



The Condensation Phenomenon of Deep Neural Networks

Yaoyu Zhang

Institute of Natural Sciences & School of Mathematical Sciences

Shanghai Jiao Tong University

MaD Seminar, New York University

饮水思源·爱国荣校

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

30

0

Ó

10

5

15

index

(e) final weight

20

25

30



(a) Loss





- 0.75

-1.00

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



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



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.



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)

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)

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



Feature distribution across the phase diagram



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



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



Typical cases across the 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)

Phase diagram in three-layer ReLU NN



 上海交通大学 Shanghai Jiao Tong UNIVERSITY

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

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.

3.Hanxu Zhou, Qixuan Zhou, Tao Luo, Yaoyu Zhang, Zhi-Qin John Xu, "Towards Understanding the Condensation of Neural Networks at Initial Training," NeurIPS 2022.

Typical training behavior with strong condensation



Width-500 tanh-NN (~1500 parameters)



Trajectory of training loss







Initial condensation



Hanxu Zhou, Qixuan Zhou, Tao Luo, Yaoyu Zhang, Zhi-Qin John Xu, Towards Understanding the Condensation of Neural Networks at Initial Training, NeurIPS 2022.



Loss landscape around 0 and Initial condensation

$$\dot{\boldsymbol{w}}_{j} = \sum_{i=1}^{m} (y_{i} - f_{\boldsymbol{\theta}}(\boldsymbol{x}_{i})) a_{j} \sigma' (\boldsymbol{w}_{j}^{\mathrm{T}} \boldsymbol{x}_{i}) \boldsymbol{x}_{i}$$

When $\theta \approx 0$, then $f_{\theta}(\cdot) \approx 0(\cdot)$:

$$\dot{\mathbf{w}}_j \approx a_j \sum_{i=1}^m y_i \sigma'(\mathbf{w}_j^{\mathrm{T}} \mathbf{x}_i) \mathbf{x}_i$$

If $\sigma'(0) \neq 0$ (e.g. tanh, swish, gelu): $\dot{w}_j \approx a_j \sigma'(0) \sum_{i=1}^m y_i x_i$

i. No coupling between w_i and $w_{i'}$! ii. 2 limiting directions: $\pm \sum_{i=1}^{m} y_i x_i$.



(a) tanh(x)



49

(b) $x \tanh(x)$







Hanxu Zhou, Qixuan Zhou, Tao Luo, Yaoyu Zhang, Zhi-Qin John Xu, Towards Understanding the Condensation of Neural Networks at Initial Training, NeurIPS 2022.

Intermediate condensation





Zhang, Zhang, Luo, Xu, NeurIPS 2021 spotlight. Zhang, Li, Zhang, Luo, Xu, Journal of Machine Learning 2022.



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



Zhang, Zhang, Luo, Xu, NeurIPS 2021 spotlight. Zhang, Li, Zhang, Luo, Xu, Journal of Machine Learning 2022.

hierarchical structure of DNN loss landscape



Simple \rightarrow Complex





Zhang, Li, Zhang, Luo, Xu, Journal of Machine Learning 2022.

Example: identification of critical points and functions

500 tanh neuron





Zhang, Li, Zhang, Luo, Xu, Journal of Machine Learning 2022.

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





Final condensation





Geometry of global-min: simpler f*, higher-dim Q*_____

- Model: $F(\theta)(x) = a_1 \sigma(w_1^T x) + a_2 \sigma(w_2^T x), x \in \mathbb{R}^2, \theta \in \mathbb{R}^6$
- Target: $f^* = \overline{a}\sigma(\overline{w}^T x)$
- **Target Set** $Q^* = F^{-1}(f^*)$ generally consists of three "branches" (sets)

 $L_{s}^{-1}(0$

(a)
$$Q_1 = \{(a_k, w_k)_{k=1}^2 : w_1 = w_2 = \overline{w}, a_1 + a_2 = \overline{a}\}$$

(b)
$$Q_2 = \{(a_k, w_k)_{k=1}^2 : w_1 = \overline{w}, a_1 = \overline{a}, a_2 = 0\}$$

(c)
$$Q_3 = \{(a_k, w_k)_{k=1}^2 : w_2 = \overline{w}, a_2 = \overline{a}, a_1 = 0\}.$$

As sample size *n* increases, how global min $L_s^{-1}(0)$ shrinks to Q^* ?

Illustration of Q^1, Q^2, Q^3

 Q_3



Geometry of global minima for final condensation





Typical convergence rate for final condensation

Gradient flows near Q^* **:** γ_1 : sublinear rate; γ_2, γ_3 : linear rate.





Stability of target branches underlies final condensation

Theorem 5.4 (recovery stability). Given $m_0 \leq r \leq m$, partition P and permutation π and separating inputs $\{x_i\}_{i=1}^n$. Then no point in $Q_{P,\pi}^r$ is recovery stable when $n \leq r + (r-l)d$ (l is the deficient number of P), and almost all points in $Q_{P,\pi}^r$ are recovery stable when $n \geq r + (m + m_0 - r)d$. Moreover, all points in Q^* are recovery stable when n > (d+1)m, namely, Q^* is recovery stable.

Sample size/Branches	Q^{m_0}	•••	Q^r		Q^m		
$\leq (d+1)m_0$	X	•••	X		X		
$\geq m + m_0 d$					\checkmark		
• • •				•••	•		
$\geq r + (m + m_0 - r)d$			\checkmark		\checkmark		
• • •			• •		•		
$\geq m_0 + md$	\checkmark	•••	\checkmark	•••	\checkmark		
> (d+1)m	\checkmark^*						
\checkmark^* : any point in Q^* is recovery stable							



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. arXiv:2406.18035, (2024).

Generalization consequence of condensation

Large initialization (no condensation)



Small initialization (Strong condensation)





Condensation improves sample efficiency



How many samples are required to recover f^* by NN_{wide}?



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



Quantification of condensation--model rank

Model:

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

Model rank:

$$r_{\theta} \coloneqq \operatorname{rank} DF(\theta) = \dim \operatorname{Im}(DF(\theta))$$

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

Intuition: effective degrees of freedom at θ

$$F(\boldsymbol{\theta} + \boldsymbol{\delta})(\cdot) \approx F(\boldsymbol{\theta})(\cdot) + \sum_{i=1}^{M} \partial_{\theta_i} F(\boldsymbol{\theta})(\cdot) \delta_i$$

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



lower model rank signifies stronger condensation

Example:

$$F(\boldsymbol{\theta})(x) = a_1 \tanh(w_1 x) + a_2 \tanh(w_2 x)$$

Model rank:

dim span{tanh($w_1 x$), a_1 tanh'($w_1 x$)x, tanh($w_2 x$), a_2 tanh'($w_2 x$)x}

• **Condensed**($w_1 = \pm w_2$):

$$r_{\theta} \leq 2$$

 $r_{0} = 4$

• Not condensed($w_1 \neq \pm w_2 \neq 0, a_1 \neq 0, a_2 \neq 0$):

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



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,



Optimistic sample size reflects practice

Theorem 5 (optimistic sample sizes for two-layer tanh-NN). Given a two-layer NN $f_{\theta}(\boldsymbol{x}) = \sum_{i=1}^{m} a_i \tanh(\boldsymbol{w}_i^T \boldsymbol{x}), \boldsymbol{x} \in \mathbb{R}^d, \boldsymbol{\theta} = (a_i, \boldsymbol{w}_i)_{i=1}^m$, 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. arXiv:2406.18035, (2024).

Optimistic sample size reflects practice

Theorem 5 (optimistic sample sizes for two-layer tanh-NN). Given a two-layer NN $f_{\theta}(\boldsymbol{x}) = \sum_{i=1}^{m} a_i \tanh(\boldsymbol{w}_i^T \boldsymbol{x}), \boldsymbol{x} \in \mathbb{R}^d, \boldsymbol{\theta} = (a_i, \boldsymbol{w}_i)_{i=1}^m$, 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. arXiv:2406.18035, (2024).





wider network is sample efficient



Yaoyu Zhang*, Leyang Zhang, Zhongwang Zhang and Zhiwei Bai, Local Linear Recovery Guarantee of Deep Neural Networks at Overparameterization. arXiv:2406.18035, (2024).

Specialty of DNN models via optimistic estimate



Sample inefficient: Adding (unnecessary) connections worsens generalization

Sample efficient:

Increasing width doesn't harm generalization



Principle of model scaling

- Freely increase width
- Refrain from adding connection



vs. Scaling of brain

- mouse: $\sim 10^8$ neurons, $\sim 10^3 10^4$ connections/neuron
- human: $\sim 10^{11}$ neurons, $\sim 10^3 10^4$ connections/neuron

Condensation improves sample efficiency

Picture of sample size requirement for nonlinear models



Ways to facilitate condensation:

smaller initialization, dropout, larger weight decay, large learning rate



The origin of condensation



Permutation symmetry: e.g., $j, j' \in [m_{l-1}]$

$$f^{[l]}(x;\theta) = \sigma\left(\sum_{j=1}^{m_{l-1}} W^{[l-1]}_{,j} \sigma\left(W^{[l-2]}_{j}f^{[l-2]}(x;\theta) + b^{[l-2]}_{j}\right) + b^{[l-1]}\right)$$

Theorem(informal):

permutation-invariant manifolds are invariant manifolds of gradient flow. e.g., $\left(W_{j}^{[l-1]}, W_{j}^{[l-2]}, b_{j}^{[l-2]}\right) = \left(W_{j'}^{[l-1]}, W_{j'}^{[l-2]}, b_{j'}^{[l-2]}\right)$

Permutation symmetry → invariant manifolds (equiv to smaller network) → optimistically as efficient as smaller networks







permutation symmetry -> optimistic sample efficiency preserving

Permutation symmetric:

Embedding dim: d_{model} Attention mat dim: dHeads: h



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





KAN:

$$f_{ heta}(oldsymbol{x}) = f_{ heta}\left(x_1, \cdots, x_d
ight) = \sum_{i=1}^m \Phi_{oldsymbol{i}}\left(\sum_{j=1}^d \phi_{oldsymbol{i},j}\left(x_j
ight)
ight)$$

width: permutation symmetric

Increase *m* **preserves sample efficiency**

$$\Phi_i(x) ext{ or } \phi_{i,j}(x) = \sum_{i=1}^G c_i B_i(x), heta = \{c_i\} ext{ is learnable}$$

grid: no symmetry

Increase *G* **reduces sample efficiency**

Experiment





Our series works on condensation



A1. Regime of condensation—phase diagram series

1. Tao Luo, Zhi-Qin John Xu, Zheng Ma, Yaoyu Zhang, <u>"Phase Diagram for Two-layer ReLU Neural Networks at Infinite-Width Limit,"</u> 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, <u>"Empirical Phase Diagram for Three-layer Neural Networks</u> with Infinite Width," NeurIPS 2022.

A2. Loss landscape structure—embedding principle series

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

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

3.Hanxu Zhou, Qixuan Zhou, Tao Luo, Yaoyu Zhang, Zhi-Qin John Xu, <u>"Towards Understanding the Condensation of Neural Networks at Initial</u> Training," NeurIPS 2022.

4.Zhiwei Bai, Tao Luo, Zhi-Qin John Xu, Yaoyu Zhang, <u>"Embedding Principle in Depth for the Loss Landscape Analysis of Deep Neural Networks,</u>" CSIAM Trans. Appl. Math., 5 (2024), pp. 350-389.

A3. Generalization advantage—optimistic estimate series

1.Yaoyu Zhang, Zhongwang Zhang, Leyang Zhang, Zhiwei Bai, Tao Luo, Zhi-Qin John Xu, <u>"Linear Stability Hypothesis and Rank Stratification for</u> 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, Zhiwei Bai, <u>"Local Linear Recovery Guarantee of Deep Neural Networks at</u> Overparameterization," arXiv:2406.18035 (2024).



Overview of our works on condensation



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





Thanks!