

Understanding Gradient Descent for Over-parameterized Deep Neural Networks

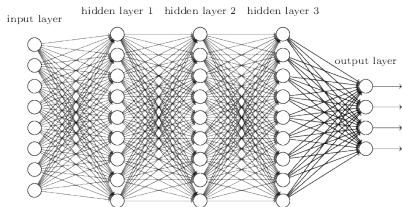
Marco Mondelli

Institute of Science and Technology Austria (IST Austria)

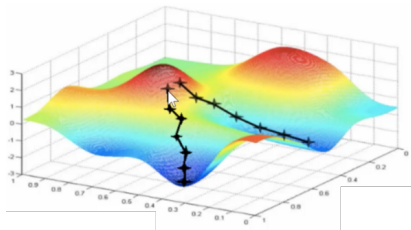
Youth in High-dimensions, 2 July 2020



Training a neural network is difficult (NP-hardness, local/disconnected minima...), but it works remarkably well!

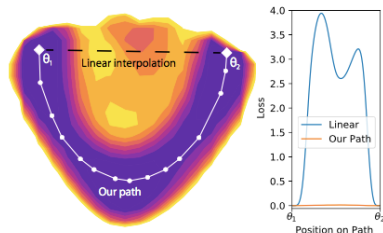


Over-parameterization



(Stochastic) gradient descent

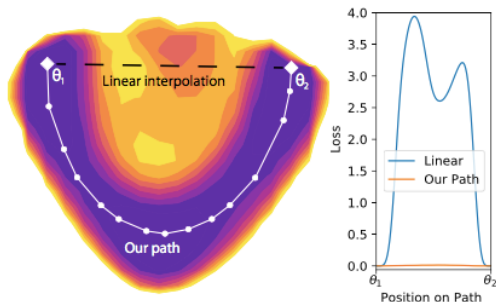
Landscape Connectivity and Dropout Stability



[SM20] A. Shevchenko and M. Mondelli, "Landscape Connectivity and Dropout Stability of SGD Solutions for Over-parameterized Neural Networks", *ICML*, 2020.

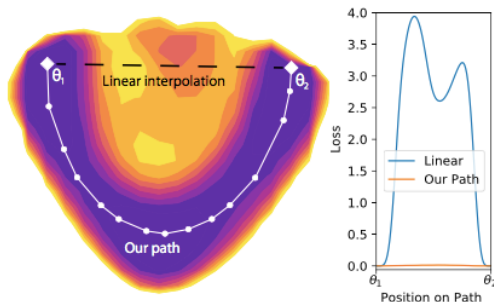
Landscape Connectivity

- SGD minima connected via piecewise linear path with constant loss [Garipov et al., 2018; Draxler et al., 2018]



Landscape Connectivity

- SGD minima connected via piecewise linear path with constant loss [Garipov et al., 2018; Draxler et al., 2018]



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

Landscape of Neural Networks

- Local minima are globally optimal for deep linear networks and networks with more neurons than training samples
- Connected sublevel sets for one wide layer and pyramidal topology
- Energy gap scaling exponentially with input dimension

Deep Learning without Poor Local Minima

Kenji Kawaguchi
Massachusetts Institute of Technology

Optimization Landscape and Expressivity of Deep CNNs

Quynh Nguyen¹ Matthias Hein²

On Connected Sublevel Sets in Deep Learning

Quynh Nguyen¹

Spurious Valleys in Two-layer Neural Network Optimization Landscapes

Luca Venturi^{*1}, Afonso S. Bandeira^{†1,2}, and Joan Bruna^{‡1,2}

The Loss Surface of Deep and Wide Neural Networks

Quynh Nguyen¹ Matthias Hein¹

ON THE LOSS LANDSCAPE OF A CLASS OF DEEP NEURAL NETWORKS WITH NO BAD LOCAL VALLEYS

Quynh Nguyen
Saarland University, Germany

Mahesh Chandra Mukkamala
Saarland University, Germany

TOPOLOGY AND GEOMETRY OF HALF-RECTIFIED NETWORK OPTIMIZATION

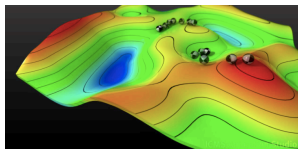
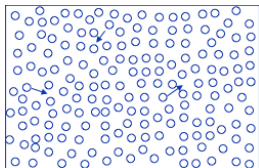
C. Daniel Freeman
Department of Physics,
UC Berkeley, Berkeley, CA 94720, USA
daniel.freeman@berkeley.edu
Joan Bruna^{*}
Courant Institute of Mathematical Sciences
New York University, New York, NY 10011, USA
bruna@cims.nyu.edu

Landscape of Neural Networks

- Local minima are globally optimal for deep linear networks and networks with more neurons than training samples
- Connected sublevel sets for one wide layer and pyramidal topology
- Energy gap scaling exponentially with input dimension

Strong assumptions on the model and **poor scaling** of parameters 😞

Mean Field View of Neural Networks



Discrete dynamics of SGD

Continuous dynamics of
gradient flow

- Two layers [Mei et al., 2018; Rotskoff et al., 2018; Chizat et al., 2018; Sirignano et al., 2018; ...]
- Multiple layers [Nguyen, 2019; Sirignano et al., 2019; Araujo et al., 2019]

Warm-Up: Two-Layer Networks

Data: $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\} \sim_{\text{i.i.d.}} \mathbb{P}(\mathbb{R}^d \times \mathbb{R})$

Warm-Up: Two-Layer Networks

Data: $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\} \sim_{\text{i.i.d.}} \mathbb{P}(\mathbb{R}^d \times \mathbb{R})$

Goal: Minimize loss $L_N(\boldsymbol{\theta}) = \mathbb{E} \left\{ \left(y - \frac{1}{N} \sum_{i=1}^N a_i \sigma(\mathbf{x}; \mathbf{w}_i) \right)^2 \right\}$, $\boldsymbol{\theta} = (\mathbf{w}, a)$

Warm-Up: Two-Layer Networks

Data: $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\} \sim_{\text{i.i.d.}} \mathbb{P}(\mathbb{R}^d \times \mathbb{R})$

Goal: Minimize loss $L_N(\boldsymbol{\theta}) = \mathbb{E} \left\{ \left(y - \frac{1}{N} \sum_{i=1}^N a_i \sigma(\mathbf{x}; \mathbf{w}_i) \right)^2 \right\}$, $\boldsymbol{\theta} = (\mathbf{w}, a)$

Online SGD: $\boldsymbol{\theta}^{k+1} = \boldsymbol{\theta}^k + \alpha N \nabla_{\boldsymbol{\theta}^k} \left(\left(y_k - \frac{1}{N} \sum_{i=1}^N a_i^k \sigma(\mathbf{x}_k; \mathbf{w}_i^k) \right)^2 \right)$

Warm-Up: Two-Layer Networks

Data: $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\} \sim_{\text{i.i.d.}} \mathbb{P}(\mathbb{R}^d \times \mathbb{R})$

Goal: Minimize loss $L_N(\boldsymbol{\theta}) = \mathbb{E} \left\{ \left(y - \frac{1}{N} \sum_{i=1}^N a_i \sigma(\mathbf{x}; \mathbf{w}_i) \right)^2 \right\}$, $\boldsymbol{\theta} = (\mathbf{w}, a)$

Online SGD: $\boldsymbol{\theta}^{k+1} = \boldsymbol{\theta}^k + \alpha N \nabla_{\boldsymbol{\theta}^k} \left(\left(y_k - \frac{1}{N} \sum_{i=1}^N a_i^k \sigma(\mathbf{x}_k; \mathbf{w}_i^k) \right)^2 \right)$

- y bounded, $\nabla_{\mathbf{w}} \sigma(\mathbf{x}; \mathbf{w})$ sub-gaussian
- σ bounded and differentiable, $\nabla \sigma$ bounded and Lipschitz
- initialization of a_i with bounded support

Dropout Stability

Dropout stability: loss does not change much if we remove part of neurons (and suitably rescale remaining neurons).

Dropout Stability

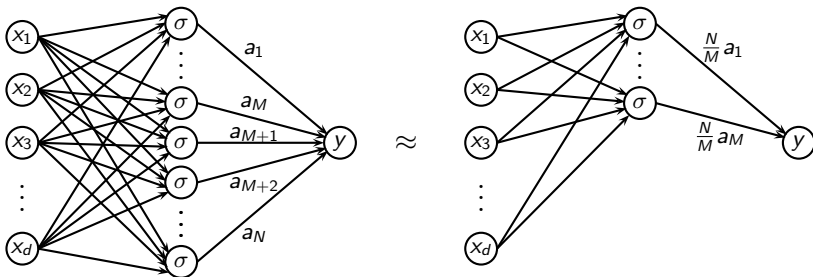
$$L_M(\boldsymbol{\theta}) = \mathbb{E} \left\{ \left(y - \frac{1}{M} \sum_{i=1}^M a_i \sigma(\mathbf{x}, \mathbf{w}_i) \right)^2 \right\}$$

$\boldsymbol{\theta}$ is ε_D -**dropout stable** if $|L_N(\boldsymbol{\theta}) - L_M(\boldsymbol{\theta})| \leq \varepsilon_D$

Dropout Stability

$$L_M(\theta) = \mathbb{E} \left\{ \left(y - \frac{1}{M} \sum_{i=1}^M a_i \sigma(\mathbf{x}, \mathbf{w}_i) \right)^2 \right\}$$

θ is ε_D -**dropout stable** if $|L_N(\theta) - L_M(\theta)| \leq \varepsilon_D$



Dropout Stability and Connectivity

$$L_M(\boldsymbol{\theta}) = \mathbb{E} \left\{ \left(y - \frac{1}{M} \sum_{i=1}^M a_i \sigma(\mathbf{x}, \mathbf{w}_i) \right)^2 \right\}$$

$\boldsymbol{\theta}$ is ε_D -**dropout stable** if $|L_N(\boldsymbol{\theta}) - L_M(\boldsymbol{\theta})| \leq \varepsilon_D$

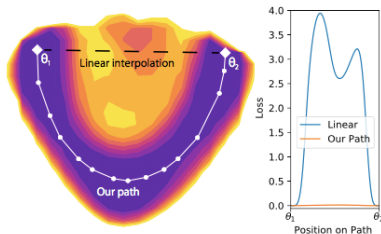
$\boldsymbol{\theta}$ and $\boldsymbol{\theta}'$ are ε_C -**connected** if there exists a continuous path connecting them where the loss does not increase more than ε_C

Dropout Stability and Connectivity

$$L_M(\theta) = \mathbb{E} \left\{ \left(y - \frac{1}{M} \sum_{i=1}^M a_i \sigma(\mathbf{x}, \mathbf{w}_i) \right)^2 \right\}$$

θ is ε_D -**dropout stable** if $|L_N(\theta) - L_M(\theta)| \leq \varepsilon_D$

θ and θ' are ε_C -**connected** if there exists a continuous path connecting them where the loss does not increase more than ε_C



Dropout Stability, Provably

- $N = \#$ neurons of full network
- $M = \#$ neurons after dropout
- $\alpha =$ step size of SGD
- $D =$ dimension of weights

Theorem [SM20]

Let θ^k be obtained after k SGD iterations. Then, with probability $1 - e^{-z^2}$, for all $k \in [T/\alpha]$, θ^k is ε_D -dropout stable with

$$\varepsilon_D = Ke^{KT^3} \left(\frac{\sqrt{\log M} + z}{\sqrt{M}} + \sqrt{\alpha}(\sqrt{D + \log N} + z) \right).$$

Dropout Stability, Provably

- $N = \#$ neurons of full network
- $M = \#$ neurons after dropout
- $\alpha =$ step size of SGD
- $D =$ dimension of weights

Theorem [SM20]

Let θ^k be obtained after k SGD iterations. Then, with probability $1 - e^{-z^2}$, for all $k \in [T/\alpha]$, θ^k is ε_D -dropout stable with

$$\varepsilon_D = Ke^{KT^3} \left(\frac{\sqrt{\log M} + z}{\sqrt{M}} + \sqrt{\alpha}(\sqrt{D + \log N} + z) \right).$$

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

Connectivity, Provably

Theorem [SM20]

Let θ^k be obtained after k SGD iterations using $\{(\mathbf{x}_j, y_j)\}_{j=0}^k \sim \mathbb{P}$, and $(\theta')^{k'}$ after k' SGD iterations using $\{(\mathbf{x}'_j, y'_j)\}_{j=0}^{k'} \sim \mathbb{P}$.

Connectivity, Provably

Theorem [SM20]

Let θ^k be obtained after k SGD iterations using $\{(\mathbf{x}_j, y_j)\}_{j=0}^k \sim \mathbb{P}$, and $(\theta')^{k'}$ after k' SGD iterations using $\{(\mathbf{x}'_j, y'_j)\}_{j=0}^{k'} \sim \mathbb{P}$. Then, with probability $1 - e^{-z^2}$, for all $k \in [T/\alpha]$ and $k' \in [T'/\alpha]$, θ^k and $(\theta')^{k'}$ are ε_C -connected with

$$\varepsilon_C = Ke^{K \max(T, T')^3} \left(\frac{\sqrt{\log N} + z}{\sqrt{N}} + \sqrt{\alpha}(\sqrt{D + \log N} + z) \right).$$

Connectivity, Provably

Theorem [SM20]

Let θ^k be obtained after k SGD iterations using $\{(\mathbf{x}_j, y_j)\}_{j=0}^k \sim \mathbb{P}$, and $(\theta')^{k'}$ after k' SGD iterations using $\{(\mathbf{x}'_j, y'_j)\}_{j=0}^{k'} \sim \mathbb{P}$. Then, with probability $1 - e^{-z^2}$, for all $k \in [T/\alpha]$ and $k' \in [T'/\alpha]$, θ^k and $(\theta')^{k'}$ are ε_C -connected with

$$\varepsilon_C = Ke^{K \max(T, T')^3} \left(\frac{\sqrt{\log N} + z}{\sqrt{N}} + \sqrt{\alpha}(\sqrt{D + \log N} + z) \right).$$

- Change in loss scales as $\sqrt{\frac{\log N}{N}} + \sqrt{\alpha(D + \log N)}$

Connectivity, Provably

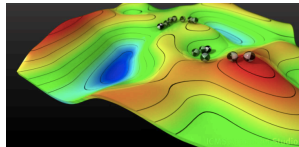
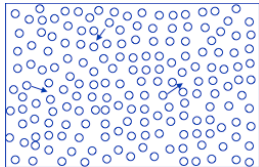
Theorem [SM20]

Let θ^k be obtained after k SGD iterations using $\{(\mathbf{x}_j, y_j)\}_{j=0}^k \sim \mathbb{P}$, and $(\theta')^{k'}$ after k' SGD iterations using $\{(\mathbf{x}'_j, y'_j)\}_{j=0}^{k'} \sim \mathbb{P}$. Then, with probability $1 - e^{-z^2}$, for all $k \in [T/\alpha]$ and $k' \in [T'/\alpha]$, θ^k and $(\theta')^{k'}$ are ε_C -connected with

$$\varepsilon_C = Ke^{K \max(T, T')^3} \left(\frac{\sqrt{\log N} + z}{\sqrt{N}} + \sqrt{\alpha}(\sqrt{D + \log N} + z) \right).$$

- Change in loss scales as $\sqrt{\frac{\log N}{N}} + \sqrt{\alpha(D + \log N)}$
- Can connect SGD solutions obtained from different training data (but same data distribution) and different initialization

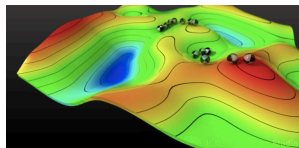
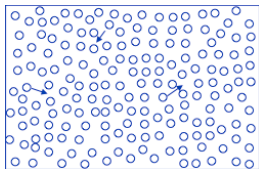
Proof Idea



Discrete dynamics of SGD

Continuous dynamics of
gradient flow

Proof Idea

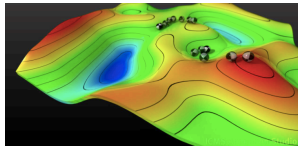
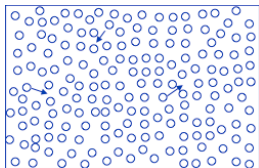


Discrete dynamics of SGD

Continuous dynamics of
gradient flow

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

Proof Idea

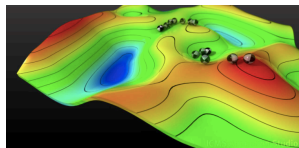
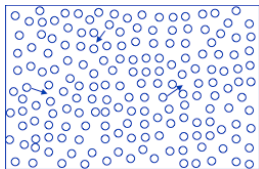


Discrete dynamics of SGD

Continuous dynamics of
gradient flow

- θ^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

Proof Idea



Discrete dynamics of SGD

Continuous dynamics of
gradient flow

- θ^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
- Dropout stability with $M = N/2 \Rightarrow$ connectivity

Multilayer Networks

- $N = \#$ neurons per layer of full network
- $M = \max. \#$ neurons per layer after dropout
- $\alpha = \text{step size of SGD}$
- $D = \max(d_x, d_y)$

Theorem [SM20]

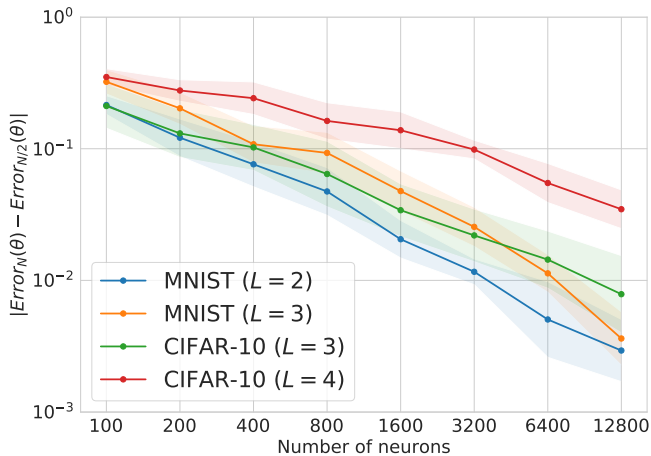
Let θ^k be obtained after k SGD iterations, with $k = T/\alpha$. Then, w. p. $1 - e^{-z^2}$, θ^k is ε_D -dropout stable with

$$\varepsilon_D = K(T, L) \left(\frac{\sqrt{D} + z}{\sqrt{M}} + \frac{\sqrt{\log N}}{\sqrt{N}} + \sqrt{\alpha}(\sqrt{D + \log N} + z) \right).$$

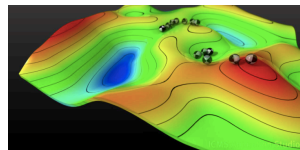
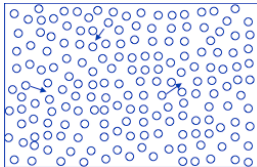
Let $(\theta')^{k'}$ be obtained after k' SGD iterations, with $k' = T'/\alpha$. Then, w. p. $1 - e^{-z^2}$, θ^k and $(\theta')^{k'}$ are ε_C -connected with

$$\varepsilon_C = K(T, T', L) \left(\frac{\sqrt{D + \log N} + z}{\sqrt{N}} + \sqrt{\alpha}(\sqrt{D + \log N} + z) \right).$$

Numerical Results



Global Convergence Without Rates



Discrete dynamics of SGD

Continuous dynamics of
gradient flow

A mean field view of the landscape of two-layer neural networks

Song Mei^a, Andrea Montanari^{b,c,1}, and Phan-Minh Nguyen^b

^aInstitute for Computational and Mathematical Engineering, Stanford University, Stanford, CA 94305; ^bDepartment of Electrical Engineering, Stanford University, Stanford, CA 94305; and ^cDepartment of Statistics, Stanford University, Stanford, CA 94305

Edited by Peter J. Bickel, University of California, Berkeley, CA, and approved June 21, 2018 (received for review April 16, 2018)

Multilayer neural networks are among the most powerful models in machine learning, yet the fundamental reasons for this success defy mathematical understanding. Learning a neural network requires optimizing a nonconvex high-dimensional objective (risk function), a problem that is usually attacked using stochastic gra-

phic descent. In the present case, this amounts to the iteration

$$\theta_i^{t+1} = \theta_i^t + 2\eta_t (y_i - \hat{y}(x_i; \theta^t)) \nabla_{\theta_i} \sigma(x_i; \theta_i^t). \quad [3]$$

On the Global Convergence of Gradient Descent for Over-parameterized Models using Optimal Transport

Lénaïc Chizat
INRIA, ENS, PSL Research University
Paris, France
lenaic.chizat@inria.fr

Francis Bach
INRIA, ENS, PSL Research University
Paris, France
francis.bach@inria.fr

Abstract

Many tasks in machine learning and signal processing can be solved by minimizing a convex function of a measure. This includes sparse spikes deconvolution or training a neural network with a single hidden layer. For these problems, we study a simple minimization method: the unknown measure is discretized into a mixture of particles and a continuous-time gradient descent is performed on their weights

Data Model

Data: $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\} \sim_{\text{i.i.d.}} \mathbb{P}(\mathbb{R}^d \times \mathbb{R})$

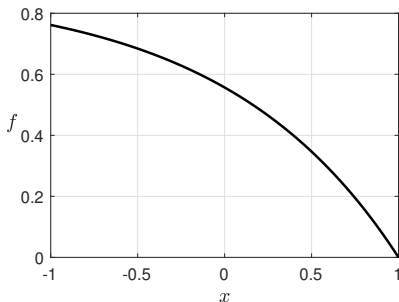
$\mathbf{x}_i \sim \text{Unif}(\Omega)$ and $y_i = f(\mathbf{x}_i)$

Data Model

Data: $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\} \sim_{\text{i.i.d.}} \mathbb{P}(\mathbb{R}^d \times \mathbb{R})$

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

- Ω bounded and convex
- f positive, smooth and γ -strongly concave



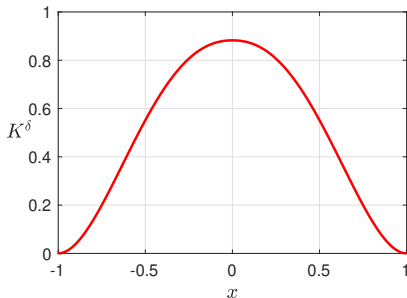
Problem Statement

Goal: Minimize loss $L_N(\mathbf{w}) = \mathbb{E}\left\{\left(y - \frac{1}{N} \sum_{i=1}^N \sigma(\mathbf{x}; \mathbf{w}_i)\right)^2\right\}$

Problem Statement

Goal: Minimize loss $L_N(\mathbf{w}) = \mathbb{E}\left\{\left(y - \frac{1}{N} \sum_{i=1}^N \sigma(\mathbf{x}; \mathbf{w}_i)\right)^2\right\}$

- σ 'bump'-like
 $\sigma(\mathbf{x}; \mathbf{w}_i) = K^\delta(\mathbf{x} - \mathbf{w}_i)$
- \mathbf{w}_i = center of bump
- δ = width of the bump
 $K^\delta(\mathbf{x}) = \delta^{-d} K(\mathbf{x}/\delta)$



Exponential and Dimension-free Convergence Rate

Theorem [JMM20]

Let f be γ -strongly concave, σ a bump of width δ , and \mathbf{w}^k obtained after k SGD iterations with step size α . Then, with high probability,

$$\underbrace{L_N(\mathbf{w}^k)}_{\text{loss at step } k} \leq L_N(\mathbf{w}^0) \underbrace{e^{-2\gamma k\alpha}}_{\text{exp. rate}} + \underbrace{\Delta(N, \alpha, \delta)}_{\text{distance to opt.}},$$

where $\Delta(N, \alpha, \delta) \rightarrow 0$ as $N \rightarrow \infty$, $\alpha \rightarrow 0$, and $\delta \rightarrow 0$.

- As $\delta \rightarrow 0$, gradient flow optimizes a **displacement convex** loss.

[JMM20] A. Javanmard, M. Mondelli, and A. Montanari, "Analysis of a Two-Layer Neural Network via Displacement Convexity", *Annals of Statistics*, 2020+.

Thank You for Your Attention

