



The Condensation Phenomenon of Deep Neural Networks

Yaoyu Zhang

Institute of Natural Sciences & School of Mathematical Sciences Shanghai Jiao Tong University MLPDES25 workshop, FAU MoD

Learning systems with increasingly large size



Suzana Herculano-Houzel, 2009

Parameters of transformer-based language models



62023 TECHTARGET. ALL RIGHTS RESERVED TechTarget



Failure of traditional wisdom

Large complexity → Large generalization gap



Traditional wisdom: complex models easily overfit



Generalization Gap

Long-standing problems





Leo Breiman Statistics Department, University of California, Berkeley, CA 94305; e-mail: leo@stat.berkeley.edu

Reflections After Refereeing Papers for NIPS

- Why don't heavily parameterized neural networks overfit the data?
- What is the effective number of parameters?
- Why doesn't backpropagation head for a poor local minima?
- When should one stop the backpropagation and use the current parameters?

How (overparameterized) neural networks control the complexity of output function during **nonlinear** training?



Condensation Phenomenon



Illustration of Condensation





Tao Luo, Zhi-Qin John Xu, Zheng Ma, Yaoyu Zhang, Phase diagram for two-layer ReLU neural networks at infinite-width limit, Journal of Machine Learning Research (2021)

1d example: condensation with small initialization









Small initialization: $a_j(0), w_j(0), b_j(0) \sim N(0, \sigma^2)$ with small σ



Evolution trajectory: change significantly









(a) epoch=100

(b) epoch=1000

(c) epoch=3000



Evolution trajectory: change significantly









(d) epoch=5000

(e) epoch=10000

(f) epoch=100000



Condensation in CNN on MNIST



Cosine similarity: $D(u_1, u_2) = \frac{u_1^{\mathsf{T}} u_2}{(u_1^{\mathsf{T}} u_1)^{1/2} (u_2^{\mathsf{T}} u_2)^{1/2}}.$

100% training and 97.62% test accuracy





Condensation in transformer



Zhi-Qin John Xu, Yaoyu Zhang, Zhangchen Zhou, "An overview of condensation phenomenon in deep learning," arXiv:2504.09484 (2025).



Regime of Condensation

1.Tao Luo, Zhi-Qin John Xu, Zheng Ma, Yaoyu Zhang, "Phase Diagram for Two-layer ReLU Neural Networks at Infinite-Width Limit," Journal of Machine Learning Research (JMLR) 22(71):1–47, (2021).

2.Hanxu Zhou, Qixuan Zhou, Zhenyuan Jin, Tao Luo, Yaoyu Zhang, Zhi-Qin John Xu, "Empirical Phase Diagram for Three-layer Neural Networks with Infinite Width," NeurIPS 2022.



Normalization and scaling parameters



$$f^{\alpha}_{\boldsymbol{\theta}}(\boldsymbol{x}) = \frac{1}{\alpha} \sum_{k=1}^{m} a_k \sigma(\boldsymbol{w}_k^{\mathsf{T}} \boldsymbol{x}) \qquad a^0_k \sim N(0, \beta_1^2), \ \boldsymbol{w}^0_k \sim N(0, \beta_2^2 \boldsymbol{I}_d) \qquad \begin{array}{l} \boldsymbol{x} = [\boldsymbol{x}^T, 1]^T \\ \boldsymbol{w}_k = [\boldsymbol{w}^T_k, \boldsymbol{b}_k]^T \end{array}$$

Normalized gradient flow

$$\bar{a}_{k} = \beta_{1}^{-1} a_{k}, \quad \bar{\boldsymbol{w}}_{k} = \beta_{2}^{-1} \boldsymbol{w}_{k}, \quad \bar{t} = \frac{1}{\beta_{1}\beta_{2}} t,$$

$$\frac{\mathrm{d}\bar{a}_{k}}{\mathrm{d}\bar{t}} = -\frac{1}{\kappa'} \frac{1}{n} \sum_{i=1}^{n} \kappa \sigma(\bar{\boldsymbol{w}}_{k}^{\mathsf{T}} \boldsymbol{x}_{i}) \left(\kappa \sum_{k'=1}^{m} \bar{a}_{k'} \sigma(\bar{\boldsymbol{w}}_{k'}^{\mathsf{T}} \boldsymbol{x}_{i}) - y_{i}\right),$$

$$\frac{\mathrm{d}\bar{\boldsymbol{w}}_{k}}{\mathrm{d}\bar{t}} = -\kappa' \frac{1}{n} \sum_{i=1}^{n} \kappa \bar{a}_{k} \sigma'(\bar{\boldsymbol{w}}_{k}^{\mathsf{T}} \boldsymbol{x}_{i}) \boldsymbol{x}_{i} \left(\kappa \sum_{k'=1}^{m} \bar{a}_{k'} \sigma(\bar{\boldsymbol{w}}_{k'}^{\mathsf{T}} \boldsymbol{x}_{i}) - y_{i}\right).$$

$$m \to +\infty$$
$$\frac{\beta_1 \beta_2}{\alpha} = m^{-\gamma}$$
$$\frac{\beta_1}{\beta_2} = m^{-\gamma'}$$

Scaling parameters and infinite-width limit

$$\kappa := \frac{\beta_1 \beta_2}{\alpha}, \quad \kappa' := \frac{\beta_1}{\beta_2}, \quad \gamma = \lim_{m \to \infty} -\frac{\log \kappa}{\log m}, \quad \gamma' = \lim_{m \to \infty} -\frac{\log \kappa'}{\log m}$$



Tao Luo, Zhi-Qin John Xu, Zheng Ma, Yaoyu Zhang, Phase diagram for two-layer ReLU neural networks at infinite-width limit, Journal of Machine Learning Research (2021)

Regime separation -- theorems

Small γ (large init) \rightarrow Linear/NTK regime

Theorem 1*. (Informal statement of Theorem 6) If $\gamma < 1$ or $\gamma' > \gamma - 1$, then with a high probability over the choice of θ^0 , we have

$$\lim_{m \to +\infty} \sup_{t \in [0, +\infty)} \operatorname{RD}(\boldsymbol{\theta}_{\boldsymbol{w}}(t)) = 0. \qquad f_{\boldsymbol{\theta}}^{\alpha}(\boldsymbol{x}) = \frac{1}{\alpha} \sum_{k=1}^{m} a_{k} \sigma(\boldsymbol{w}_{k}^{\mathsf{T}} \boldsymbol{x})$$
$$\boldsymbol{\theta}_{\boldsymbol{w}} = \operatorname{vec}(\{\boldsymbol{w}_{k}\}_{k=1}^{m})$$
$$\operatorname{RD}(\boldsymbol{\theta}_{\boldsymbol{w}}(t)) = \frac{\|\boldsymbol{\theta}_{\boldsymbol{w}}(t) - \boldsymbol{\theta}_{\boldsymbol{w}}(0)\|_{2}}{\|\boldsymbol{\theta}_{\boldsymbol{w}}(0)\|_{2}}.$$

Large γ (small init) \rightarrow Condensed regime

Theorem 2*. (Informal statement of Theorem 8) If $\gamma > 1$ and $\gamma' < \gamma - 1$, then with a high probability over the choice of θ^0 , we have

$$\lim_{m \to +\infty} \sup_{t \in [0, +\infty)} \operatorname{RD}(\boldsymbol{\theta}_{\boldsymbol{w}}(t)) = +\infty.$$
(21)



Tao Luo, Zhi-Qin John Xu, Zheng Ma, Yaoyu Zhang, Phase diagram for two-layer ReLU neural networks at infinite-width limit, Journal of Machine Learning Research (2021)

Initialization scheme



Name (related works)	α	eta_1	eta_2	$rac{\kappa}{\left(rac{eta_1eta_2}{lpha} ight)}$	$\kappa' \ \left(rac{eta_1}{eta_2} ight)$	$\gamma_{\left(\lim_{m\to\infty}\frac{\log 1/\kappa}{\log m}\right)}$	$\gamma' \ (\lim_{m \to \infty} rac{\log 1/\kappa'}{\log m})$
LeCun (LeCun et al., 2012)	1	$\sqrt{\frac{1}{m}}$	$\sqrt{\frac{1}{d}}$	$\sqrt{rac{1}{md}}$	$\sqrt{rac{d}{m}}$	$\frac{1}{2}$	$\frac{1}{2}$
He (He et al., 2015)	1	$\sqrt{\frac{2}{m}}$	$\sqrt{\frac{2}{d}}$	$\sqrt{rac{4}{md}}$	$\sqrt{rac{d}{m}}$	$\frac{1}{2}$	$\frac{1}{2}$
Xavier (Glorot and Bengio, 2010)	1	$\sqrt{\frac{2}{m+1}}$	$\sqrt{\frac{2}{m+d}}$	$\sqrt{\frac{4}{(m+1)(m+d)}}$	$\sqrt{\frac{m+d}{m+1}}$	1	0
NTK (Jacot et al., 2018)	\sqrt{m}	1	1	$\sqrt{\frac{1}{m}}$	1	$\frac{1}{2}$	0
Mean-field (Mei et al., 2018) (Sirignano and Spiliopoulos, 2020)	m	1	1	$\frac{1}{m}$	1	1	0
(Rotskoff and Vanden-Eijnden, 2018) E et al. (E et al., 2020)	1	eta	1	eta	eta	$\lim_{m \to \infty} \frac{\log 1/\beta}{\log m}$	$\lim_{m \to \infty} \frac{\log 1/\beta}{\log m}$





When condensation happens (at infinite width limit)?

Phase Diagram





Tao Luo, Zhi-Qin John Xu, Zheng Ma, Yaoyu Zhang, Phase diagram for two-layer ReLU neural networks at infinite-width limit, Journal of Machine Learning Research (2021)

Feature distribution across the phase diagram



Blue: $m = 10^3$ Red: $m = 10^4$ Yellow: $m = 10^6$



Tao Luo, Zhi-Qin John Xu, Zheng Ma, Yaoyu Zhang, Phase diagram for two-layer ReLU neural networks at infinite-width limit, Journal of Machine Learning Research (2021)

Loss landscape structure underlying condensation

1.Yaoyu Zhang, Zhongwang Zhang, Tao Luo, Zhi-Qin John Xu, "Embedding Principle of Loss Landscape of Deep Neural Networks," NeurIPS 2021 spotlight.

2.Yaoyu Zhang, Yuqing Li, Zhongwang Zhang, Tao Luo, Zhi-Qin John Xu, "Embedding Principle: a hierarchical structure of loss landscape of deep neural networks," Journal of Machine Learning, 1(1), pp. 60-113, 2022.



Typical training behavior (small init)



Width-500 tanh-NN (~1500 parameters)



Trajectory of training loss







Intermediate condensation







Condensed critical points for intermediate stage _____



Embedding Principle (informal Theorem) The loss landscape of any network ``contains" all critical points of all narrower networks.

Equivalent Statement $\mathcal{F}_{narr}^{c} \subseteq \mathcal{F}_{wide}^{c}$, where $\mathcal{F}^{c} \coloneqq \{f_{\theta}(\cdot) | \nabla R_{S}(\theta) = 0\}$.

Observation: Width similarity

Implication of theory: simple condensed critical points are common



Example: identification of critical points and functions

500 tanh neuron





Embedding principle

One-step splitting embedding $T: \mathbb{R}^{M_{\text{narr}}} \to \mathbb{R}^{M_{\text{wide}}}$



Theorem: One-step splitting embedding *T* with $\theta_{wide} = T(\theta_{narr})$ satisfies: (i) **output preserving**: $f_{\theta_{narr}}(x) = f_{\theta_{wide}}(x)$; (ii) **criticality preserving**: If $\nabla R_S(\theta_{narr}) = \mathbf{0}$, then $\nabla R_S(\theta_{wide}) = \mathbf{0}$.



Existance of condensed critical points---embedding principle





Generalization advantage of condensation

1. Yaoyu Zhang, Zhongwang Zhang, Leyang Zhang, Zhiwei Bai, Tao Luo, Zhi-Qin John Xu, Linear Stability Hypothesis and Rank Stratification for Nonlinear Models. arXiv:2211.11623, (2022).

2. Yaoyu Zhang, Zhongwang Zhang, Leyang Zhang, Zhiwei Bai, Tao Luo, Zhi-Qin John Xu, Optimistic Estimate Uncovers the Potential of Nonlinear Models. arXiv:2307.08921, (2023).

3. Yaoyu Zhang, Leyang Zhang, Zhongwang Zhang and Zhiwei Bai, Local Linear Recovery Guarantee of Deep Neural Networks at Overparameterization. JMLR (2025) Accepted

Generalization consequence of condensation

Large initialization (no condensation)



Small initialization (Strong condensation)





How many samples are required to recover *f**?





Model:

$$F:\mathbb{R}^M\to \boldsymbol{\mathcal{F}}\subset C(R^d)$$

Model rank:

$$r_{\boldsymbol{\theta}} = \dim \operatorname{span} \left\{ \partial_{\theta_i} F(\boldsymbol{\theta})(\cdot) \right\}_{i=1}^{M}$$

Optimistic sample size (
$$f^* \in \mathcal{F}$$
) :
 $O_{f^*} = \min_{\theta \in F^{-1}(f^*)} r_{\theta}$ $F^{-1}(f^*)$: Target set

Intuitive procedure:



Yaoyu Zhang, Leyang Zhang, Zhongwang Zhang and Zhiwei Bai,

Local Linear Recovery Guarantee of Deep Neural Networks at Overparameterization. JMLR (2025) Accepted



Optimistic estimate vs. experiment

Theorem 5 (optimistic sample sizes for two-layer tanh-NN). Given a two-layer NN $f_{\theta}(x) = \sum_{i=1}^{m} a_i \tanh(w_i^T x), x \in \mathbb{R}^d, \theta = (a_i, w_i)_{i=1}^m$, for any target function $f^* \in \mathcal{F}_k^{\text{NN}} \setminus \mathcal{F}_{k-1}^{\text{NN}}$ with $0 \leq k \leq m$, the optimistic sample size



Yaoyu Zhang, Leyang Zhang, Zhongwang Zhang and Zhiwei Bai,

Local Linear Recovery Guarantee of Deep Neural Networks at Overparameterization. JMLR (2025) Accepted





Yaoyu Zhang, Leyang Zhang, Zhongwang Zhang and Zhiwei Bai, Local Linear Recovery Guarantee of Deep Neural Networks at Overparameterization. JMLR (2025) Accepted





Adding (unnecessary) connections reduces sample efficiency

MNIST *k*-kernel, kernel size: 3x3

1000X worse!

CNN (no sharing) : 6760k samples
FNN: 530660k samples

Yaoyu Zhang, Leyang Zhang, Zhongwang Zhang and Zhiwei Bai,

Local Linear Recovery Guarantee of Deep Neural Networks at Overparameterization. JMLR (2025) Accepted

 \succ CNN: 685k samples



Condensation—hen with golden eggs







FAU MoD Course



Friedrich-Alexander-Universität Research Center for Mathematics of Data | MoD

FAU MoD Course



Towards a mathematical foundation of Deep Learning: From phenomena to theory

Yaoyu Zhang

SHANGHAI JIAO TONG UNIVERSITY



WHEN Fri.-Thu. May 2-8, 2025 10:00H (Berlin time)

WHERE On-site / Online

Friedrich-Alexander-Universität Erlangen-Nürnberg (FAU) Room H11 / H16 Felix-Klein building Cauerstraße 11, 91058 Erlangen. Bavaria, Germany

Live-streaming: www.fau.tv/fau-mod-livestream-2025

*Check room/day on website

Establishing a mathematical foundation for deep learning is a significant and challenging endeavor in mathematics. Recent theoretical advancements are transforming deep learning from a black box into a more transparent and understandable framework. This course offers an in-depth exploration of these developments, emphasizing a promising phenomenological approach. It is designed for those seeking an intuitive understanding of how neural networks learn from data, as well as an appreciation of their theoretical underpinnings. (...)

Session Titles: 1. Mysteries of Deep Learning 2. Frequency Principle/Spectral Bias 3. Condensation Phenomenon 4. From Condensation to Loss Landscape Analysis 5. From Condensation to Generalization Theory

Overall, this course serves as a gateway to the vibrant field of deep learning theory, inspiring participants to contribute fresh perspectives to its advancement and application.

Towards a Mathematical Foundation of Deep Learning: From Phenomena to Theory

Date Fri. – Thu. May 2 – 8, 2025

Session Titles

- 1. Mysteries of Deep Learning
- 2. Frequency Principle/Spectral Bias
- 3. Condensation Phenomenon
- 4. From Condensation to Loss Landscape Analysis
- 5. From Condensation to Generalization Theory





Thanks!