

Gradient Descent for Deep Neural Networks: New Perspectives from Mean-field and NTK

Marco Mondelli

Institute of Science and Technology Austria (IST Austria)

Math ML seminar MPI MiS + UCLA, 3 March 2022



Understanding Gradient Descent for Over-parameterized Deep Neural Networks

Marco Mondelli

Institute of Science and Technology Austria (IST Austria)

Mathematics of Data Seminar, MPI MiS, 4 August 2020



New Perspectives from Mean-field and NTK

[Picture credit: F. Draxler et al.]

mode connectivity

[Picture credit: F. Draxler et al.]

[Picture credit: K. Tanir, *Getty Images*]

mode connectivity

memorization

[Picture credit: F. Draxler et al.]

[Picture credit: K. Tanir, *Getty Images*]

mode connectivity

memorization

convergence

[Picture credit: F. Draxler et al.]

[Picture credit: K. Tanir, *Getty Images*]

mode connectivity

memorization

convergence

implicit bias

A mean-field view

!

Two-layer networks

Data: $f(\mathbf{x}_1; y_1); \dots; (\mathbf{x}_n; y_n) \mathcal{G}$ i.i.d. $\mathcal{P}(\mathbb{R}^d \times \mathbb{R})$

Goal: Minimize loss $L_N(\cdot) = \mathbb{E} \sum_{i=1}^N y_i \left(\frac{1}{N} \sum_{j=1}^N a_j(\mathbf{x}; \mathbf{w}_j) \right)^2$, $\mathcal{L} = \mathcal{L}(\mathbf{w}; a)$

Online SGD: $\mathbf{w}^{k+1} = \mathbf{w}^k - \eta \sum_{i=1}^N y_k a_i^k(\mathbf{x}_k; \mathbf{w}_i^k)$

SGD minimizes loss $L_N(\theta)$

$$L_N(\theta) = \mathbb{E} \sum_{i=1}^N \frac{1}{N} \sum_{j=1}^M a_i(\mathbf{x}; \mathbf{w}_i)^2; \quad \theta = (\mathbf{w}; \mathbf{a})$$

Microscopic discrete dynamics of a gas with N particles

SGD minimizes loss $L_N(\cdot)$

$$L_N(\cdot) = \mathbb{E} \sum_{i=1}^n \frac{1}{N} \sum_{\mathbf{x}^i} a_i(\mathbf{x}; \mathbf{w}_i)^2; \quad \mathbf{w} = (\mathbf{w}; a)$$

$L_N(\cdot)$ depends only on empirical distribution $b^{(N)} = \frac{1}{N} \sum_{i=1}^n \delta_{\mathbf{x}^i}$

SGD minimizes loss $L_N(\cdot)$

$$L_N(\cdot) = \mathbb{E} \sum_{i=1}^n y_i \left(\frac{1}{N} \sum_{j=1}^N a_j(\mathbf{x}; \mathbf{w}_j) \right)^2; \quad \mathbf{w} = (\mathbf{w}; a)$$

$L_N(\cdot)$ depends only on empirical distribution $b^{(N)} = \frac{1}{N} \sum_{i=1}^N \delta_{\mathbf{x}_i}$

$$\frac{1}{N} \sum_{i=1}^N a_i(\mathbf{x}; \mathbf{w}_i) = \int_{\mathcal{Z}} a(\mathbf{x}; \mathbf{w}) b^{(N)}(d\mathbf{w})$$

|-----{Z}-----|
expected a w.r.t. empirical distribution $b^{(N)}$

SGD minimizes loss $L_N(\cdot)$

$$L_N(\cdot) = \mathbb{E}_y \frac{1}{N} \sum_{i=1}^N a_i(\mathbf{x}; \mathbf{w}_i)^2; \quad = (\mathbf{w}; a)$$

expected a
w.r.t. empirical distribution $b^{(N)}$

relaxation

$$L(\cdot) = \mathbb{E}_y \int_{\mathcal{Z}} a(\mathbf{x}; \mathbf{w})^2 (d\mu d\mathbf{w})$$

expected a
w.r.t.

L takes as input a probability distribution

Already exploited in [Bartlett, 1998; Bengio et al., 2006]

SGD ! gradient flow PDE

As $N \rightarrow \infty$ and $\eta \rightarrow 0$, we have:

$$b_k^{(N)} = \frac{1}{N} \sum_{i=1}^N \langle \cdot, \cdot \rangle_k \quad k'' = \text{gradient flow PDE for } L(\cdot)$$

Propagation-of-chaos argument [Sznitman, 1991; Mei et al., 2018]

SGD ! gradient flow PDE

[Picture credit: Peletier]

!

Discrete dynamics of SGD**Continuous dynamics** of
gradient flow

Two layers [Mei et al., 2018, 2019; Rotskoff et al., 2018;
Chizat et al., 2018; Sirignano et al., 2018; :::]

Multiple layers [Nguyen, 2019; Sirignano et al., 2019; Araujo
et al., 2019; Pham et al., 2020; Fang et al., 2021; :::]

Three vignettes

1. SGD weights close to **i.i.d.** particles
) *dropout stability and connectivity*

A. Shevchenko and M. Mondelli, “Landscape Connectivity and Dropout Stability of SGD Solutions for Over-parameterized Neural Networks”, *ICML*, 2020.

Three vignettes

1. SGD weights close to **i.i.d.** particles
) *dropout stability and connectivity*
2. Particles evolve with **gradient flow**
) *exponential convergence rate*

A. Javanmard, M. Mondelli, and A. Montanari, “Analysis of a Two-Layer Neural Network via Displacement Convexity”, *Annals of Statistics*, 2020.

Three vignettes

1. SGD weights close to **i.i.d.** particles
 -) *dropout stability and connectivity*
2. Particles evolve with **gradient flow**
 -) *exponential convergence rate*
3. **Stationary distribution** has Gibbs form
 -) *implicit bias*

A. Shevchenko, V. Kungurtsev, and M. Mondelli, “Mean-field Analysis of Piecewise Linear Solutions for Wide ReLU Networks”, arXiv preprint, 2021.

Vignette 1: Dropout stability and connectivity

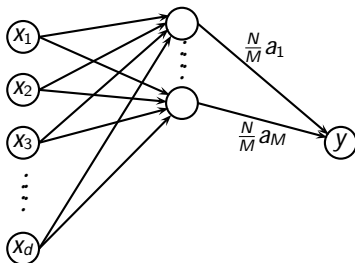
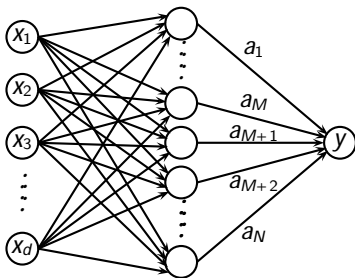
SGD minima connected via piecewise linear path with low loss [Garipov et al., 2018; Draxler et al., 2018; Fort et al., 2019; Anokhin et al., 2020]

Mode connectivity proved assuming properties of well-trained networks (dropout stability and noise stability) [Kuditipudi et al., 2019]

Dropout stability

$$L_M(\cdot) = \mathbb{E} \left[\frac{1}{M} \sum_{i=1}^M a_i(\mathbf{x}; \mathbf{w}_i) \right]^2$$

is **"D-dropout stable** if $\|L_N(\cdot) - L_M(\cdot)\|_D$



Dropout stability and connectivity

$$L_M(\theta) = \mathbb{E}_y \left[\frac{1}{M} \sum_{i=1}^M a_i(\mathbf{x}; \mathbf{w}_i) \right]^2$$

is **" \mathcal{D} -dropout stable** if $\|L_N(\theta) - L_M(\theta)\| \leq \epsilon_{\mathcal{D}}$

and θ^0 are **" \mathcal{C} -connected** if there exists a continuous path connecting them where the loss does not increase more than $\epsilon_{\mathcal{C}}$

Dropout stability, provably

N = # neurons of full network

η = step size of SGD

M = # neurons after dropout

D = dimension of weights

Theorem (Shevchenko and M., 2020)

Let θ^k be obtained after k SGD iterations. Then, with probability $1 - e^{-z^2}$, for all $k \geq \lceil T \rceil$, θ^k is " D -dropout stable with

$$\epsilon_D = Ke^{kT^3} \left(\frac{\sqrt{\log M} + z}{\sqrt{M}} + \frac{\sqrt{D + \log N} + z}{\sqrt{D + \log N}} \right) :$$

Change in loss scales as $\frac{\sqrt{\log M}}{\sqrt{M}} + \frac{\sqrt{D + \log N}}{\sqrt{D + \log N}}$

Dropout stability, provably

N = # neurons of full network

η = step size of SGD

M = # neurons after dropout

D = dimension of weights

Theorem (Shevchenko and M., 2020)

Let θ^k be obtained after k SGD iterations. Then, with probability $1 - e^{-z^2}$, for all $k \geq \lceil T \rceil$, θ^k is " D -dropout stable with

$$\epsilon_D = Ke^{KT^3} \left(\frac{P \sqrt{\log M + z}}{P \overline{M}} + P - P \frac{P}{D + \log N + z} \right) :$$

k close to N i.i.d. particles that evolve with gradient flow

$L_N(\theta^k)$ and $L_M(\theta^k)$ concentrate to the same limit

Connectivity, provably

Theorem (Shevchenko and M., 2020)

Let \mathbf{x}^k be obtained after k SGD iterations using $f(\mathbf{x}_j; \mathbf{y}_j)_{\mathcal{G}_{j=0}^k}$ \mathbb{P} , and $(\mathbf{x}^0)^{k^0}$ after k^0 SGD iterations using $f(\mathbf{x}_j^0; \mathbf{y}_j^0)_{\mathcal{G}_{j=0}^{k^0}}$ \mathbb{P} . Then, with probability $1 - e^{-z^2}$, for all $k \geq \lceil T = \rceil$ and $k^0 \geq \lceil T^0 = \rceil$, \mathbf{x}^k and $(\mathbf{x}^0)^{k^0}$ are ϵ_C -connected with

$$\epsilon_C = Ke^{K \max(T; T^0)^3} \frac{\mathbb{P} \overline{\log N} + z}{\mathbb{P} \overline{N}} + \mathbb{P} \overline{D + \log N} + z :$$

Change in loss scales as $\mathbb{P} \overline{\frac{\log N}{N}} + \mathbb{P} \overline{\frac{1}{(D + \log N)}}$

Dropout stability with $M = N=2$) connectivity

Connectivity, provably

Theorem (Shevchenko and M., 2020)

Let \mathbf{x}^k be obtained after k SGD iterations using $f(\mathbf{x}_j; y_j)_{j=0}^k$ \mathbb{P} , and $(\mathbf{y}^k)^{k^0}$ after k^0 SGD iterations using $f(\mathbf{x}_j^0; y_j^0)_{j=0}^{k^0}$ \mathbb{P} . Then, with probability $1 - e^{-z^2}$, for all $k \geq \lceil T = \lceil \frac{P}{\epsilon} \log N + z \rceil \rceil$ and $k^0 \geq \lceil T^0 = \lceil \frac{P}{\epsilon} \log N + z \rceil \rceil$, \mathbf{x}^k and $(\mathbf{y}^k)^{k^0}$ are ϵ_C -connected with

$$\epsilon_C = Ke^{K \max(T; T^0)^3} \frac{P \overline{\log N} + z}{P \overline{N}} + P - P \frac{P \overline{\log N} + z}{D + \log N + z} :$$

Bounds similar in spirit for the **multilayer** case

Vignette 2: Exponential convergence rate

Data model

Data: $f(\mathbf{x}_1; y_1); \dots; (\mathbf{x}_n; y_n) \mathcal{G}$ i.i.d: $\mathbb{P}(\mathbb{R}^d \times \mathbb{R})$

$\mathbf{x}_i \sim \text{Unif}(\cdot)$ and $y_i = f(\mathbf{x}_i) + \epsilon_i$

Data model

Data: $f(\mathbf{x}_1; y_1); \dots; (\mathbf{x}_n; y_n) \mathcal{G}$ i.i.d: $\mathbb{P}(\mathbb{R}^d \times \mathbb{R})$

$\mathbf{x}_i \sim \text{Unif}(\Omega)$ and $y_i = f(\mathbf{x}_i) + \epsilon_i$

Ω bounded and convex

f positive and smooth

f μ -strongly concave
(lower bound on curvature)

ϵ_i sub-Gaussian

Problem statement

Goal: Minimize loss $L_N(\mathbf{w}) = \mathbb{E}_y \frac{1}{N} \sum_{i=1}^N \ell(\mathbf{x}; \mathbf{w}_i)$

Problem statement

Goal: Minimize loss $L_N(\mathbf{w}) = \mathbb{E}_y \frac{1}{N} \sum_{i=1}^N \ell(\mathbf{x}; \mathbf{w}_i)$

'bump'-like

$$\ell(\mathbf{x}; \mathbf{w}_i) = K(\|\mathbf{x} - \mathbf{w}_i\|)$$

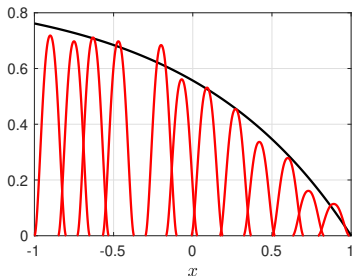
\mathbf{w}_i = center of bump

σ = width of the bump

$$K(\|\mathbf{x}\|) = \frac{1}{\sigma^d} K(\|\mathbf{x}\|/\sigma)$$

Problem statement

Goal: Minimize loss $L_N(\mathbf{w}) = \mathbb{E}_n \left[\frac{1}{N} \sum_{i=1}^N K(\mathbf{x}_i, \mathbf{w}_i) \right]^2$



Fit data with combination of bumps

Learn N centers \mathbf{w} via gradient descent

Non-convex optimization problem in $N \cdot d$ dimensions

Exponential and dimension-free convergence rate

Theorem (Javanmard, M. and Montanari, 2020)

Let f be μ -strongly concave, σ a bump of width δ , and w^k obtained after k SGD iterations with constant step size η . Then, with high probability,

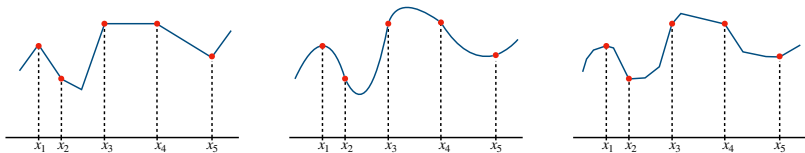
$$\underbrace{L_N(w^k)}_{\text{loss at step } k} \leq \underbrace{L_N(w^0)}_{\text{exp. rate}} e^{-\frac{\mu}{\sigma} k} + \underbrace{\frac{\sigma}{\mu} \left(\frac{N}{\delta} \right)}_{\text{distance to opt.}};$$

where $\left(\frac{N}{\delta} \right) \rightarrow 0$ as $N \rightarrow \infty$, $\delta \rightarrow 0$, and $\frac{\sigma}{\mu} \rightarrow 0$.

As $\frac{\sigma}{\mu} \rightarrow 0$, gradient flow optimizes a **displacement convex** loss.

Vignette 3: Implicit bias

Which solution does SGD choose?



Related work

Solution **implicitly biased** towards:::

Total variation of curvature [Savarese et al., 2019; Ongie et al., 2020]

Max-margin [Soudry et al., 2018; Chizat et al., 2020]

Cubic splines [Williams et al., 2019; Jin et al., 2020]

Linear splines [Williams et al., 2019; Ergen et al., 2021]

Minimum number of tangent changes [Blanc et al., 2020]

Best fit closest to initialization [Zhang et al., 2020]

SGD + ReLU = piecewise linear solution

Theorem (Shevchenko, Kungurtsev and M., 2021)

(Informal) Consider training a two-layer ReLU network via noisy-SGD, in the low noise regime. Then, the resulting estimator converges to a **piecewise linear solution**, and the number of **knot points** – i.e., points at which distinct linear pieces connect – between two consecutive training inputs is **at most 3**.

SGD + ReLU = piecewise linear solution

Theorem (Shevchenko, Kungurtsev and M., 2021)

(Informal) Consider training a two-layer ReLU network via noisy-SGD, in the low noise regime. Then, the resulting estimator converges to a **piecewise linear solution**, and the number of **knot points** – i.e., points at which distinct linear pieces connect – between two consecutive training inputs is **at most 3**.

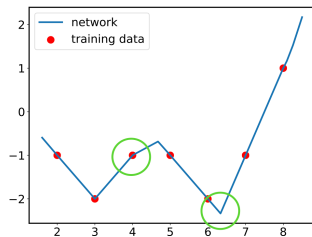
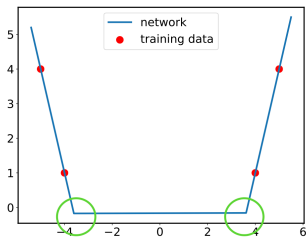
Qualitatively different behaviour from [Williams et al., 2019; Blanc et al., 2020; Jin et al., 2020; Ergen et al., 2021].

Analysis holds for both constant and vanishing λ_2 regularization.

SGD + ReLU = piecewise linear solution

Theorem (Shevchenko, Kungurtsev and M., 2021)

(*Informal*) Consider training a two-layer ReLU network via noisy-SGD, in the low noise regime. Then, the resulting estimator converges to a **piecewise linear solution**, and the number of **knot points** – i.e., points at which distinct linear pieces connect – between two consecutive training inputs is **at most 3**.



Two-layer network trained on a 1D dataset

Data: $f(x_1; y_1); \dots; (x_n; y_n) \mathcal{g}$ i.i.d: $\mathbb{P}(\mathbb{R} \times \mathbb{R})$

Data distribution supported on M points: $\mathbb{P} = \frac{1}{M} \sum_{j=1}^M \delta_{(x_j; y_j)}$

Goal: Minimize loss

$$L_N(\cdot) = \mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N a_i (x; \mathbf{w}_i)^2 \right]; \quad \cdot = (\mathbf{w}; \mathbf{a})$$

Online SGD:

$$\mathbf{w}^{k+1} = \mathbf{w}^k - \frac{1}{N} \sum_{i=1}^N a_i^k (x_k; \mathbf{w}_i^k)^2$$

Two-layer network trained on a 1D dataset

Data: $f(x_1; y_1); \dots; (x_n; y_n) \sim \text{i.i.d. } P(\mathbb{R} \times \mathbb{R})$

Data distribution supported on M points: $P = \frac{1}{M} \sum_{j=1}^M \delta_{(x_j; y_j)}$

Goal: Minimize **regularized** loss

$$L_N(\mathbf{w}; \mathbf{a}) = \frac{1}{M} \sum_{j=1}^M y_j \left(\frac{1}{N} \sum_{i=1}^N a_i (x_j; \mathbf{w}_i) \right)^2 + \frac{1}{N} \sum_{i=1}^N \|\mathbf{w}_i\|^2; \quad \mathbf{w}; \mathbf{a}$$

Online SGD:

$$\mathbf{w}^{k+1} = \mathbf{w}^k - \frac{1}{N} \sum_{i=1}^N a_i^k (x_k; \mathbf{w}_i^k)^2$$

Two-layer network trained on a 1D dataset

Data: $f(x_1; y_1); \dots; (x_n; y_n) \mathcal{g}$ i.i.d: $\mathcal{P}(\mathbb{R} \times \mathbb{R})$

Data distribution supported on M points: $\mathcal{P} = \frac{1}{M} \sum_{j=1}^M \delta_{(x_j; y_j)}$

Goal: Minimize **regularized** loss

$$L_N(\mathbf{w}; \mathbf{a}) = \frac{1}{M} \sum_{j=1}^M y_j - \frac{1}{N} \sum_{i=1}^N a_i (x_j; \mathbf{w}_i)^2 + \frac{1}{N} \sum_{i=1}^N k_i^2; \quad \mathbf{g} = (\mathbf{w}; \mathbf{a})$$

Online noisy-SGD:

$$\mathbf{g}^{k+1} = (1 - \eta) \mathbf{g}^k - \eta \left(\frac{1}{N} \sum_{i=1}^N a_i^k (x_k; \mathbf{w}_i^k)^2 - \mathbf{g}^k \right) + \mathbf{g}^k$$

$$\mathbf{g}^k \sim N(0; \mathbf{I})$$

Mean-field analysis

As $N \rightarrow \infty$ and $\eta \rightarrow 0$, μ^k close to N i.i.d. particles
 solution of gradient flow $\mu^k =$

Mean-field analysis

As $N \rightarrow \infty$ and $\beta \rightarrow 0$, μ^k close to N i.i.d. particles $\mu^k =$
 solution of gradient flow

As $t = k \rightarrow \infty$, μ_t close to μ^* = minimizer of free-energy

$$\mu^k = \frac{1}{Z} \exp(\dots)$$

$$\mu^k = \frac{1}{M} \prod_{j=1}^M \int a^j(x_j; \mathbf{w}^j) (da^j d\mathbf{w}^j) \prod_j a(x_j; \mathbf{w}) + \frac{1}{2} k \frac{1}{k^2}$$

Main technical result

Theorem (Shevchenko, Kungurtsev and M., 2021)

Consider the **curvature** of the neural network estimator:

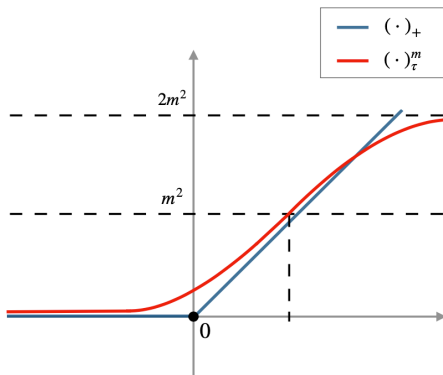
$$C(x) = \frac{d^2}{dx^2} \mathbb{E} a(x; \mathbf{w}) \quad (\text{dadw}):$$

Then, as $n \rightarrow \infty$, $C(x) = 0$ for all x except a set of **vanishing measure** which concentrates on **at most 3 points** per consecutive training inputs.

Characterization of $C(x)$ capturing the location of the knot points

Proof ideas (1)

$a(x; \mathbf{w}) = a(x \mathbf{w} + b)_+$ neither bounded nor smooth)
 need a sequence of approximations



Proof ideas (2)

Taking two derivatives of ReLU gives Dirac delta

$$C(x) = \frac{d^2}{dx^2} \int a(x; \mathbf{w}) (da/d\mathbf{w}) \int a w^2 (a; \mathbf{w}; \mathbf{w} x) da d\mathbf{w}$$

Proof ideas (2)

Taking two derivatives of ReLU gives Dirac delta) analysis of 2D 'Gaussian-like' integral

$$C(x) = \frac{d^2}{dx^2} \int a(x; \mathbf{w}) \, (da d\mathbf{w}) \int a w^2 (a; \mathbf{w}; \mathbf{w} x) da d\mathbf{w}$$

$$(a; \mathbf{w}; \mathbf{w} x) = \frac{\exp\left[-\frac{1}{2} (2aw(A^j x - B^j) + a^2 + w^2 + w^2 x^2)\right]}{Z}$$

$w > 0$, $x \in [x_j, x_{j+1}]$

A^j and B^j coefficients (independent of $x; a; w$)

Z normalization constant

Proof ideas (3)

Taking two derivatives of ReLU gives Dirac delta) analysis of 2D 'Gaussian-like' integral

$$C(x) = \int \frac{\exp\left(-\frac{n}{2} (2aw(A^j x - B^j) + a^2 + w^2 + w^2 x^2)\right)}{Z} da dw$$

Proof ideas (3)

Taking two derivatives of ReLU gives Dirac delta) analysis of 2D 'Gaussian-like' integral

$$C(x) = \int \frac{\exp\left(-\frac{n}{2} (2aw(A^j x - B^j) + a^2 + w^2 + w^2 x^2)\right)}{Z} da dw$$

$C(x)$ blows up when covariance of Gaussian not PSD;
otherwise $C(x) \neq 0$

A Neural Tangent Kernel (NTK) view

$$\mathbf{K}^{(L)} = \sum_{l=1}^L \frac{\partial \mathbf{F}_L}{\partial \text{vec}(\mathbf{W}_l)} \frac{\partial \mathbf{F}_L}{\partial \text{vec}(\mathbf{W}_l)}^T \in \mathbb{R}^{N \times N}$$

Neural Tangent Kernel (NTK)

N training samples in \mathbb{R}^d

Neural network of depth L

$\mathbf{W}_l \in \mathbb{R}^{n_{l-1} \times n_l}$ weights of the network ($n_0 = d$)

$\mathbf{F}_L \in \mathbb{R}^N$ network output corresponding to training data

$$\mathbf{K}^{(L)} = \sum_{l=1}^L \frac{\partial \mathbf{F}_L}{\partial \text{vec}(\mathbf{W}_l)} \frac{\partial \mathbf{F}_L}{\partial \text{vec}(\mathbf{W}_l)}^T \in \mathbb{R}^{N \times N}$$

NTK properties

Convergence to **non-random limit** [Jacot et al., 2018]

$$\mathbb{E} \left[\frac{1}{n} \sum_{l=1}^L \mathbf{K}^{(l)} \right] \approx \mathbf{K}^{(L)}$$

Connection to **memorization** and **generalization bounds** [Arora et al., 2019; Montanari et al., 2020]

Connection to **convergence of gradient descent**

$$k_{\text{min}} \leq k_2 \leq k_{\text{max}}$$

NTK properties

Convergence to **non-random limit** [Jacot et al., 2018]

$$\frac{1}{n} \sum_{l=1}^L \mathbf{K}^{(l)}$$

Connection to **memorization** and **generalization bounds** [Arora et al., 2019; Montanari et al., 2020]

Connection to **convergence of gradient descent**

$$\frac{1}{n} \sum_{l=1}^L \mathbf{K}^{(l)}$$

Only **two layers** or **all layers with poly(N) neurons** /

Two vignettes

1. Global convergence for deep networks with one wide layer + pyramidal topology

Q. Nguyen and M. Mondelli, “Global Convergence of Deep Networks with One Wide Layer Followed by Pyramidal Topology”, *NeurIPS*, 2020.

Two vignettes

1. Global convergence for deep networks with one wide layer + pyramidal topology
2. Bounds on NTK spectrum for deep ReLU networks with one wide layer

Q. Nguyen, M. Mondelli and G. Montufar, “Tight Bounds on the Smallest Eigenvalue of the Neural Tangent Kernel for Deep ReLU Networks”, *ICML*, 2021.

Convergence for (very) wide networks

	Deep?	Activation	Layer Width	Train All Layers?	# Wide Layers
[Oymak et al., '20]	No	Smooth	$(N^2 \text{ }_0^2)$	No	x
[Song et al., '21]	No	Smooth	$(N^{3=2})$	Yes	x
[Allen-Zhu et al., '18]	Yes	General	$(N^{24} L^{12} \text{ }_4)$	No	All
[Allen-Zhu et al., '18]	Yes	General	$(N^{24} L^{12} \text{ }_4)$	No	All
[Zou et al., '19]	Yes	ReLU	$(N^8 L^{12} \text{ }_4)$	No	All
[Du et al., '19]	Yes	Smooth	$(\frac{N^4 2^{O(L)}}{\min(\mathbf{K}^{(L)})})$	Yes	All

Vignette 1: Global convergence

	Deep?	Activation	Layer Width	Train All Layers?	# Wide Layers
[Oymak et al., '20]	No	Smooth	$(N^2 \ 0 \ 2)$	No	x
[Song et al., '21]	No	Smooth	$(N^{3=2})$	Yes	x
[Allen-Zhu et al., '18]	Yes	General	$(N^{24} L^{12} \ 4)$	No	All
[Zou et al., '19]	Yes	ReLU	$(N^8 L^{12} \ 4)$	No	All
[Du et al., '19]	Yes	Smooth	$(\frac{N^4 2^{O(L)}}{4 \min(K^{(L)})})$	Yes	All
[Nguyen & M., '20] (Gen.)	Yes	Smooth	N	Yes	One
[Nguyen & M., '20] (LeCun)	Yes	Smooth	$\Omega(N^2 2^{O(L)})$	Yes	One

Need only one wide layer + pyramidal topology

Connection to the NTK

General recipe

Lipschitz gradient + PL inequality \Rightarrow Linear convergence of GD

$$kr \leq k_2^2 \quad \min \mathbf{K}^{(L)} \geq 2$$

Connection to the NTK

General recipe

Lipschitz gradient + PL inequality \Rightarrow Linear convergence of GD

$$kr \leq k_2^2 \min \mathbf{K}^{(L)} \quad + \quad (\text{pyramidal topology})$$

$$kr \leq k \min (\mathbf{F}_1) \prod_{l=3}^L \min (\mathbf{W}_l)^{p_l} \quad (\text{weak PL inequality})$$

r = loss function

\mathbf{F}_1 feature matrix of first layer

\mathbf{W}_l weight matrix of layer l

Maintaining the (weak) PL inequality

General recipe

Lipschitz gradient + PL inequality \Rightarrow Linear convergence of GD

$$kr \quad k \quad \min_{l=3} (\mathbf{F}_l)^{\Psi} \quad \min_{l=3} (\mathbf{W}_l)^{P-} \quad (\text{weak PL inequality})$$

$\mathbf{F}_1; \mathbf{W}_3; \dots; \mathbf{W}_L$ are **full rank** at initialization

Maintaining the (weak) PL inequality

General recipe

Lipschitz gradient + PL inequality \Rightarrow Linear convergence of GD

$$k r \quad k \quad \min_{l=3} (\mathbf{F}_l) \quad \prod_{l=3}^L \min (\mathbf{W}_l)^{p_l} \quad (\text{weak PL inequality})$$

$\mathbf{F}_1; \mathbf{W}_3; \dots; \mathbf{W}_L$ stay **full rank** during training



$$k r \quad () \quad r \quad (\theta) k \quad Q(; \theta) k \quad \theta k \quad (\text{locally Lipschitz gradient})$$

Open questions

ReLU activations?

Arbitrary topologies?

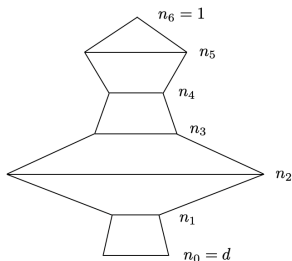
Trade-off width for depth (e.g., two layers with $\binom{P}{N}$ neurons, rather than one layer with (N) neurons)?

Vignette 2: Bounds on NTK spectrum

Main results (Nguyen, M. and Montufar, 2021)

Tight **lower bounds on the smallest eigenvalue** of NTK matrices for deep ReLU networks:

- 1 **Asymptotic** setting with infinite widths
- 2 Finite widths: (at least) **one wide layer** with (N) neurons + all other layers with arbitrary sublinear widths



$$n_2 = \Omega(N)$$

$$d, n_1, n_3, n_4, n_5 \ll N$$

NTK with one wide layer

Theorem 4.1 (Nguyen, M. and Montufar, 2021)

Let $(W_l)_{i,j}$ i.i.d. $N(0, \frac{2}{d})$. Assume that the data points $f(x_i)g_{i=1}^N$ are i.i.d. with $k(x_i)$ and satisfying Lipschitz concentration. Let $n_l = 1$ if $n_l = (N \text{ polylog}(N))$ and $n_l = 0$, otherwise. Then, with high probability,

$$O \left(\sum_{l=1}^L \frac{1}{d} \frac{1}{n_l} \sum_{l=1}^L \frac{1}{l^2} \min_{l=1}^L K^L \right)$$

NTK with one wide layer

Theorem 4.1 (Nguyen, M. and Montufar, 2021)

Let $(W_l)_{i,j}$ i.i.d. $N(0, \frac{2}{d})$. Assume that the data points $f(x_i)g_{i=1}^N$ are i.i.d. with $k(x_i)$ and satisfying Lipschitz concentration. Let $n_l = 1$ if $n_l = (N \text{ polylog}(N))$ and $n_l = 0$, otherwise. Then, with high probability,

$$O \left(\prod_{l=1}^L \frac{1}{d} \frac{1}{n_l} \right) \min_{\mathbf{K}^L} \left(\prod_{l=1}^L \frac{1}{n_l} \right) \left(\prod_{l=1}^L \frac{1}{d} \frac{1}{n_l} \right)$$

\mathcal{G} layer with $\sim (N)$ neurons) **NTK bounded away from 0**

NTK with one wide layer

$$\mathcal{O} \left(\prod_{l=1}^L d_l n_l \right) \min \mathbf{K}^L$$

9 layer with $\sim (N)$ neurons) **NTK bounded away from 0**

Regardless of position of wide layer

Allow for arbitrary (polynomial) bottlenecks

NTK with one wide layer

$$\mathcal{O} \left(\sum_{l=1}^L d \sum_{i=1}^{n_l} \sum_{j=1}^{n_l} \sum_{k=1}^{n_l} \sum_{l=1}^L \min \mathbf{K}^L \right)$$

Immediate consequence on **memorization capacity**

Matching bounds if \mathcal{G} layer with $\sim (N)$ neurons + LeCun initialization $\left(\sum_{l=1}^L n_l = N \right)$

Numerical results

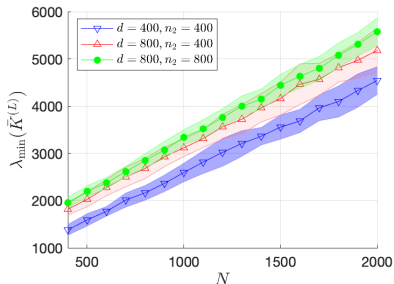
$$\min K^L = \max_{1 \leq l \leq L} n_l$$

$$L = 3$$

$$(W_l)_{ij} \sim N(0; 2/n_{l-1})$$

3 choices of $(d; n_2)$

$$n_1 = 8N$$



Smallest eigenvalue of NTK **scales linearly** in N (and in n_1)

