



III. Condensation phenomenon

Yaoyu Zhang

Institute of Natural Sciences & School of Mathematical Sciences

Shanghai Jiao Tong University

FAU MoD Course 饮水思源 • 爱国荣校

Deep learning is no longer a black-box



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



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

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

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

Date

Fri. – Thu. May 2 – 8, 2025

Session Titles

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



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.

Illustration of Condensation





Initial:random

Training: condense

Effect: equiv to small net



1d example: condensation with small initialization

train loss

105

104

9



 10^{-1}

10-2

10-3

 10^{-4}

100

10¹

10²

epoch

10³

g2u

0

Test

True

-1

1.0

0.8

0.6

0.4

0.2

0.0



270°

 \mathcal{W}

Small initialization: $a_i(0), w_i(0), b_i(0) \sim N(0, \sigma^2)$ with small σ



Evolution trajectory: change significantly





Evolution trajectory: change significantly







Condensation in CNN on MNIST





(a) Loss





10

5

15

index

(e) final weight

20

25

30

-1.00

0+0 0

Cosine similarity:

$$D(\boldsymbol{u}_1, \boldsymbol{u}_2) = \frac{\boldsymbol{u}_1^{\mathsf{T}} \boldsymbol{u}_2}{(\boldsymbol{u}_1^{\mathsf{T}} \boldsymbol{u}_1)^{1/2} (\boldsymbol{u}_2^{\mathsf{T}} \boldsymbol{u}_2)^{1/2}}.$$

100% training and 97.62% test accuracy



Condensation in transformer





$$A_{ heta}(X) = \sum_{i=1}^{h} \operatorname{softmax}_{\operatorname{row}} \left(rac{XW_{Q_i}W_{K_i}^{ op}X^{ op}}{\sqrt{d}}
ight) XW_{V_i}W_{O_i}^{ op}$$

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.



- Data: $\left\{x_i \in \mathbb{R}^d, y_i \in \mathbb{R}\right\}_{i=1}^n$
- Two layer ReLU network

$$f^{\alpha}_{\boldsymbol{\theta}}(\boldsymbol{x}) = \frac{1}{\alpha} \sum_{k=1}^{m} a_k \sigma(\boldsymbol{w}_k^{\mathsf{T}} \boldsymbol{x})$$

$$a_k^0 \sim N(0, \beta_1^2), \ \boldsymbol{w}_k^0 \sim N(0, \beta_2^2 \boldsymbol{I}_d)$$

• Loss

$$R_S(\boldsymbol{\theta}) = rac{1}{2n} \sum_{i=1}^n (f^{\alpha}_{\boldsymbol{\theta}}(\boldsymbol{x}_i) - y_i)^2$$

Gradient flow dynamics

$$\frac{\mathrm{d}\boldsymbol{\theta}}{\mathrm{d}t} = -\nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}).$$

Overparameterized setup:

 $M=m(d+1)\gg n,$

Properties:

 Global minima is M-n dimensional (proved by Yaim Cooper 2018)
 Often non-overfitting
 Evalution of O(t) and O(t)

 $x = [x^T, 1]^T$

 $w_k = \begin{bmatrix} w_k^T, b_k \end{bmatrix}^T$

3. Evolution of $\theta(t)$ and $\theta(\infty)$ depend on α, β_1, β_2

Goal:

Identify dynamical regimes of training over α , β_1 , β_2 at infinite-width limit.

Name (related works)	α	eta_1	β_2	$rac{\kappa}{\left(rac{eta_1eta_2}{lpha} ight)}$	$\kappa' \ (rac{eta_1}{eta_2})$	$\gamma \ \left(\lim_{m o \infty} rac{\log 1/\kappa}{\log m} ight)$	$\gamma' \ (\lim_{m o \infty} rac{\log 1/\kappa'}{\log m}$
LeCun (LeCun et al., 2012)	1	$\sqrt{\frac{1}{m}}$	$\sqrt{\frac{1}{d}}$	$\sqrt{rac{1}{md}}$	$\sqrt{\frac{d}{m}}$	$\frac{1}{2}$	$\frac{1}{2}$
He (He et al., 2015)	1	$\sqrt{\frac{2}{m}}$	$\sqrt{\frac{2}{d}}$	$\sqrt{rac{4}{md}}$	$\sqrt{\frac{d}{m}}$	$\frac{1}{2}$	$\frac{1}{2}$
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	β	$\lim_{m \to \infty} \frac{\log 1/\beta}{\log m}$	$\lim_{m \to \infty} \frac{\log 1/\beta}{\log m}$

Normalization and scaling parameters

• Two layer ReLU network

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

17

Normalized gradient flow

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

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

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

Scaling parameters and infinite-width limit

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



Phase diagram





Phase diagram for matter distinctive states of matter <-> environment (phase transition happens at infinite size limit) solid, liquid, gas <-> pressure, temperature

• Phase diagram for two-layer ReLU NN training dynamics <-> initialization $(m \rightarrow \infty)$? <->?

 $\gamma = \lim_{m o \infty} - rac{\log eta_1 eta_2 / lpha}{\log m}, \;\; \gamma' = \lim_{m o \infty} - rac{\log eta_1 / eta_2}{\log m}$

Identification of coordinates of phase diagram (in analogy to pressure, temperature)

- 1. Effectively independent
- 2. Dynamical similarity
- 3. Differentiation capability







Typical cases across the phase diagram





• Linear regime (with ASI)

$$f_{\boldsymbol{\theta}}^{\text{lin}} = \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}(0)} \cdot (\boldsymbol{\theta}(t) - \boldsymbol{\theta}(0)).$$

Relative distance

$$\operatorname{RD}(\boldsymbol{\theta}_{\boldsymbol{w}}(t)) = \frac{\|\boldsymbol{\theta}_{\boldsymbol{w}}(t) - \boldsymbol{\theta}_{\boldsymbol{w}}(0)\|_{2}}{\|\boldsymbol{\theta}_{\boldsymbol{w}}(0)\|_{2}}.$$

$$f^{\alpha}_{\boldsymbol{\theta}}(\boldsymbol{x}) = \frac{1}{\alpha} \sum_{k=1}^{m} a_k \sigma(\boldsymbol{w}_k^{\mathsf{T}} \boldsymbol{x})$$

$$\boldsymbol{\theta}_{\boldsymbol{w}} = \operatorname{vec}(\{\boldsymbol{w}_k\}_{k=1}^m)$$

As $m \to \infty$,

• Linear regime: $\sup_{t \in [0, +\infty)} \operatorname{RD}(\boldsymbol{\theta}_{\boldsymbol{w}}(t)) \to 0$

- Condensed regime: $\sup_{t \in [0,+\infty)} \operatorname{RD}(\boldsymbol{\theta}_{\boldsymbol{w}}(t)) \rightarrow +\infty$
- Critical regime:

$$\sup_{t \in [0,+\infty)} \operatorname{RD}(\boldsymbol{\theta}_{\boldsymbol{w}}(t)) \to O(1),$$





• Two layer ReLU network at infinite-width limit

$$\begin{split} f^{\alpha}_{\boldsymbol{\theta}}(\boldsymbol{x}) = & \frac{1}{\alpha} \sum_{k=1}^{m} a_{k} \sigma(\boldsymbol{w}_{k}^{\mathsf{T}} \boldsymbol{x}) \qquad a^{0}_{k} \sim N(0, \beta_{1}^{2}), \ \boldsymbol{w}_{k}^{0} \sim N(0, \beta_{2}^{2} \boldsymbol{I}_{d}) \qquad \stackrel{\boldsymbol{x} = \begin{bmatrix} \boldsymbol{x}^{T}, 1 \end{bmatrix}^{T}}{\boldsymbol{w}_{k} = \begin{bmatrix} \boldsymbol{w}_{k}^{T}, \boldsymbol{b}_{k} \end{bmatrix}^{T}} \\ \kappa := & \frac{\beta_{1} \beta_{2}}{\alpha}, \quad \kappa' := \frac{\beta_{1}}{\beta_{2}}, \end{split}$$

• "capability" of NN: $C = m\beta_1\beta_2/\alpha = m\kappa \gtrsim O(1)$

$$\kappa\gtrsim 1/m$$

• output-layer dominant: $C = m\beta_2 \mathbb{E}(|a|)/\alpha \ll m\beta_2^2/\alpha = m\kappa/\kappa'$.

$$1/\kappa' \gg 1/m\kappa$$



Regime identification through experiments









 $A = |a| \|\boldsymbol{w}\|_2$





Blue: $m = 10^3$ red: $m = 10^4$ **Yellow:** $m = 10^{6}$









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

Regime separation -- theorems

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

$$\lim_{m \to +\infty} \sup_{t \in [0, +\infty)} \operatorname{RD}(\boldsymbol{\theta}_{\boldsymbol{w}}(t)) = 0.$$
(20)

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)







Typical cases across the phase diagram









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

Condensation facilitates reasoning



Zhongwang Zhang, Pengxiao Lin, Zhiwei Wang, Yaoyu Zhang, Zhi-Qin John Xu, "Initialization is Critical to Whether Transformers Fit Composite Functions by Inference or Memorizing," NeurIPS 2024.











Mechanism 1: learn symmetric structure









Phase diagram of symmetric solution



Phase diagram of inferential solution



Condensation of $W^{Q(1)}$ by column







Mechanism 1: learn symmetric structure



Not only one pair But ten pairs in training 11;12,21; 13, 31; 14, 41; 23, 32;... Mechanism 2: infer single anchor mappings



Need to learn four functions









36



Symmetric solution





Inferential solution













Condensation





See more works on my personal website: https://yaoyuzhang1.github.io/





How can condensation be facilitated in a neural network?

Is it valid to compare the performance of wide and narrow networks when the initialization variance is fixed?

What initialization strategy can be used for a three-layer network to induce condensation in the first hidden layer but not in the second?

Where condensation can happen within a transformer?





Thanks!

