \left(\rho_t \nabla_{\boldsymbol{\theta}} \Psi(\boldsymbol{\theta}; \rho_t) \right) + \frac{1}{\beta} \Delta \rho_t. \tag{*}$$

Wasserstein grad. flow of the free energy

$$F_{oldsymbol{eta}}(
ho) = R(
ho) + rac{1}{oldsymbol{eta}}\int
ho(oldsymbol{ heta})\log
ho(oldsymbol{ heta})\mathrm{d}oldsymbol{ heta}.$$

Theorem (M., Montanari, Nguyen, 2018)

(*) converges to the minimizer of $F_{eta}(
ho)$ as $t \to \infty$.

Song Mei (Stanford University)

1

(日)

General convergence for noisy SGD

Theorem (M., Montanari, Nguyen, 2018)

Under certain assumptions. Initialization $(\theta_i^0)_{i \leq N} \sim_{iid} \rho_0$. Then there exists $\beta_0 = \beta_0(D, U, V, \eta)$, such that, for $\beta \geq \beta_0$, there exists $T = T(D, U, V, \beta, \eta)$ such that for any $k \in [T/\varepsilon, 10T/\varepsilon]$, $N \geq C_0 D \log D$, $\varepsilon \leq 1/(C_0 D)$, we have, w.h.p.

$$R_{\lambda,N}(\boldsymbol{\theta}^{k}) \leq \inf_{\boldsymbol{\theta} \in \mathbb{R}^{D \times N}} R_{\lambda,N}(\boldsymbol{\theta}) + \eta.$$

Cautious: no polynomial convergence rate.

< ロ > < 同 > < 回 > < 回 > < 回 > <

General convergence for noisy SGD

Theorem (M., Montanari, Nguyen, 2018)

Under certain assumptions. Initialization $(\theta_i^0)_{i \leq N} \sim_{iid} \rho_0$. Then there exists $\beta_0 = \beta_0(D, U, V, \eta)$, such that, for $\beta \geq \beta_0$, there exists $T = T(D, U, V, \beta, \eta)$ such that for any $k \in [T/\varepsilon, 10T/\varepsilon]$, $N \geq C_0 D \log D$, $\varepsilon \leq 1/(C_0 D)$, we have, w.h.p.

$$R_{\lambda,N}(\boldsymbol{\theta}^k) \leq \inf_{\boldsymbol{\theta} \in \mathbb{R}^{D \times N}} R_{\lambda,N}(\boldsymbol{\theta}) + \eta.$$

Cautious: no polynomial convergence rate.

Other references

Other references

- ▶ [Chizat, Bach, 2018a]: Global convergence with homogenuity.
- ▶ [Ma, Lee, Liu, Wei, 2019]: Quantitative rate for modified PDE.
- [Rotskoff, Jelassi, Bruna, Vanden-Eijnden, 2019]: Semi-quantitative rate for the birth-death process.
- [Nguyen, 2019], [Sirignano, Spiliopoulos, 2019], [Araújo, Oliveira, Yukimura, 2019]: Extension to multilayers.

Still many open problems: establish global convergence, extend to multi-layers.



1. Mean field theory

2. Tangent kernel theory

3. Transitions between mean field and tangent kernel regime

Song Mei (Stanford University)

Mean Field and Tangent Kernel

- 34

<ロ> (日) (日) (日) (日) (日)

• Multi-layers neural network $f(x; \theta)$

$$f(\boldsymbol{x}; \boldsymbol{\theta}) = \sigma(\boldsymbol{W}_L \sigma(\cdots \boldsymbol{W}_2 \sigma(\boldsymbol{W}_1 \boldsymbol{x}))).$$

 \blacktriangleright Linearization around random parameter θ_0

$$f(\boldsymbol{x};\boldsymbol{\theta}) = f(\boldsymbol{x};\boldsymbol{\theta}_0) + \underline{\langle \boldsymbol{\theta} - \boldsymbol{\theta}_0, \nabla_{\boldsymbol{\theta}} f(\boldsymbol{x};\boldsymbol{\theta}_0) \rangle} + o(\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|_2).$$

 \blacktriangleright NT model: the linear part of f

$$f_{\mathsf{NT}}(\boldsymbol{x};\boldsymbol{ heta}) = \langle \boldsymbol{ heta} - \boldsymbol{ heta}_0,
abla_{\boldsymbol{ heta}} f(\boldsymbol{x};\boldsymbol{ heta}_0)
angle.$$

1

• Multi-layers neural network $f(x; \theta)$

$$f(\boldsymbol{x}; \boldsymbol{\theta}) = \sigma(\boldsymbol{W}_L \sigma(\cdots \boldsymbol{W}_2 \sigma(\boldsymbol{W}_1 \boldsymbol{x}))).$$

• Linearization around random parameter θ_0

$$f(\boldsymbol{x};\boldsymbol{\theta}) = f(\boldsymbol{x};\boldsymbol{\theta}_0) + \underline{\langle \boldsymbol{\theta} - \boldsymbol{\theta}_0, \nabla_{\boldsymbol{\theta}} f(\boldsymbol{x};\boldsymbol{\theta}_0) \rangle} + o(\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|_2).$$

 \blacktriangleright NT model: the linear part of f

$$f_{\mathsf{NT}}(\boldsymbol{x};\boldsymbol{ heta}) = \langle \boldsymbol{ heta} - \boldsymbol{ heta}_0,
abla_{\boldsymbol{ heta}} f(\boldsymbol{x};\boldsymbol{ heta}_0)
angle.$$

-

イロト イポト イヨト イヨ

• Multi-layers neural network $f(x; \theta)$

$$f(\boldsymbol{x}; \boldsymbol{\theta}) = \sigma(\boldsymbol{W}_L \sigma(\cdots \boldsymbol{W}_2 \sigma(\boldsymbol{W}_1 \boldsymbol{x}))).$$

• Linearization around random parameter θ_0

$$f(\boldsymbol{x};\boldsymbol{\theta}) = f(\boldsymbol{x};\boldsymbol{\theta}_0) + \underline{\langle \boldsymbol{\theta} - \boldsymbol{\theta}_0, \nabla_{\boldsymbol{\theta}} f(\boldsymbol{x};\boldsymbol{\theta}_0) \rangle} + o(\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|_2).$$

 \blacktriangleright NT model: the linear part of f

$$f_{\mathsf{NT}}(\boldsymbol{x};\boldsymbol{\theta}) = \langle \boldsymbol{\theta} - \boldsymbol{\theta}_0, \nabla_{\boldsymbol{\theta}} f(\boldsymbol{x};\boldsymbol{\theta}_0) \rangle.$$

3

-

NT model: the linear part of f

$$f_{\mathsf{NT}}(\boldsymbol{x};\boldsymbol{\theta}) = \langle \boldsymbol{\theta} - \boldsymbol{\theta}_0, \nabla_{\boldsymbol{\theta}} f(\boldsymbol{x};\boldsymbol{\theta}_0) \rangle.$$

- ▶ Random feature map: $\nabla_{\theta} f(\boldsymbol{x}; \boldsymbol{\theta}_0)$.
- Neural tangent kernel: $\mathcal{K}_{NT}(x, y) = \langle \nabla_{\theta} f(x; \theta_0), \nabla_{\theta} f(y; \theta_0) \rangle$. [Jacot, Gabriel, Hongler, 2018], [Du, Zhai, Poczos, Singh, 2018], [Chizat, Bach, 2018b],
- Successful optimization: under certain conditions (different from the mean field theory), the trajectory (of GF on empirical risk) of NT and NN is uniformly close.

Does NT models fully explain the success of neural networks?

イロト イポト イヨト イヨ
The neural tangent model

▶ NT model: the linear part of f

$$f_{\mathsf{NT}}(\boldsymbol{x};\boldsymbol{\theta}) = \langle \boldsymbol{\theta} - \boldsymbol{\theta}_0, \nabla_{\boldsymbol{\theta}} f(\boldsymbol{x};\boldsymbol{\theta}_0) \rangle.$$

▶ Random feature map: $\nabla_{\theta} f(x; \theta_0)$.

▶ Neural tangent kernel: $\mathcal{K}_{NT}(x, y) = \langle \nabla_{\theta} f(x; \theta_0), \nabla_{\theta} f(y; \theta_0) \rangle$. [Jacot, Gabriel, Hongler, 2018], [Du, Zhai, Poczos, Singh, 2018], [Chizat, Bach, 2018b],

Successful optimization: under certain conditions (different from the mean field theory), the trajectory (of GF on empirical risk) of NT and NN is uniformly close.

Does NT models fully explain the success of neural networks?

イロト イボト イヨト イヨト

The neural tangent model

▶ NT model: the linear part of f

$$f_{\mathsf{NT}}(\boldsymbol{x};\boldsymbol{\theta}) = \langle \boldsymbol{\theta} - \boldsymbol{\theta}_0, \nabla_{\boldsymbol{\theta}} f(\boldsymbol{x};\boldsymbol{\theta}_0) \rangle.$$

▶ Random feature map: $\nabla_{\theta} f(x; \theta_0)$.

- ▶ Neural tangent kernel: $\mathcal{K}_{NT}(x, y) = \langle \nabla_{\theta} f(x; \theta_0), \nabla_{\theta} f(y; \theta_0) \rangle$. [Jacot, Gabriel, Hongler, 2018], [Du, Zhai, Poczos, Singh, 2018], [Chizat, Bach, 2018b],
- Successful optimization: under certain conditions (different from the mean field theory), the trajectory (of GF on empirical risk) of NT and NN is uniformly close.

Does NT models fully explain the success of neural networks?

イロト イポト イヨト イヨト

The neural tangent model

▶ NT model: the linear part of f

$$f_{\mathsf{NT}}(\boldsymbol{x};\boldsymbol{\theta}) = \langle \boldsymbol{\theta} - \boldsymbol{\theta}_0, \nabla_{\boldsymbol{\theta}} f(\boldsymbol{x};\boldsymbol{\theta}_0) \rangle.$$

▶ Random feature map: $\nabla_{\theta} f(x; \theta_0)$.

- ▶ Neural tangent kernel: $\mathcal{K}_{NT}(x, y) = \langle \nabla_{\theta} f(x; \theta_0), \nabla_{\theta} f(y; \theta_0) \rangle$. [Jacot, Gabriel, Hongler, 2018], [Du, Zhai, Poczos, Singh, 2018], [Chizat, Bach, 2018b],
- Successful optimization: under certain conditions (different from the mean field theory), the trajectory (of GF on empirical risk) of NT and NN is uniformly close.

Does NT models fully explain the success of neural networks?

<ロ> (四) (四) (三) (三) (三)

Empirically, the generalization of NT models are not as good as $\ensuremath{\operatorname{NN}}$

Table: Cifar10 experiments

Architecture	Classification error
Best convolutional NN	5%-
Best convolutional NT	23%
CNN of best CNT	19%

[Arora, Du, Hu, Li, Salakhutdinov, Wang, 2019]

-

・ロト ・ 一下・ ・ ヨト・

Theoretical analysis of generalization gap

Two-layers neural network

$$f_N({m x};{m \Theta}) = \sum_{i=1}^N {m a}_i \sigma(\langle {m w}_i, {m x}
angle), \quad {m \Theta} = (a_1, {m w}_1, \dots, a_N, {m w}_N).$$

- ▶ Input vector $x \in \mathbb{R}^d$.
- ▶ Bottom layer weights $w_i \in \mathbb{R}^d$, i = 1, 2, ..., N.
- Top layer weights $\mathbf{a}_i \in \mathbb{R}, i = 1, 2, \dots, N$.

E ∕Q (~

イロト イポト イヨト イヨト

Random features model and Neural tangent model

Linearization

$$f_N(x; \Theta) = f_N(x; \Theta^0) + \underbrace{\sum_{i=1}^N \Delta a_i \sigma(\langle w_i^0, x \rangle)}_{\text{Second layer linearization}} + \underbrace{\sum_{i=1}^N a_i^0 \sigma'(\langle w_i^0, x \rangle) \langle \Delta w_i, x \rangle}_{\text{First layer linearization}} + o(\cdot).$$

Linearized neural network: $(w_i \sim \text{Unif}(\mathbb{S}^{d-1}))$

$$\mathcal{F}_{\mathsf{RF},N}(\boldsymbol{W}) = \Big\{ f = \sum_{i=1}^{N} \boldsymbol{a}_{i} \sigma(\langle \boldsymbol{w}_{i}, \boldsymbol{x} \rangle) : \boldsymbol{a}_{i} \in \mathbb{R}, i \in [N] \Big\},$$
$$\mathcal{F}_{\mathsf{NT},N}(\boldsymbol{W}) = \Big\{ f = \sum_{i=1}^{N} \sigma'(\langle \boldsymbol{w}_{i}, \boldsymbol{x} \rangle) \langle \boldsymbol{a}_{i}, \boldsymbol{x} \rangle : \boldsymbol{a}_{i} \in \mathbb{R}^{d}, i \in [N] \Big\}.$$

Blue: random and fixed. Red: parameters to be optimized.

Approximation error

Data distribution:

$$oldsymbol{x} \sim \mathrm{Unif}(\mathbb{S}^{d-1}(\sqrt{d})), \quad f_{\star} \in L^2(\mathbb{S}^{d-1}(\sqrt{d})).$$

Minimum risk (approximation error):

$$R_{\mathsf{M},N}(f_{\star}) = \inf_{f \in \mathcal{F}_{\mathsf{M},N}(oldsymbol{W})} \mathbb{E}_{oldsymbol{x}}\Big[\Big(f_{\star}(oldsymbol{x}) - f(oldsymbol{x})\Big)^2\Big], \hspace{0.2cm} \mathsf{M} \in \{\mathsf{RF},\mathsf{NT}\}.$$

▲ロト ▲御 ▶ ▲ 臣 ▶ ▲ 臣 ▶ ○ 臣 ○ の Q @

Staircase lower bound

1

・ロト ・個ト ・モト ・モト

Lower bound for random features regression

$$\mathcal{F}_{\mathsf{RF}, oldsymbol{N}}(oldsymbol{W}) = \Big\{ f = \sum_{i=1}^{oldsymbol{N}} a_i \sigma(\langle oldsymbol{w}_i, oldsymbol{x}
angle) : a_i \in \mathbb{R}, i \in [oldsymbol{N}] \Big\}.$$

Theorem (Ghorbani, M., Misiakiwics, Montanari, 2019) Assume $N = O_d(d^{\ell+1-\delta})$, and $(w_i)_{i \in [N]} \sim \text{Unif}(\mathbb{S}^{d-1})$, we have $\inf_{f \in \mathcal{F}_{\text{RF},N}(W)} \mathbb{E}_x[(f_\star(x) - f(x))^2] \ge \|\mathsf{P}_{>\ell}f_\star\|_{L^2}^2 + o_{d,\mathbb{P}}(\|f_\star\|_2^2).$

 $\mathsf{P}_{>\ell}$: projection orthogonal to the space of degree- ℓ polynomials.

With ${\it N}={\it O}_d(d^k)$ parameters, one can only fit a degree k polynomial.

Lower bound for random features regression

$$\mathcal{F}_{\mathsf{RF}, oldsymbol{N}}(oldsymbol{W}) = \Big\{ f = \sum_{i=1}^{oldsymbol{N}} a_i \sigma(\langle oldsymbol{w}_i, oldsymbol{x}
angle) : a_i \in \mathbb{R}, i \in [oldsymbol{N}] \Big\}.$$

Theorem (Ghorbani, M., Misiakiwics, Montanari, 2019) Assume $N = O_d(d^{\ell+1-\delta})$, and $(w_i)_{i \in [N]} \sim \text{Unif}(\mathbb{S}^{d-1})$, we have $\inf_{f \in \mathcal{F}_{\mathsf{RF},N}(W)} \mathbb{E}_x[(f_\star(x) - f(x))^2] \ge \|\mathsf{P}_{>\ell}f_\star\|_{L^2}^2 + o_{d,\mathbb{P}}(\|f_\star\|_2^2).$

 $\mathsf{P}_{>\ell} {:}\ projection\ orthogonal\ to\ the\ space\ of\ degree-\ell\ polynomials.$

With $N = O_d(d^k)$ parameters, one can only fit a degree k polynomial.

<ロ> (四) (四) (三) (三) (三) (三)

Similar result for NT

$$\mathcal{F}_{\mathsf{NT},\mathbf{N}}(\mathbf{W}) = \left\{ f = \sum_{i=1}^{\mathbf{N}} \sigma'(\langle \mathbf{w}_i, \mathbf{x} \rangle) \langle \mathbf{a}_i, \mathbf{x} \rangle : \mathbf{a}_i \in \mathbb{R}^d, i \in [\mathbf{N}]
ight\}.$$

Theorem (Ghorbani, M., Misiakiwics, Montanari, 2019) Assume $N = O_d(d^{\ell+1-\delta})$, and $(w_i)_{i \in [N]} \sim \text{Unif}(\mathbb{S}^{d-1})$, we have $\inf_{f \in \mathcal{F}_{\mathsf{NT},N}(W)} \mathbb{E}_{\boldsymbol{x}}[(f_{\star}(\boldsymbol{x}) - f(\boldsymbol{x}))^2] \geq ||\mathsf{P}_{>\ell+1}f_{\star}||_{L^2}^2 + o_{d,\mathbb{P}}(||f_{\star}||_2^2),$

 $\mathsf{P}_{>\ell+1} {:}\ \text{projection orthogonal to the space of degree-} (\ell+1) \ \text{polynomials.}$

With $O_d(d^{k+1})$ parameters, one can only fit a degree k+1 polynomial.

・ロト ・ 一下・ ・ ヨト ・ ヨー・

Similar result for NT

$$\mathcal{F}_{\mathsf{NT},\mathbf{N}}(\mathbf{W}) = \left\{ f = \sum_{i=1}^{\mathbf{N}} \sigma'(\langle \mathbf{w}_i, \mathbf{x} \rangle) \langle \mathbf{a}_i, \mathbf{x} \rangle : \mathbf{a}_i \in \mathbb{R}^d, i \in [\mathbf{N}]
ight\}.$$

Theorem (Ghorbani, M., Misiakiwics, Montanari, 2019) Assume $N = O_d(d^{\ell+1-\delta})$, and $(w_i)_{i \in [N]} \sim \text{Unif}(\mathbb{S}^{d-1})$, we have $\inf_{f \in \mathcal{F}_{\mathsf{NT},N}(\mathbf{W})} \mathbb{E}_{\mathbf{x}}[(f_{\star}(\mathbf{x}) - f(\mathbf{x}))^2] \ge \|\mathsf{P}_{>\ell+1}f_{\star}\|_{L^2}^2 + o_{d,\mathbb{P}}(\|f_{\star}\|_2^2),$

 $P_{>\ell+1}$: projection orthogonal to the space of degree- $(\ell+1)$ polynomials.

With $O_d(d^{k+1})$ parameters, one can only fit a degree k+1 polynomial.

The staircase lower bound

$$f = \mathsf{P}_0 f + \mathsf{P}_1 f + \mathsf{P}_2 f + \mathsf{P}_3 f + \cdots$$



Song Mei (Stanford University)

October 22, 2019 31/55

3

<ロト <回ト < 回ト < 三

Implication

Function
$$f: \mathbb{S}^{d-1} \to \mathbb{R}$$
, $f(x) = P_k(x_1)$.

- ► NT: $N \ge \Theta_d(d^{k-1});$
- ▶ NN: $N = O_d(1)$.
- ▶ Different from the RKHS theory [Bach, 2017], [E, Ma, Wu, 2018].

Difference 1:

$$f_{\star} \in L^2$$
, v.s. $f_{\star} \in \text{RKHS}$.

Difference 2:

 $N = d^k$ as $d \to \infty$, v.s. fixed d as $N \to \infty$.

Implication

Function
$$f: \mathbb{S}^{d-1} \to \mathbb{R}, f(x) = P_k(x_1)$$
.

- ► NT: $N \ge \Theta_d(d^{k-1});$
- ▶ NN: $N = O_d(1)$.
- ▶ Different from the RKHS theory [Bach, 2017], [E, Ma, Wu, 2018].

▶ Difference 1:

$$f_{\star} \in L^2$$
, v.s. $f_{\star} \in \mathsf{RKHS}$.

Difference 2:

 $N = d^k$ as $d \to \infty$, v.s. fixed d as $N \to \infty$.

double descent curve

1

・ロト ・個ト ・モト ・モト

▶ MNIST dataset: $(x, y) \in \mathbb{R}^{100} \times [10]$. Training/test data size: 50000/10000.

Two-layers neural networks:

$$f_N(x; \Theta) = \sum_{i=1}^N a_i \sigma(\langle \boldsymbol{w}_i, x \rangle), \quad \Theta = (\boldsymbol{a}_1, \boldsymbol{w}_1, \dots, \boldsymbol{a}_N, \boldsymbol{w}_N).$$

Bottom layer weights $oldsymbol{w}_i \in \mathbb{R}^{100}$. Top layer weights $oldsymbol{a}_i \in \mathbb{R}.$

- Loss function: cross entropy. Training algorithm: SGD.
- ▶ N: model complexity, to be varied.

- 1

イロト イボト イヨト イヨト

- ▶ MNIST dataset: $(x, y) \in \mathbb{R}^{100} \times [10]$. Training/test data size: 50000/10000.
- ▶ Two-layers neural networks:

$$f_N(x; \Theta) = \sum_{i=1}^N a_i \sigma(\langle w_i, x \rangle), \quad \Theta = (a_1, w_1, \dots, a_N, w_N).$$

Bottom layer weights $w_i \in \mathbb{R}^{100}$. Top layer weights $a_i \in \mathbb{R}$.

▶ Loss function: cross entropy. Training algorithm: SGD.

- ▶ MNIST dataset: $(x, y) \in \mathbb{R}^{100} \times [10]$. Training/test data size: 50000/10000.
- Two-layers neural networks:

$$f_N(x; \Theta) = \sum_{i=1}^N a_i \sigma(\langle w_i, x \rangle), \quad \Theta = (a_1, w_1, \dots, a_N, w_N).$$

Bottom layer weights $w_i \in \mathbb{R}^{100}$. Top layer weights $a_i \in \mathbb{R}$.

▶ Loss function: cross entropy. Training algorithm: SGD.

▶ N: model complexity, to be varied.

- ▶ MNIST dataset: $(x, y) \in \mathbb{R}^{100} \times [10]$. Training/test data size: 50000/10000.
- Two-layers neural networks:

$$f_N(\boldsymbol{x};\boldsymbol{\Theta}) = \sum_{i=1}^N \boldsymbol{a}_i \sigma(\langle \boldsymbol{w}_i, \boldsymbol{x} \rangle), \quad \boldsymbol{\Theta} = (\boldsymbol{a}_1, \boldsymbol{w}_1, \dots, \boldsymbol{a}_N, \boldsymbol{w}_N).$$

Bottom layer weights $w_i \in \mathbb{R}^{100}$. Top layer weights $a_i \in \mathbb{R}$.

- ▶ Loss function: cross entropy. Training algorithm: SGD.
- ▶ N: model complexity, to be varied.

Experimental results



Figure: Experiments on MNIST. Left: [Spigler, Geiger, Ascoli, Sagun, Biroli, Wyart, 2018]. Right: [Belkin, Hsu, Ma, Mandal, 2018]. See also: [Neyshabur, Tomioka, Srebro, 2014a].

イロト イボト イヨト イヨト

Double descent



Figure: A cartoon by [Belkin, Hsu, Ma, Mandal, 2018].

- \checkmark Peak at the interpolation threshold.
- \checkmark Monotonic decreasing in the overparameterized regime.
- \checkmark Global minimum when the number of parameters is infinity.

< /₽ > <

Linear model with random covariates



By [Hastie, Montanari, Rosset, Tibshirani, 2019]. See also [Belkin, Hsu, Xu, 2019].

$$\begin{array}{ll} \text{Model: } y = \langle x, \beta_\star \rangle + \varepsilon, \ x \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d).\\ \text{Loss: } L(\beta) = \hat{\mathbb{E}}[(y - \langle x, \beta \rangle)^2] \end{array}$$

Song Mei (Stanford University)

October 22, 2019 37 / 55

э

(a)

Why singularity?

► Model:
$$x_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d), \ y_i = \langle x_i, \beta_\star \rangle + \varepsilon_i, \ \beta_\star = \mathbf{0}, \ i \in [n].$$

- $\blacktriangleright \text{ Test risk } \propto \mathbb{E}[\|\hat{\boldsymbol{\beta}}\|_2^2] \propto \mathbb{E}[\|\boldsymbol{X}^{\dagger}\boldsymbol{y}\|_2^2] \propto \mathbb{E}[\operatorname{tr}((\boldsymbol{X}^{\mathsf{T}}\boldsymbol{X})^{\dagger})].$
- When $n \neq d$, X is well conditioned.
- When $n \approx d$, X is infinitely ill conditioned.

• The model has marginally enough parameters to interpolate all the data, hence it interpolates in an awkward way.

• To fit the noise, the coefficients $\|\hat{\boldsymbol{\beta}}\|_2^2 = \|\boldsymbol{X}^{\dagger}\boldsymbol{y}\|_2^2$ blows up.

- 1

イロト イボト イヨト イヨト

Why singularity?

► Model:
$$x_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d), \ y_i = \langle x_i, \beta_\star \rangle + \varepsilon_i, \ \beta_\star = \mathbf{0}, \ i \in [n].$$

- $\blacktriangleright \text{ Test risk } \propto \mathbb{E}[\|\hat{\boldsymbol{\beta}}\|_2^2] \propto \mathbb{E}[\|\boldsymbol{X}^{\dagger}\boldsymbol{y}\|_2^2] \propto \mathbb{E}[\operatorname{tr}((\boldsymbol{X}^{\mathsf{T}}\boldsymbol{X})^{\dagger})].$
- When $n \neq d$, X is well conditioned.
- When $n \approx d$, X is infinitely ill conditioned.

The model has marginally enough parameters to interpolate all the data, hence it interpolates in an awkward way.

• To fit the noise, the coefficients $\|\hat{\beta}\|_2^2 = \|X^{\dagger}y\|_2^2$ blows up.

∃ \000

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Comparison



- \checkmark Peak at the interpolation threshold.
- ? Monotonic decreasing in the overparameterized regime.
- ? Global minimum when the number of parameters is infinity.

・ロト ・ 一下・ ・ 日 ト

Goal: find a tractable model that exhibits all the features of the double descent curve.



Figure: By [Belkin, Hsu, Ma, Mandal, 2018].

Image: A math a math

A simple model

The random features model

$$f_{\mathsf{RF}}({m{x}};{m{a}}) = \sum_{j=1}^N {m{a}_j \sigma(\langle {m{w}}_j, {m{x}}
angle)}.$$

Random weights

 $w_j \sim_{iid} \operatorname{Unif}(\mathbb{S}^{d-1}).$

3

・ロト ・四ト ・ヨト ・ヨト

▶ Random features regression: $\hat{a}_{\lambda} = \arg \min_{a} L_{\lambda}(a)$,

$$L_{\lambda}(\boldsymbol{a}) = \frac{1}{n} \sum_{i=1}^{n} \left[\left(y_i - \sum_{j=1}^{N} a_j \sigma(\langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle) \right)^2 \right] + \frac{\lambda N}{d} \|\boldsymbol{a}\|_2^2,$$
$$R(\boldsymbol{a}; f_{\star}) = \mathbb{E}_{\boldsymbol{x}, y} \left[\left(f_{\star}(\boldsymbol{x}) - \sum_{j=1}^{N} a_j \sigma(\langle \boldsymbol{x}, \boldsymbol{w}_j \rangle) \right)^2 \right].$$

► Data: $(x_i)_{i \in [n]} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d})), \ y_i = f_\star(x_i) + \varepsilon_i, \ \mathbb{E}[\varepsilon_i^2] = \tau^2.$

- Weights: $(\boldsymbol{w}_j)_{j \in [N]} \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}).$
- Activation σ : $\|\mathsf{P}_1\sigma\|_{L^2}^2 = \mu_1^2$, and $\|\mathsf{P}_{>1}\sigma\|_{L^2}^2 = \mu_{\star}^2$.
- ▶ Tech. ass. on f_{\star} and σ . Almost every nonlinear f_{\star} and σ .
- ▶ *n* data, *N* features, *d* dimension. $N/d \rightarrow \psi_1$, $n/d \rightarrow \psi_2$, as $d \rightarrow \infty$.

・ロト ・得ト ・ヨト ・ヨト - ヨ

▶ Random features regression: $\hat{a}_{\lambda} = \arg \min_{a} L_{\lambda}(a)$,

$$L_{\lambda}(\boldsymbol{a}) = \frac{1}{n} \sum_{i=1}^{n} \left[\left(y_i - \sum_{j=1}^{N} a_j \sigma(\langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle) \right)^2 \right] + \frac{\lambda N}{d} \|\boldsymbol{a}\|_2^2,$$
$$R(\boldsymbol{a}; f_{\star}) = \mathbb{E}_{\boldsymbol{x}, y} \left[\left(f_{\star}(\boldsymbol{x}) - \sum_{j=1}^{N} a_j \sigma(\langle \boldsymbol{x}, \boldsymbol{w}_j \rangle) \right)^2 \right].$$

► Data: $(\boldsymbol{x}_i)_{i \in [n]} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d})), \ y_i = f_\star(\boldsymbol{x}_i) + \varepsilon_i, \ \mathbb{E}[\varepsilon_i^2] = \tau^2.$

• Weights: $(\boldsymbol{w}_j)_{j \in [N]} \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}).$

• Activation σ : $\|\mathsf{P}_1\sigma\|_{L^2}^2 = \mu_1^2$, and $\|\mathsf{P}_{>1}\sigma\|_{L^2}^2 = \mu_{\star}^2$.

▶ Tech. ass. on f_* and σ . Almost every nonlinear f_* and σ .

▶ *n* data, *N* features, *d* dimension. $N/d \rightarrow \psi_1$, $n/d \rightarrow \psi_2$, as $d \rightarrow \infty$.

▶ Random features regression: $\hat{a}_{\lambda} = \arg \min_{a} L_{\lambda}(a)$,

$$L_{\lambda}(\boldsymbol{a}) = \frac{1}{n} \sum_{i=1}^{n} \left[\left(y_i - \sum_{j=1}^{N} a_j \sigma(\langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle) \right)^2 \right] + \frac{\lambda N}{d} \|\boldsymbol{a}\|_2^2,$$
$$R(\boldsymbol{a}; f_{\star}) = \mathbb{E}_{\boldsymbol{x}, y} \left[\left(f_{\star}(\boldsymbol{x}) - \sum_{j=1}^{N} a_j \sigma(\langle \boldsymbol{x}, \boldsymbol{w}_j \rangle) \right)^2 \right].$$

- ► Data: $(x_i)_{i \in [n]} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d})), y_i = f_\star(x_i) + \varepsilon_i, \mathbb{E}[\varepsilon_i^2] = \tau^2.$
- ▶ Weights: $(w_j)_{j \in [N]} \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}).$
- Activation σ : $\|\mathsf{P}_1\sigma\|_{L^2}^2 = \mu_1^2$, and $\|\mathsf{P}_{>1}\sigma\|_{L^2}^2 = \mu_{\star}^2$.

▶ Tech. ass. on f_* and σ . Almost every nonlinear f_* and σ .

▶ *n* data, *N* features, *d* dimension. $N/d \rightarrow \psi_1$, $n/d \rightarrow \psi_2$, as $d \rightarrow \infty$.

▶ Random features regression: $\hat{a}_{\lambda} = \arg \min_{a} L_{\lambda}(a)$,

$$L_{\lambda}(\boldsymbol{a}) = \frac{1}{n} \sum_{i=1}^{n} \left[\left(y_i - \sum_{j=1}^{N} a_j \sigma(\langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle) \right)^2 \right] + \frac{\lambda N}{d} \|\boldsymbol{a}\|_2^2,$$
$$R(\boldsymbol{a}; f_{\star}) = \mathbb{E}_{\boldsymbol{x}, y} \left[\left(f_{\star}(\boldsymbol{x}) - \sum_{j=1}^{N} a_j \sigma(\langle \boldsymbol{x}, \boldsymbol{w}_j \rangle) \right)^2 \right].$$

- ► Data: $(\boldsymbol{x}_i)_{i \in [n]} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d})), \ y_i = f_\star(\boldsymbol{x}_i) + \varepsilon_i, \ \mathbb{E}[\varepsilon_i^2] = \tau^2.$
- Weights: $(w_j)_{j \in [N]} \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}).$
- Activation σ : $\|\mathsf{P}_1\sigma\|_{L^2}^2 = \mu_1^2$, and $\|\mathsf{P}_{>1}\sigma\|_{L^2}^2 = \mu_{\star}^2$.

Tech. ass. on f_* and σ . Almost every nonlinear f_* and σ .

▶ *n* data, *N* features, *d* dimension. $N/d \rightarrow \psi_1$, $n/d \rightarrow \psi_2$, as $d \rightarrow \infty$.

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ ● ● ●

▶ Random features regression: $\hat{a}_{\lambda} = \arg \min_{a} L_{\lambda}(a)$,

$$L_{\lambda}(\boldsymbol{a}) = \frac{1}{n} \sum_{i=1}^{n} \left[\left(y_i - \sum_{j=1}^{N} a_j \sigma(\langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle) \right)^2 \right] + \frac{\lambda N}{d} \|\boldsymbol{a}\|_2^2,$$
$$R(\boldsymbol{a}; f_{\star}) = \mathbb{E}_{\boldsymbol{x}, y} \left[\left(f_{\star}(\boldsymbol{x}) - \sum_{j=1}^{N} a_j \sigma(\langle \boldsymbol{x}, \boldsymbol{w}_j \rangle) \right)^2 \right].$$

- ► Data: $(\boldsymbol{x}_i)_{i \in [n]} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d})), \ y_i = f_\star(\boldsymbol{x}_i) + \varepsilon_i, \ \mathbb{E}[\varepsilon_i^2] = \tau^2.$
- Weights: $(w_j)_{j \in [N]} \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}).$
- Activation σ : $\|\mathsf{P}_1\sigma\|_{L^2}^2 = \mu_1^2$, and $\|\mathsf{P}_{>1}\sigma\|_{L^2}^2 = \mu_{\star}^2$.
- ▶ Tech. ass. on f_{\star} and σ . Almost every nonlinear f_{\star} and σ .

▶ *n* data, *N* features, *d* dimension. $N/d \rightarrow \psi_1$, $n/d \rightarrow \psi_2$, as $d \rightarrow \infty$.

▶ Random features regression: $\hat{a}_{\lambda} = \arg \min_{a} L_{\lambda}(a)$,

$$L_{\lambda}(\boldsymbol{a}) = \frac{1}{n} \sum_{i=1}^{n} \left[\left(y_i - \sum_{j=1}^{N} \boldsymbol{a}_j \sigma(\langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle) \right)^2 \right] + \frac{\lambda N}{d} \|\boldsymbol{a}\|_2^2,$$
$$R(\boldsymbol{a}; f_{\star}) = \mathbb{E}_{\boldsymbol{x}, y} \left[\left(f_{\star}(\boldsymbol{x}) - \sum_{j=1}^{N} \boldsymbol{a}_j \sigma(\langle \boldsymbol{x}, \boldsymbol{w}_j \rangle) \right)^2 \right].$$

- ► Data: $(\boldsymbol{x}_i)_{i \in [n]} \sim \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d})), \ y_i = f_\star(\boldsymbol{x}_i) + \varepsilon_i, \ \mathbb{E}[\varepsilon_i^2] = \tau^2.$
- Weights: $(w_j)_{j \in [N]} \sim_{iid} \text{Unif}(\mathbb{S}^{d-1}).$
- Activation σ : $\|\mathsf{P}_1\sigma\|_{L^2}^2 = \mu_1^2$, and $\|\mathsf{P}_{>1}\sigma\|_{L^2}^2 = \mu_{\star}^2$.
- ▶ Tech. ass. on f_{\star} and σ . Almost every nonlinear f_{\star} and σ .
- ▶ *n* data, *N* features, *d* dimension. $N/d \rightarrow \psi_1$, $n/d \rightarrow \psi_2$, as $d \rightarrow \infty$.

Precise asymptotics

Random features regression: $\hat{a}_{\lambda} = \arg \min_{a} L_{\lambda}(a)$,

$$R(\boldsymbol{a};f_{\star})=\mathbb{E}_{\boldsymbol{x},\boldsymbol{y}}\Big[\Big(f_{\star}(\boldsymbol{x})-\sum_{j=1}^{N}\boldsymbol{a}_{j}\sigma(\langle \boldsymbol{x},\boldsymbol{w}_{j}
angle)\Big)^{2}\Big].$$

Theorem (M. and Montanari, 2019)

Under above assumptions, for any $\lambda > 0$, we have

$$egin{aligned} R(\hat{a}_\lambda;f_\star) &= \|\mathsf{P}_{\mathrm{lin}}f_\star\|_{L^2}^2 \cdot \mathscr{B}(\zeta,\psi_1,\psi_2,\lambda/\mu_\star^2) \ &+ (au^2+\|\mathsf{P}_{\mathrm{nl}}f_\star\|_{L^2}^2) \cdot \mathscr{V}(\zeta,\psi_1,\psi_2,\lambda/\mu_\star^2) + o_{d,\mathbb{P}}(1), \end{aligned}$$

where functions \mathscr{B} and \mathscr{V} are given explicitly below.

(Similar result for the training error.)

∃ \0 \0

イロト イポト イヨト イヨト

Explicit formulae

Let the functions $\nu_1,\nu_2:\mathbb{C}_+\to\mathbb{C}_+$ be the unique solution of

$$\begin{split} \nu_1 &= \psi_1 \left(-\xi - \nu_2 - \frac{\zeta^2 \nu_2}{1 - \zeta^2 \nu_1 \nu_2} \right)^{-1}, \\ \nu_2 &= \psi_2 \left(-\xi - \nu_1 - \frac{\zeta^2 \nu_1}{1 - \zeta^2 \nu_1 \nu_2} \right)^{-1}; \end{split}$$

Let

$$\chi\equiv
u_1(oldsymbol{i}(\psi_1\psi_2\overline{\lambda})^{1/2})\cdot
u_2(oldsymbol{i}(\psi_1\psi_2\overline{\lambda})^{1/2}),$$

and

$$\begin{split} \mathscr{E}_{0}\left(\zeta,\psi_{1},\psi_{2},\overline{\lambda}\right) &\equiv -\chi^{5}\zeta^{6} + 3\chi^{4}\zeta^{4} + (\psi_{1}\psi_{2} - \psi_{2} - \psi_{1} + 1)\chi^{3}\zeta^{6} - 2\chi^{3}\zeta^{4} - 3\chi^{3}\zeta^{2} \\ &+ (\psi_{1} + \psi_{2} - 3\psi_{1}\psi_{2} + 1)\chi^{2}\zeta^{4} + 2\chi^{2}\zeta^{2} + \chi^{2} + 3\psi_{1}\psi_{2}\chi\zeta^{2} - \psi_{1}\psi_{2} , \\ \mathscr{E}_{1}\left(\zeta,\psi_{1},\psi_{2},\overline{\lambda}\right) &\equiv \psi_{2}\chi^{3}\zeta^{4} - \psi_{2}\chi^{2}\zeta^{2} + \psi_{1}\psi_{2}\chi\zeta^{2} - \psi_{1}\psi_{2} , \\ \mathscr{E}_{2}\left(\zeta,\psi_{1},\psi_{2},\overline{\lambda}\right) &\equiv \chi^{5}\zeta^{6} - 3\chi^{4}\zeta^{4} + (\psi_{1} - 1)\chi^{3}\zeta^{6} + 2\chi^{3}\zeta^{4} + 3\chi^{3}\zeta^{2} + (-\psi_{1} - 1)\chi^{2}\zeta^{4} - 2\chi^{2}\zeta^{2} - \chi^{2} \end{split}$$

We then have

$$\mathscr{B}(\zeta,\psi_1,\psi_2,\overline{\lambda})\equivrac{\mathscr{E}_1(\zeta,\psi_1,\psi_2,\overline{\lambda})}{\mathscr{E}_0(\zeta,\psi_1,\psi_2,\overline{\lambda})}\,,\qquad \mathscr{V}(\zeta,\psi_1,\psi_2,\overline{\lambda})\equivrac{\mathscr{E}_2(\zeta,\psi_1,\psi_2,\overline{\lambda})}{\mathscr{E}_0(\zeta,\psi_1,\psi_2,\overline{\lambda})}\,.$$

<ロ> (四) (四) (三) (三) (三) (三)
Insights

- 2

・ロト ・御ト ・ヨト ・ヨト



- \checkmark Peak at the interpolation threshold.
- ✓ Monotonic decreasing in the overparameterized regime.
- \checkmark Global minimum when the number of parameters is infinity.

- 34

イロト イポト イヨト イヨト



For any λ, the min prediction error is achieved at N/n → ∞.
 For optimal λ, the prediction error is monotonically decreasing.

- 34

・ロト ・四ト ・ヨト ・ヨト



▶ High SNR: minimum at $\lambda = 0+$;

▶ Low SNR: minimum at $\lambda > 0$.

- 34

<ロト <問ト < 注ト < 注ト

1. Mean field theory

2. Tangent kernel theory

3. Transitions between mean field and tangent kernel regime

- 31

<ロ> (日) (日) (日) (日) (日)

Connections of mean field and tangent kernel

Setup: α controls the speed of change of emp. dist.

$$\begin{array}{ll} \text{Prediction function:} \quad \hat{f}_{\alpha,N}(x;\boldsymbol{\theta}) = & \frac{\alpha}{N} \sum_{j=1}^{N} \sigma_{\star}(x;\boldsymbol{\theta}_{j}),\\\\ \text{Risk function:} \quad & R_{\alpha,N}(\boldsymbol{\theta}) = & \mathbb{E}_{\boldsymbol{x}} \left[\left(f(\boldsymbol{x}) - \hat{f}_{\alpha,N}(\boldsymbol{x};\boldsymbol{\theta}) \right)^{2} \right],\\\\ \text{Gradient flow:} \quad & \frac{\mathrm{d}\boldsymbol{\theta}_{j}^{t}}{\mathrm{d}t} = - \frac{N}{2\alpha^{2}} \nabla_{\boldsymbol{\theta}_{j}} R_{\alpha,N}(\boldsymbol{\theta}^{t}). \end{array}$$

- 32

イロト イヨト イヨト イヨ

The coupled dynamics

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

$$\partial_t
ho_t^{lpha,N} = (1/lpha)
abla_{m{ heta}} \cdot (
ho_t^{lpha,N}
abla_{m{ heta}} \Psi(m{ heta};
ho_t^{lpha,N})).$$

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

$$\partial_t \|u^{lpha,N}_t\|^2_{L^2} = -\langle u^{lpha,N}_t, \mathcal{H}_{
ho^{lpha,N}_t}u^{lpha,N}_t
angle.$$

Here

[Rotskoff, Vanden-Eijnden, 2018], [Chizat, Bach, 2018b]

- 34

・ロト ・四ト ・ヨト ・ヨト

The coupled dynamics

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

$$\partial_t
ho_t^{lpha,N} = (1/lpha)
abla_{oldsymbol{ heta}} \cdot (
ho_t^{lpha,N}
abla_{oldsymbol{ heta}} \Psi(oldsymbol{ heta};
ho_t^{lpha,N})).$$

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

$$\partial_t \|u^{lpha,N}_t\|^2_{L^2} = -\langle u^{lpha,N}_t, \mathcal{H}_{
ho^{lpha,N}_t}u^{lpha,N}_t
angle.$$

Here

[Rotskoff, Vanden-Eijnden, 2018], [Chizat, Bach, 2018b]

◆□ ▶ ◆□ ▶ ▲ ■ ▶ ▲ ■ ▶ ● ● ●

The coupled dynamics

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

$$\partial_t
ho_t^{lpha,N} = (1/lpha)
abla_{oldsymbol{ heta}} \cdot (
ho_t^{lpha,N}
abla_{oldsymbol{ heta}} \Psi(oldsymbol{ heta};
ho_t^{lpha,N})).$$

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

$$\partial_t \|u^{lpha,N}_t\|^2_{L^2} = -\langle u^{lpha,N}_t, \mathcal{H}_{
ho^{lpha,N}_t}u^{lpha,N}_t
angle.$$

Here

$$egin{aligned} \mathcal{H}_{
ho}(oldsymbol{x},oldsymbol{z}) &\equiv \int \langle
abla_{oldsymbol{ heta}} \sigma_{\star}(oldsymbol{x};oldsymbol{ heta}),
abla_{oldsymbol{ heta}} \sigma_{\star}(oldsymbol{z};oldsymbol{ heta})
angle
ho(\mathrm{d}oldsymbol{ heta}) \ \Psi_{lpha}(oldsymbol{ heta};
ho^{lpha,N}) &= - \mathbb{E}_{oldsymbol{x}}[u_t^{lpha,N}(oldsymbol{x})\sigma_{\star}(oldsymbol{x};oldsymbol{ heta})]. \end{aligned}$$

[Rotskoff, Vanden-Eijnden, 2018], [Chizat, Bach, 2018b]

◆□ ▶ ◆□ ▶ ▲ ■ ▶ ▲ ■ ▶ ● ● ●

The mean field limit and tangent kernel limit

$$egin{aligned} \partial_t
ho_t^{lpha,N} &= (1/lpha)
abla_{m{ heta}} \cdot (
ho_t^{lpha,N} [
abla_{m{ heta}} \Psi(m{ heta};
ho_t^{lpha,N})]), \ \partial_t \|u_t^{lpha,N}\|_{L^2}^2 &= -\langle u_t^{lpha,N}, \mathcal{H}_{
ho_t^{lpha,N}} u_t^{lpha,N}
angle. \end{aligned}$$

▶ The mean field limit: fix $\alpha = O(1)$ and let $N \to \infty$.

• The tangent kernel limit: let $\alpha = \sqrt{N} \to \infty$.

In tangent kernel limit: the kernel will not change. The res. dynamics becomes self contained. The emp. risk converges to 0.

イロト イヨト イヨト イヨ

The mean field limit and tangent kernel limit

$$egin{aligned} \partial_t
ho_t^{lpha,N} =& (1/lpha)
abla_{oldsymbol{ heta}} \cdot (
ho_t^{lpha,N} [
abla_{oldsymbol{ heta}} \Psi(oldsymbol{ heta};
ho_t^{lpha,N})]), \ \partial_t \|u_t^{lpha,N}\|_{L^2}^2 =& -\langle u_t^{lpha,N}, \mathcal{H}_{
ho_t^{lpha,N}} u_t^{lpha,N}
angle. \end{aligned}$$

- The mean field limit: fix $\alpha = O(1)$ and let $N \to \infty$.
- The tangent kernel limit: let $\alpha = \sqrt{N} \to \infty$.

In tangent kernel limit: the kernel will not change. The res. dynamics becomes self contained. The emp. risk converges to 0.

< ロ > < 同 > < 回 > < 回 >

The mean field limit and tangent kernel limit

$$egin{aligned} \partial_t
ho_t^{lpha,N} =& (1/lpha)
abla_{oldsymbol{ heta}} \cdot (
ho_t^{lpha,N} [
abla_{oldsymbol{ heta}} \Psi(oldsymbol{ heta};
ho_t^{lpha,N})]), \ \partial_t \|u_t^{lpha,N}\|_{L^2}^2 =& -\langle u_t^{lpha,N}, \mathcal{H}_{
ho_t^{lpha,N}} u_t^{lpha,N}
angle. \end{aligned}$$

- The mean field limit: fix $\alpha = O(1)$ and let $N \to \infty$.
- The tangent kernel limit: let $\alpha = \sqrt{N} \to \infty$.
- In tangent kernel limit: the kernel will not change. The res. dynamics becomes self contained. The emp. risk converges to 0.

< 日 > < 同 > < 回 > < 回

Benefits and limitations of the mean field theory

$$\partial_t \rho_t = \nabla_{\boldsymbol{\theta}} \cdot \left(\rho_t \nabla_{\boldsymbol{\theta}} \Psi(\boldsymbol{\theta}; \rho_t) \right).$$

Benefits:

- It captures the non-linear behavior of neural networks that potentially give better generalization.
- People believe in practice we are in this regime.

Limitations:

- Hard to prove convergence of PDE.
- ▶ Hard to generalize to multi-layers.

< ロ > < 同 > < 回 > < 回 > < 回

Benefits and limitations of the mean field theory

$$\partial_t \rho_t = \nabla_{\boldsymbol{\theta}} \cdot \left(\rho_t \nabla_{\boldsymbol{\theta}} \Psi(\boldsymbol{\theta}; \rho_t) \right).$$

Benefits:

- It captures the non-linear behavior of neural networks that potentially give better generalization.
- People believe in practice we are in this regime.

Limitations:

- Hard to prove convergence of PDE.
- Hard to generalize to multi-layers.

・ロト ・同ト ・ヨト ・ヨ

Benefits and limitations of the tangent kernel theory

Benefits:

- ▶ Give provable convergence results for multi-layers neural networks.
- Exhibit many similar behaviors of neural networks: double-descent.
- Easy to use in many scenarios: ResNet, convolutional NN, graph NN, recurrent NN.

Limitations

- ▶ The intuition of fixed weights distribution is not realistic.
- ▶ The generalization is not as good as fully trained neural networks.
- Everything is just a kernel?

イロト イボト イヨト イヨト

Benefits and limitations of the tangent kernel theory

Benefits:

- ▶ Give provable convergence results for multi-layers neural networks.
- Exhibit many similar behaviors of neural networks: double-descent.
- Easy to use in many scenarios: ResNet, convolutional NN, graph NN, recurrent NN.

Limitations

- ▶ The intuition of fixed weights distribution is not realistic.
- ▶ The generalization is not as good as fully trained neural networks.
- Everything is just a kernel?

- -

・ロッ ・雪 ・ ・ ヨ ・

Thanks!

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへで