

IV. From condensation to loss landscape analysis

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



WWW.MOD.FAU.EU
#FAUMoDCourse

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



Phenomenon as a key family of facts to uncover

Frequency principle/spectral bias

Condensation

Double descent

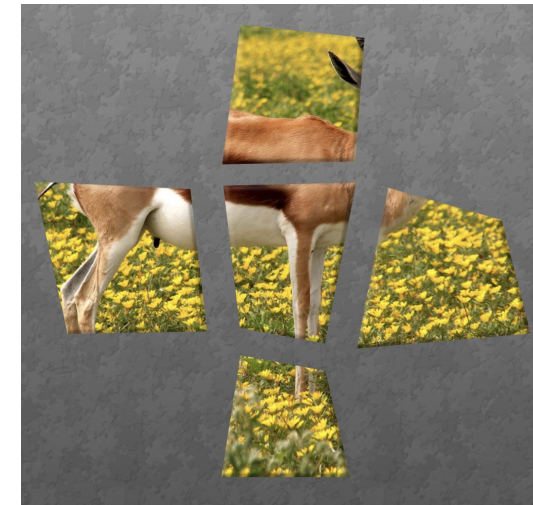
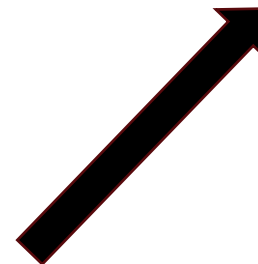
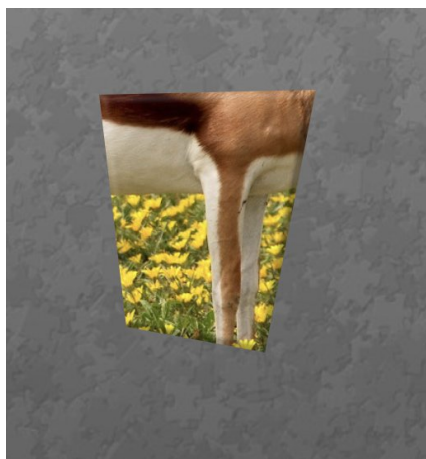
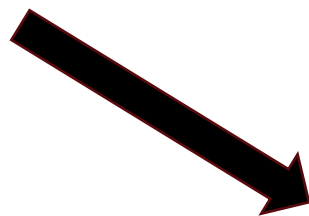
Edge of stability

Lottery ticket

Neural collapse

Grokking

.....





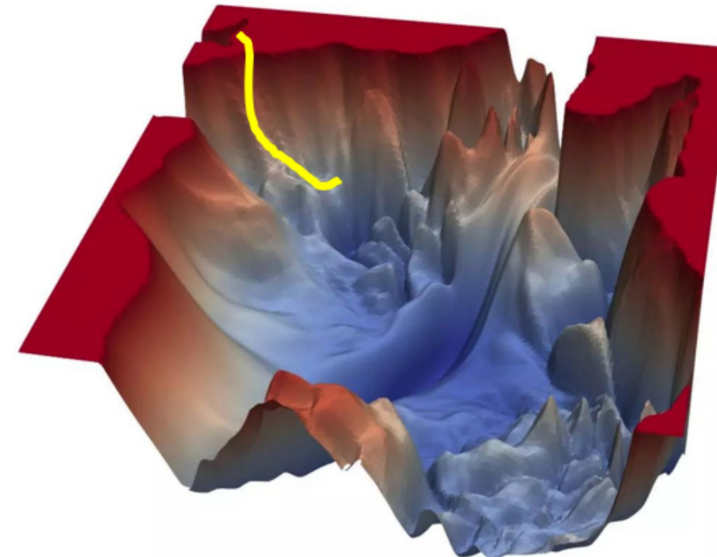
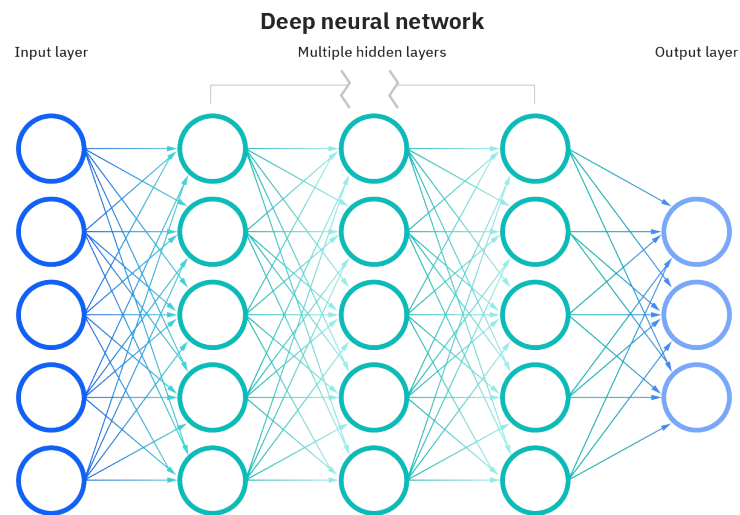
Deep learning loss landscape

$$R_S(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{f}(\mathbf{x}_i, \boldsymbol{\theta}), \mathbf{y}_i)$$

Model: $\mathbf{f}(\mathbf{x}_i, \boldsymbol{\theta})$

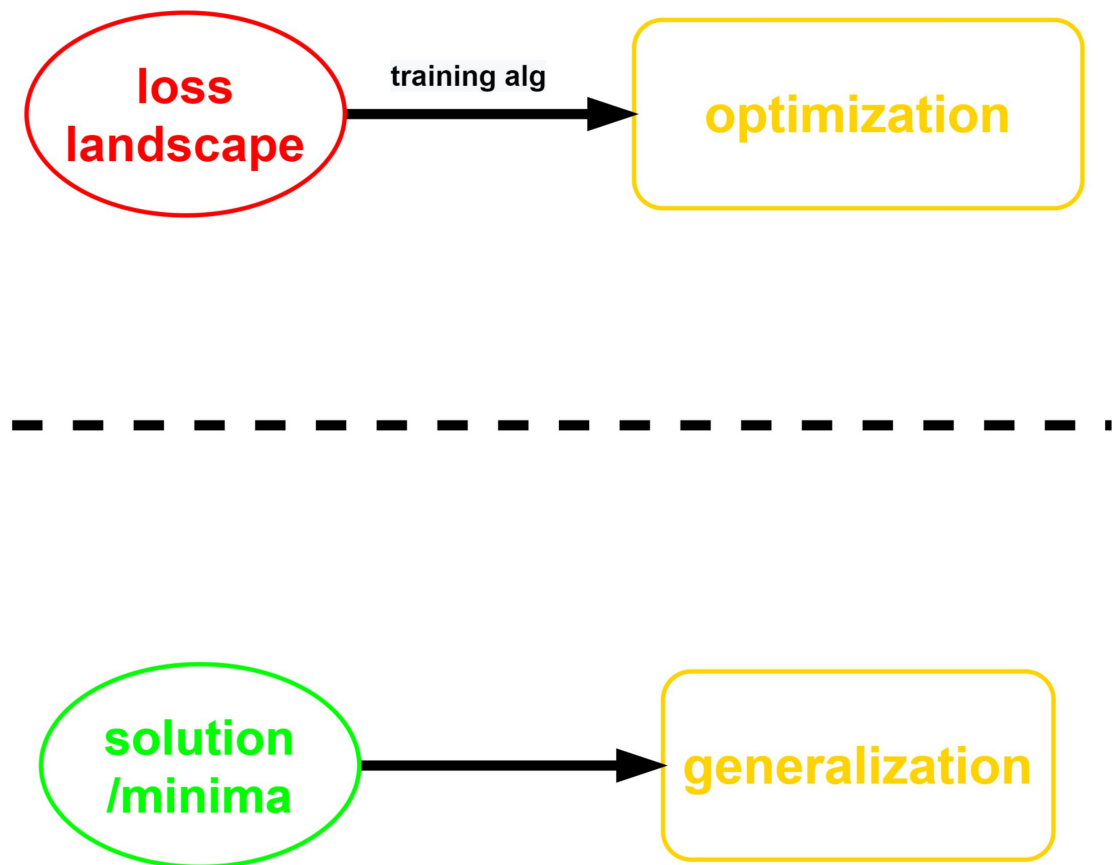
Data: $S = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$

Loss: $\ell(\cdot, \cdot)$

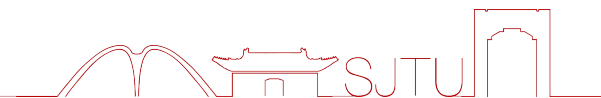
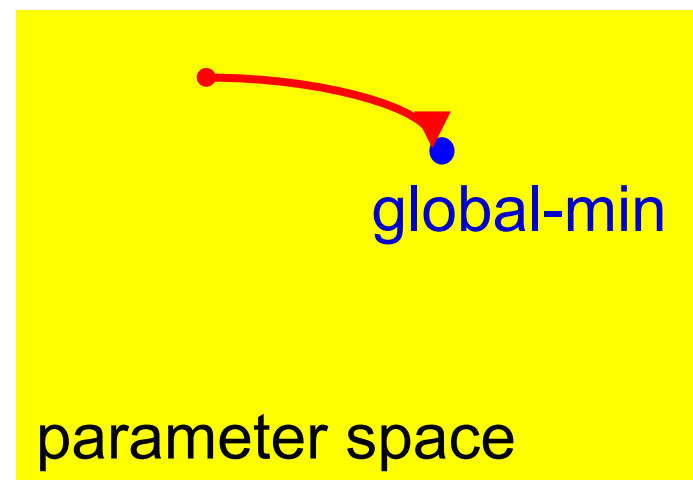




Role of loss landscape (conventional ML)

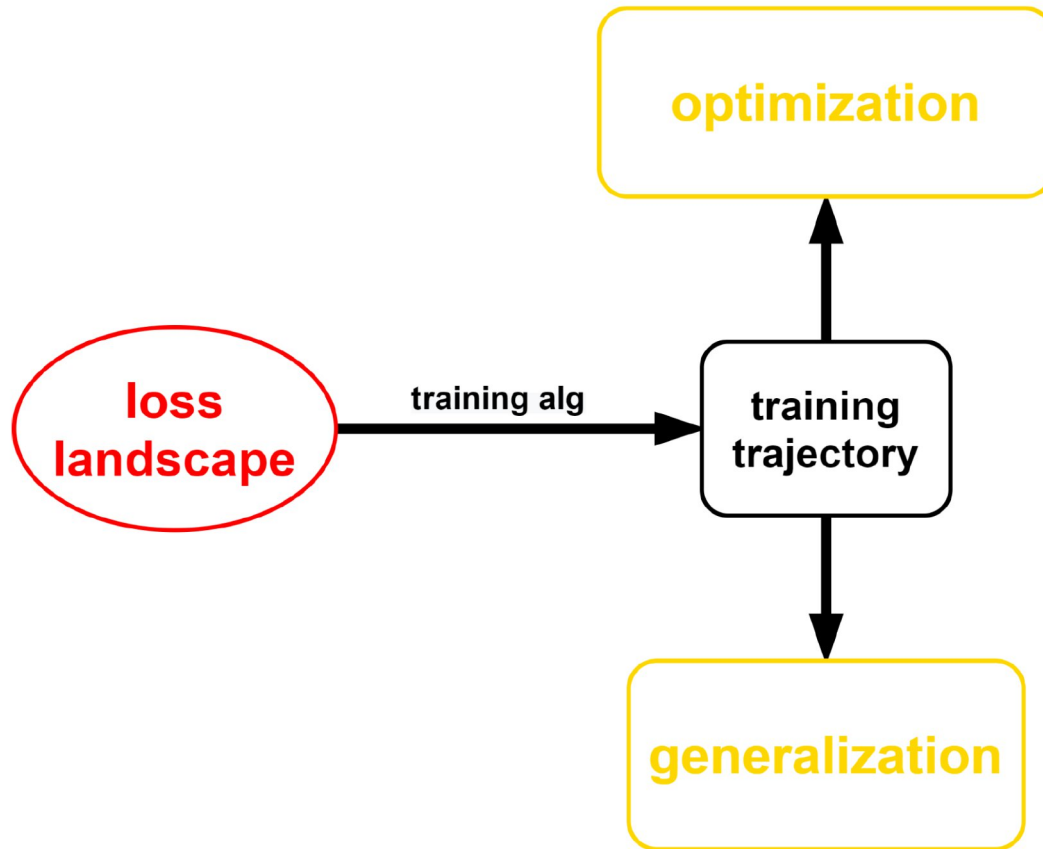


$$\# \text{param} \leq \# \text{data}$$

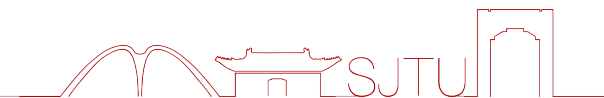
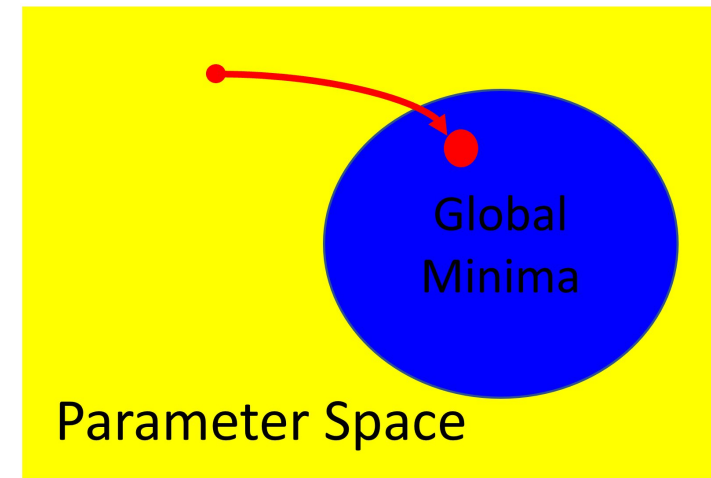




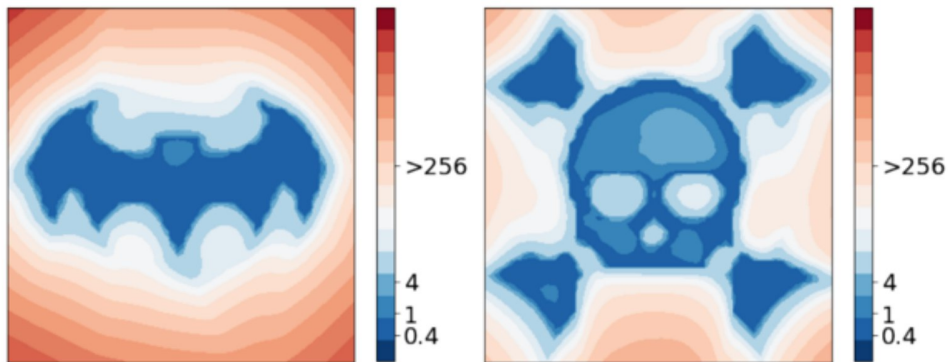
Role of loss landscape (deep learning)



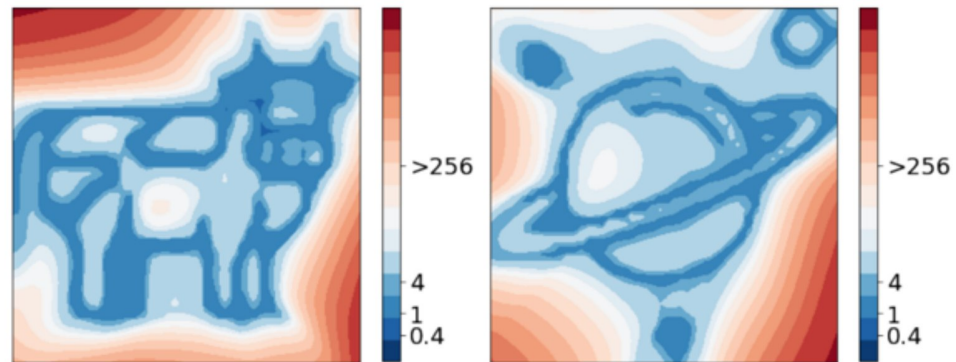
$\#param > \#data$



DNN loss landscape is complex



(a) Loss surface on FashionMNIST dataset



(b) Loss surface on CIFAR10 dataset

I. Skorokhodov, M. Burtsev, 2019



Which picture captures NN loss landscape?

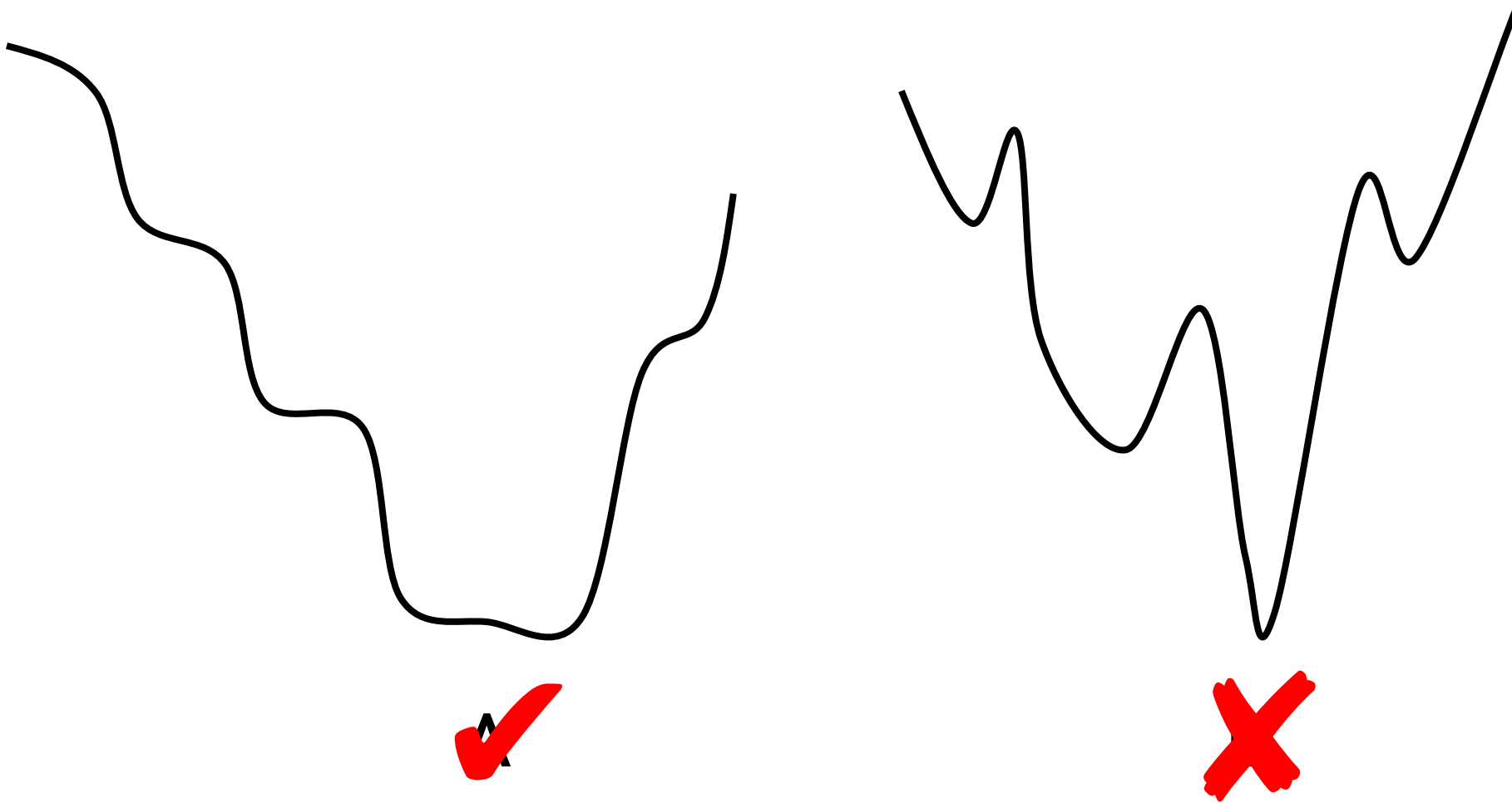
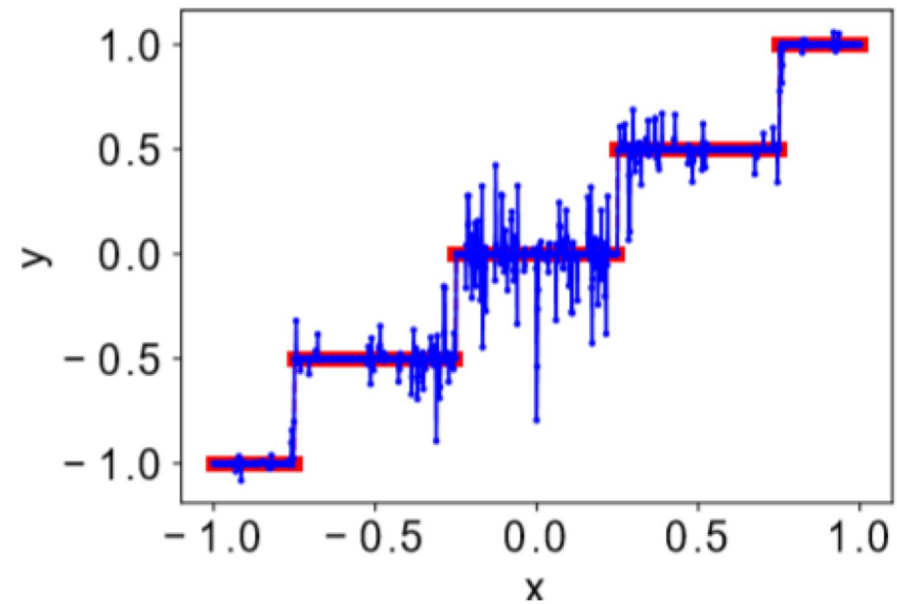
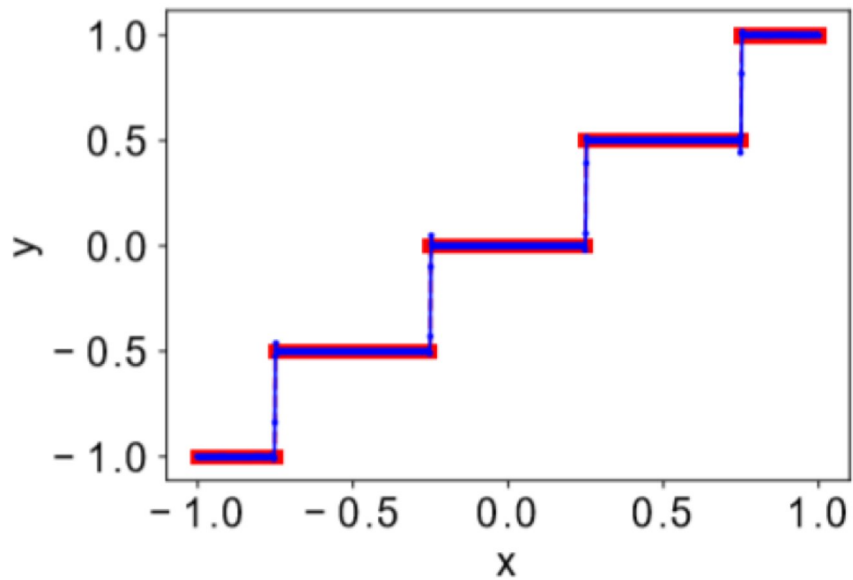


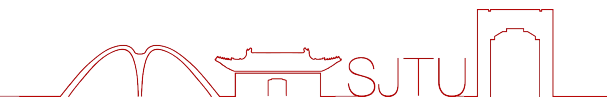


Illustration of different global-minima



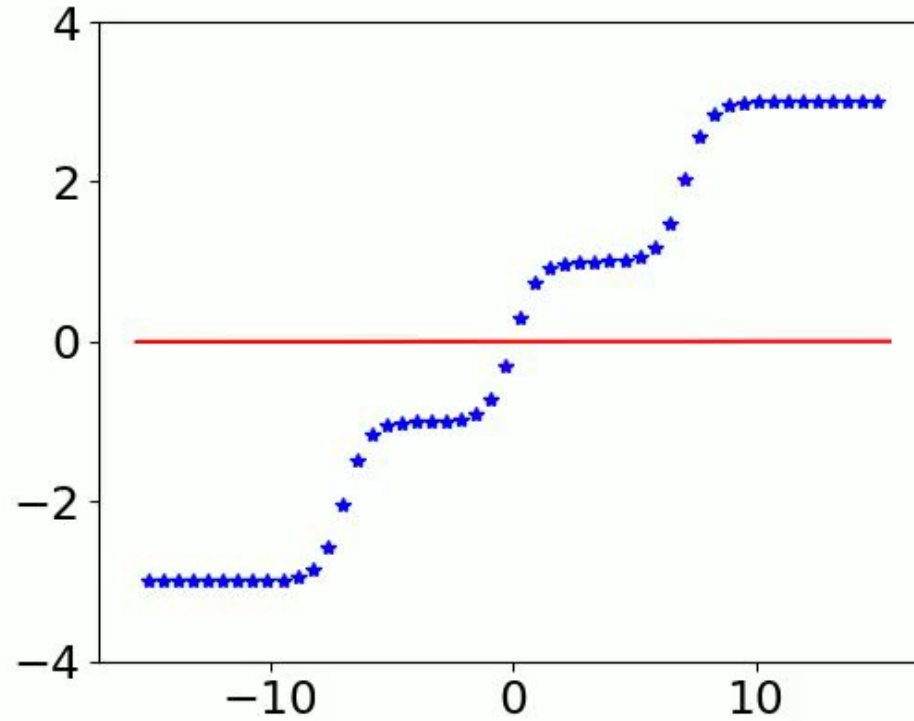
Loss landscape structure underlying condensation

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.
4. Tao Luo, Leyang Zhang, Yaoyu Zhang, "Structure and Gradient Dynamics Near Global Minima of Two-layer Neural Networks," arXiv:2309.00508 (2023).





Typical training behavior with strong condensation

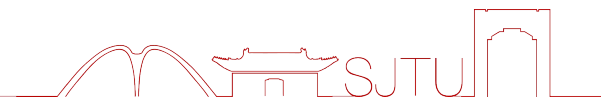
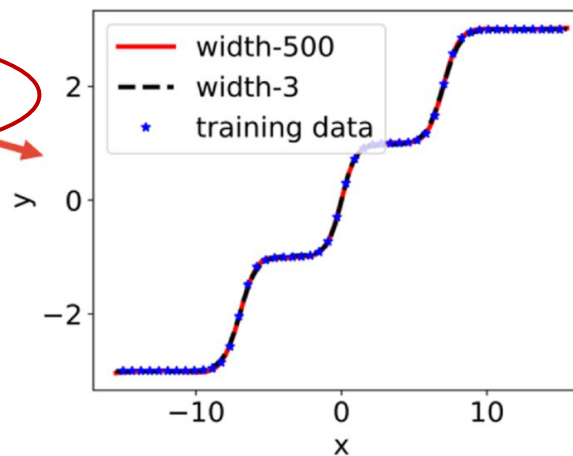
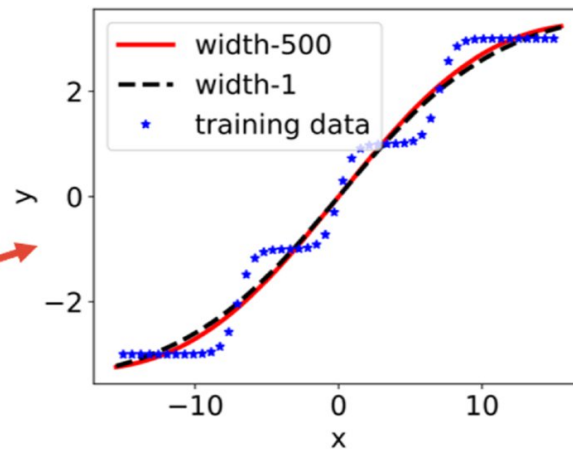
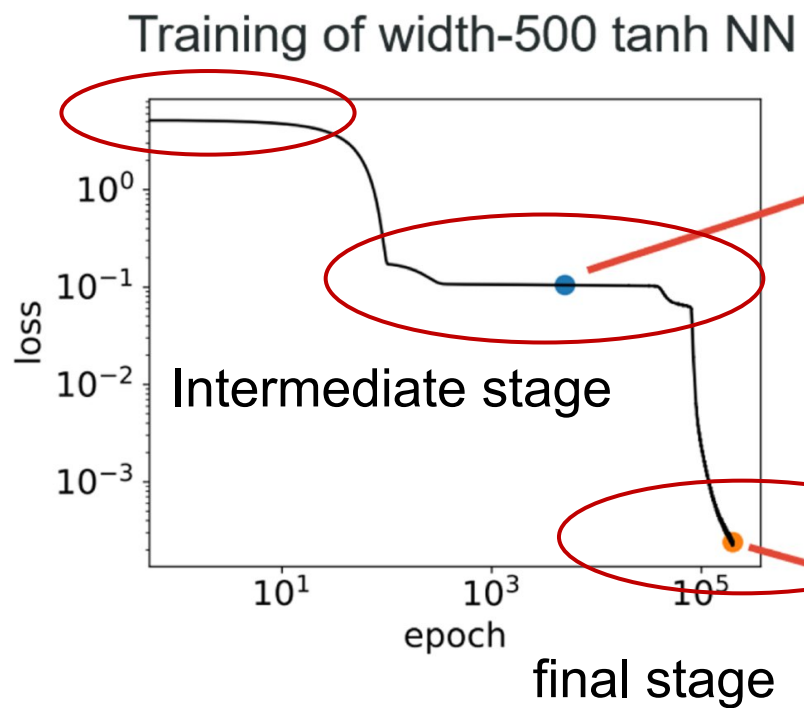


Width-500 tanh-NN (~1500 parameters)



Trajectory of training loss

Initial stage

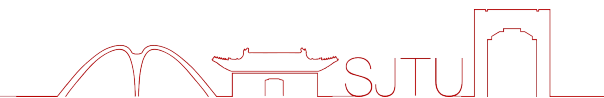
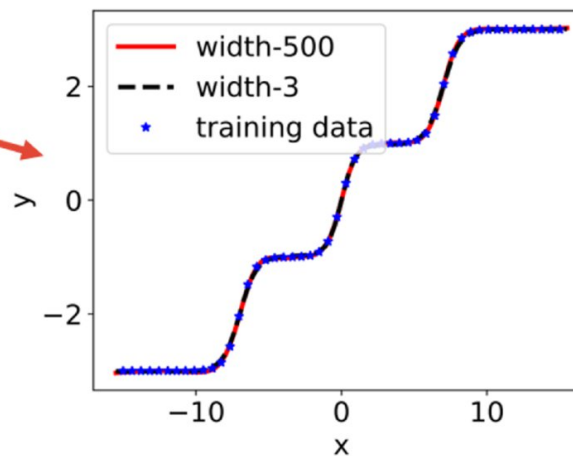
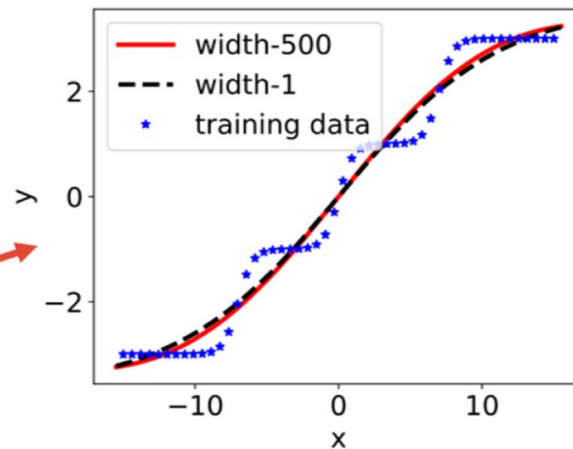
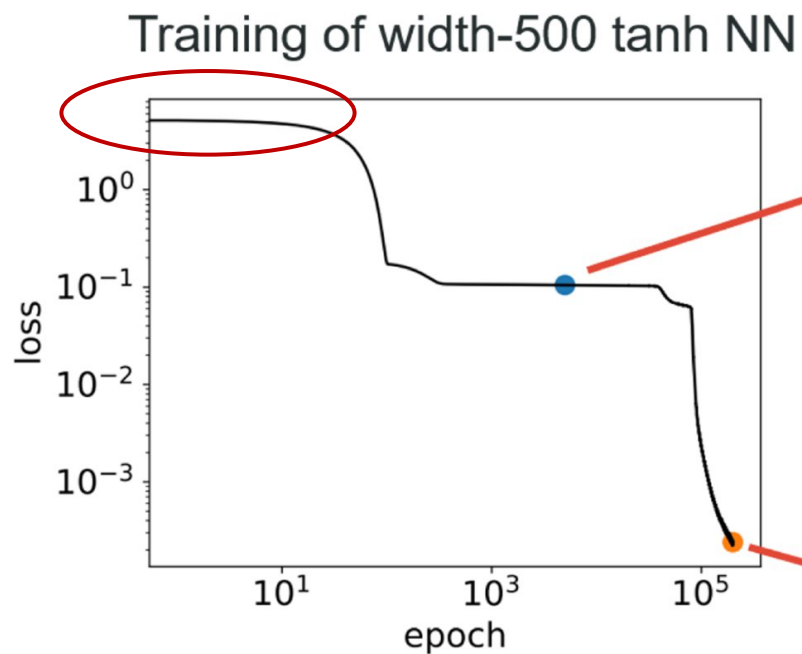


Initial condensation



Initial condensation

Initial stage





Loss landscape around 0 and Initial condensation

$$\dot{w}_j = \sum_{i=1}^m (y_i - f_{\theta}(x_i)) a_j \sigma'(w_j^T x_i) x_i$$

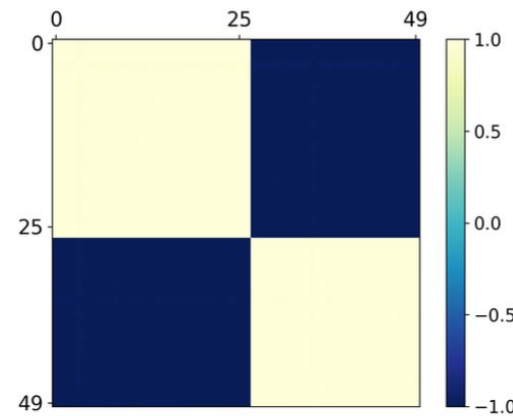
When $\theta \approx 0$, then $f_{\theta}(\cdot) \approx 0(\cdot)$:

$$\dot{w}_j \approx a_j \sum_{i=1}^m y_i \sigma'(w_j^T x_i) x_i$$

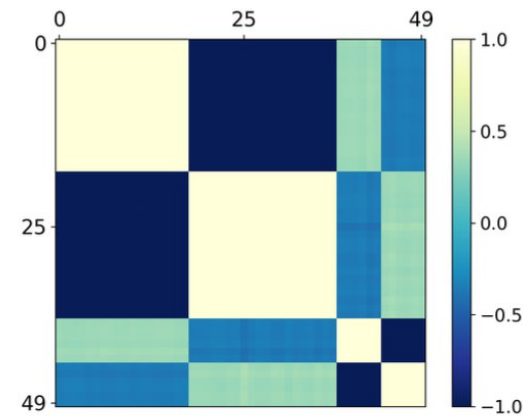
If $\sigma'(0) \neq 0$ (e.g. tanh, swish, gelu):

$$\dot{w}_j \approx a_j \sigma'(0) \sum_{i=1}^m y_i x_i$$

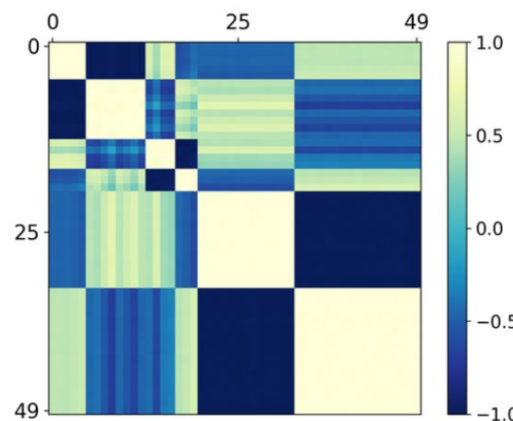
- i. No coupling between w_j and $w_{j'}$!
- ii. 2 limiting directions: $\pm \sum_{i=1}^m y_i x_i$.



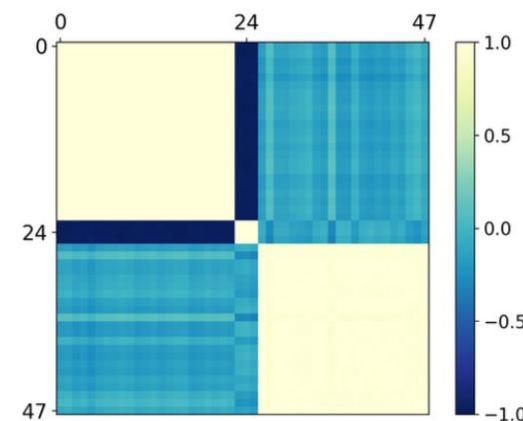
(a) $\tanh(x)$



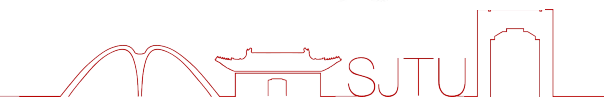
(b) $x \tanh(x)$



(c) $x^2 \tanh(x)$



(d) $\text{ReLU}(x)$



Intermediate condensation and embedding principle



Exercise

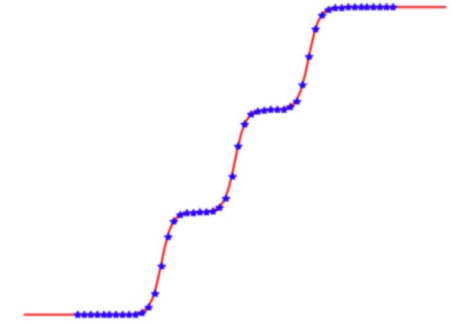


Target: $f^*(x) = \tanh(x-7) + \tanh(x) + \tanh(x+7)$

Data: $S = \{(x_i, f^*(x_i))\}_{i=1}^{50}$

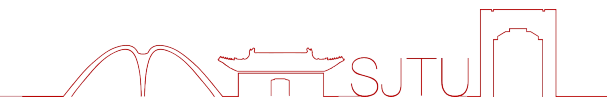
Model: $f_\theta(x) = \sum_{j=1}^m a_j \tanh(w_j x + b_j)$ ($\theta = [a_j, w_j, b_j]_{j=1}^m$)

Loss: $R_S(\theta) = \frac{1}{50} \sum_{i=1}^{50} |f_\theta(x_i) - f^*(x_i)|^2$



$m=3$,

1. Dimension of $R_S(\theta)$?
2. Existence of zero loss global minima?
3. Is $R_S(\theta)$ convex?
4. Are there non-global critical points? How many?
5. How many critical functions $\mathcal{F}^c := \{f_\theta(\cdot) | \nabla R_S(\theta) = 0\}$?
What are they?





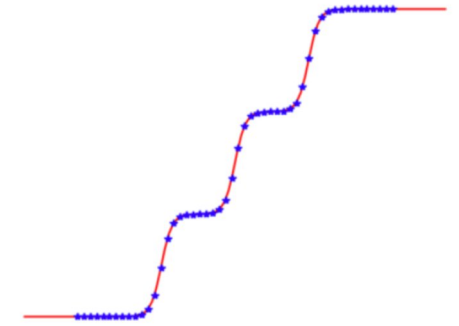
Answer

Target: $f^*(x) = \tanh(x-7) + \tanh(x) + \tanh(x+7)$

Data: $S = \{(x_i, f^*(x_i))\}_{i=1}^{50}$

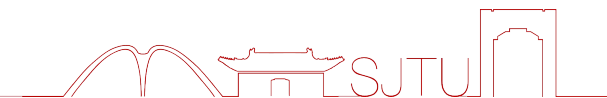
Model: $f_\theta(x) = \sum_{j=1}^m a_j \tanh(w_j x + b_j)$ ($\theta = [a_j, w_j, b_j]_{j=1}^m$)

Loss: $R_S(\theta) = \frac{1}{50} \sum_{i=1}^{50} |f_\theta(x_i) - f^*(x_i)|^2$



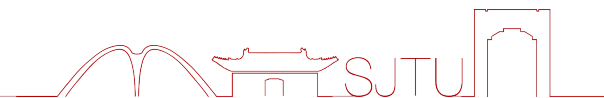
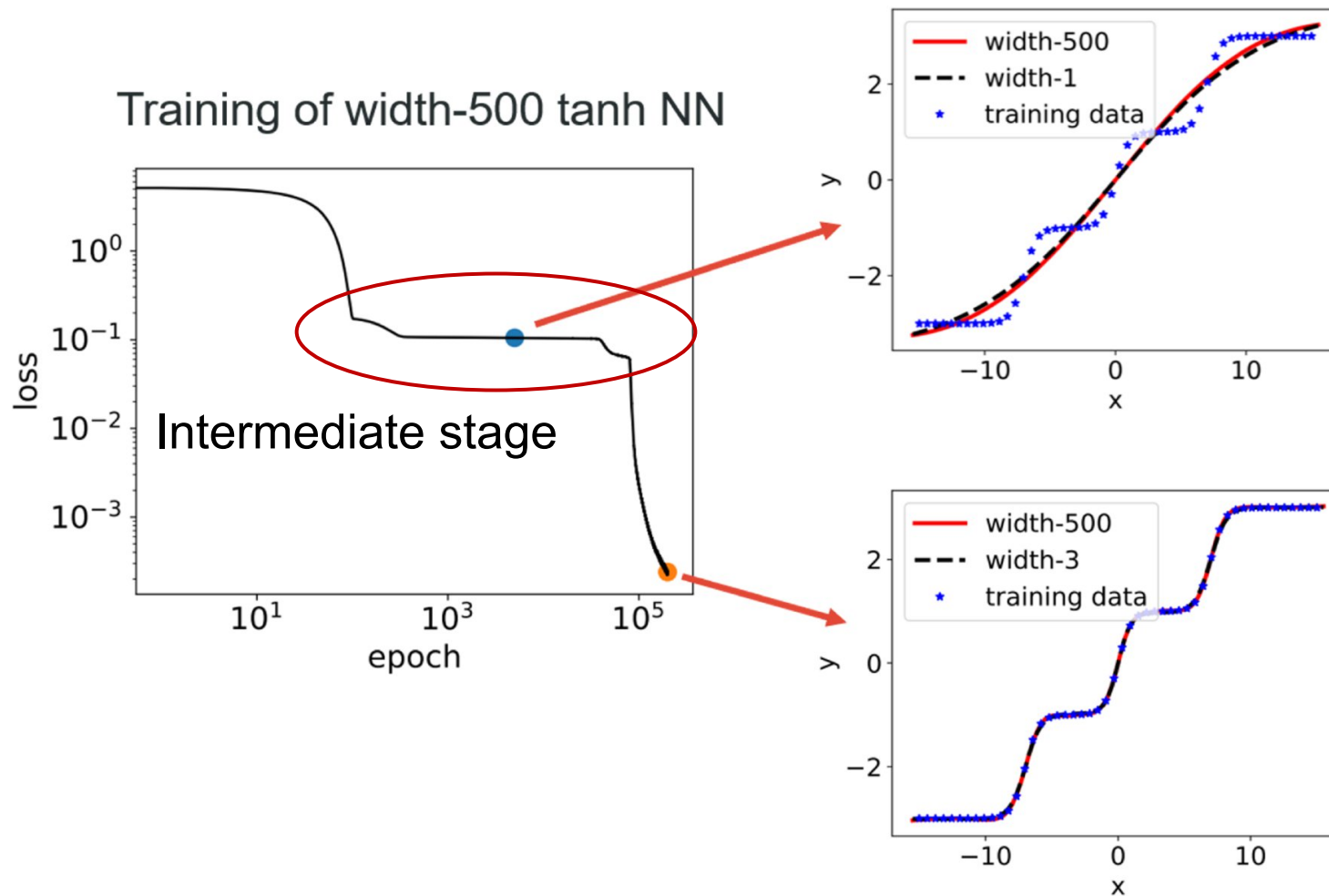
$m=3$,

1. Dimension of $R_S(\theta)$? **9**
2. Existence of zero loss global minima? **exist**
3. Is $R_S(\theta)$ convex? **no**
4. Are there non-global critical points? How many? **infinite**
5. How many critical functions $\mathcal{F}^c := \{f_\theta(\cdot) | \nabla R_S(\theta) = 0\}$?
What are they? **≥ 5**



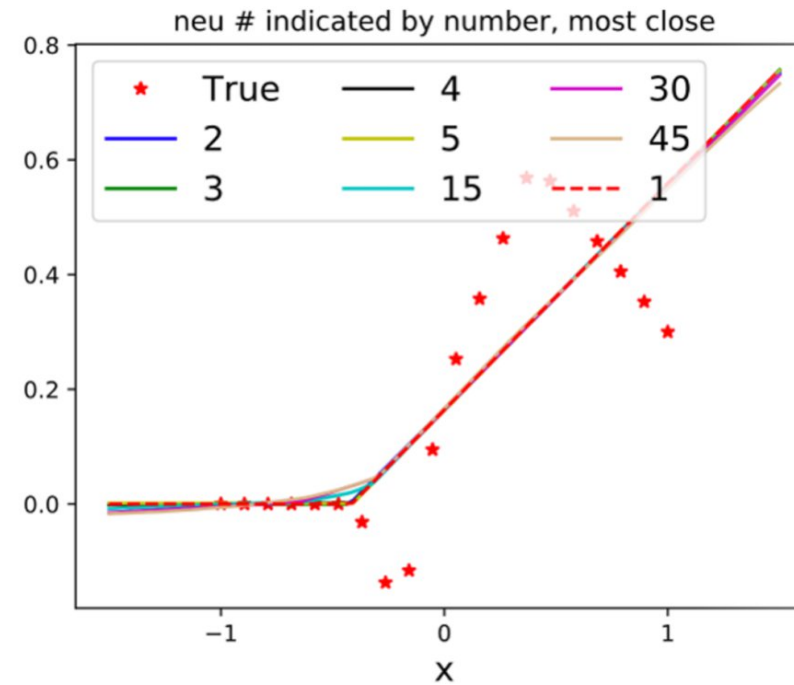
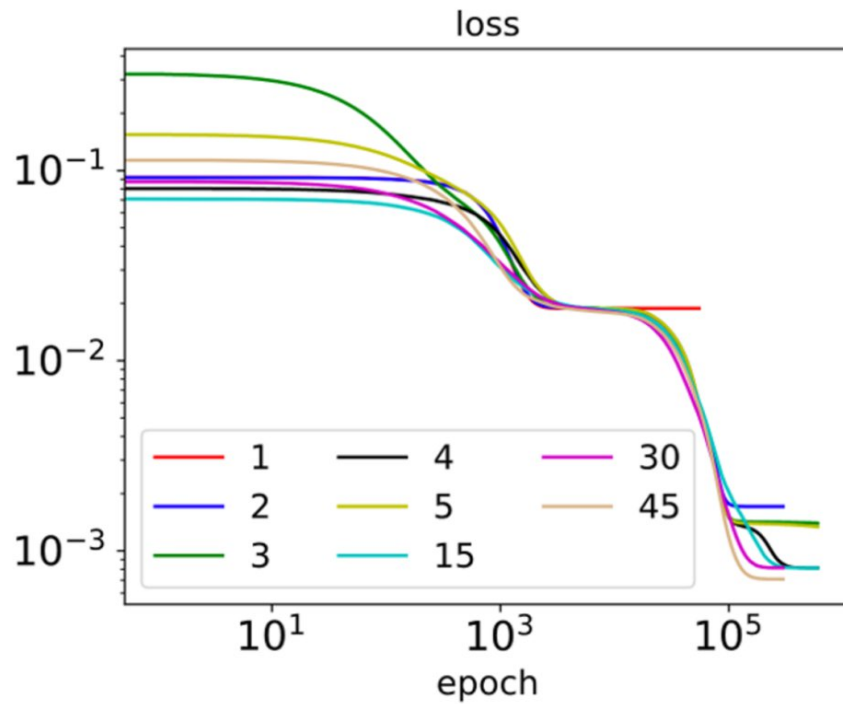


Intermediate condensation





Training similarity of NNs with different widths





Embedding Principle in width

Embedding Principle

The loss landscape of any network ``contains'' all critical points of all narrower networks.

Theorem

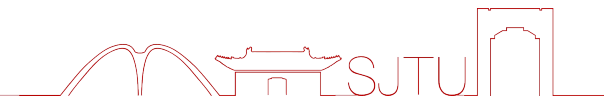
Critical functions of narrow NNs are critical functions of any wider NNs, i.e.,

$$\mathcal{F}_{\text{narr}}^c \subset \mathcal{F}_{\text{wide}}^c$$

where $\mathcal{F}^c := \{f_{\theta}(\cdot) | \nabla R_S(\theta) = 0\}$.

Implication

“simple” critical points always exist!



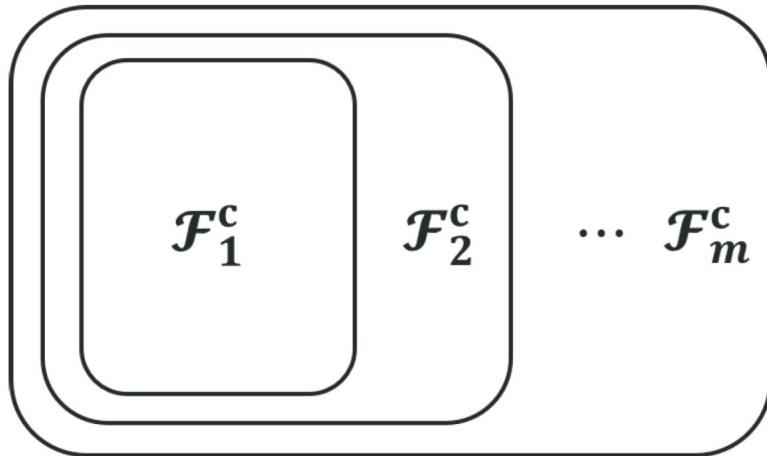
[1] Zhang, Zhang, Luo, Xu. *Embedding Principle of Loss Landscape of Deep Neural Networks*. NeurIPS 2021 Spotlight.

[2] Zhang, Li, Zhang, Luo, Xu. *Embedding Principle: a hierarchical structure of loss landscape of deep neural networks*. Journal of Machine Learning, 2022.

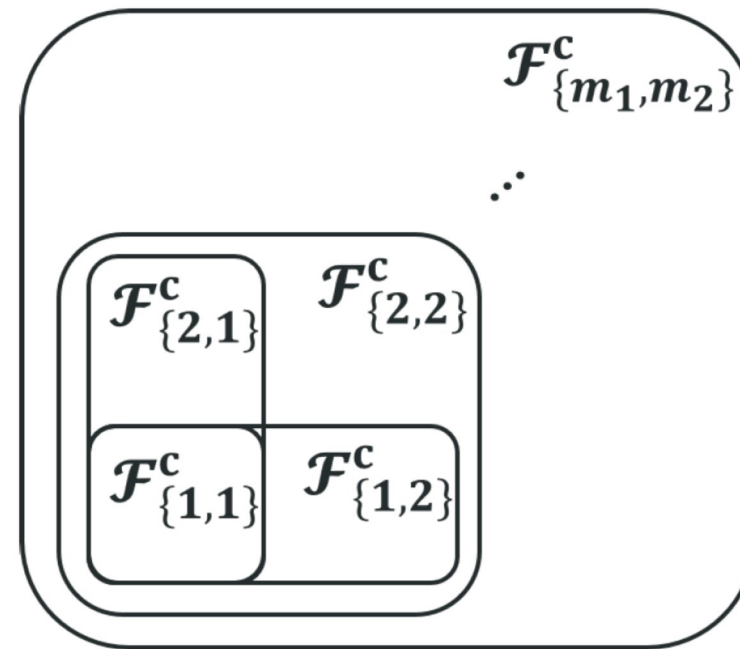


hierarchical structure of DNN loss landscape

Simple \rightarrow Complex



Simple \rightarrow Complex





Exercise



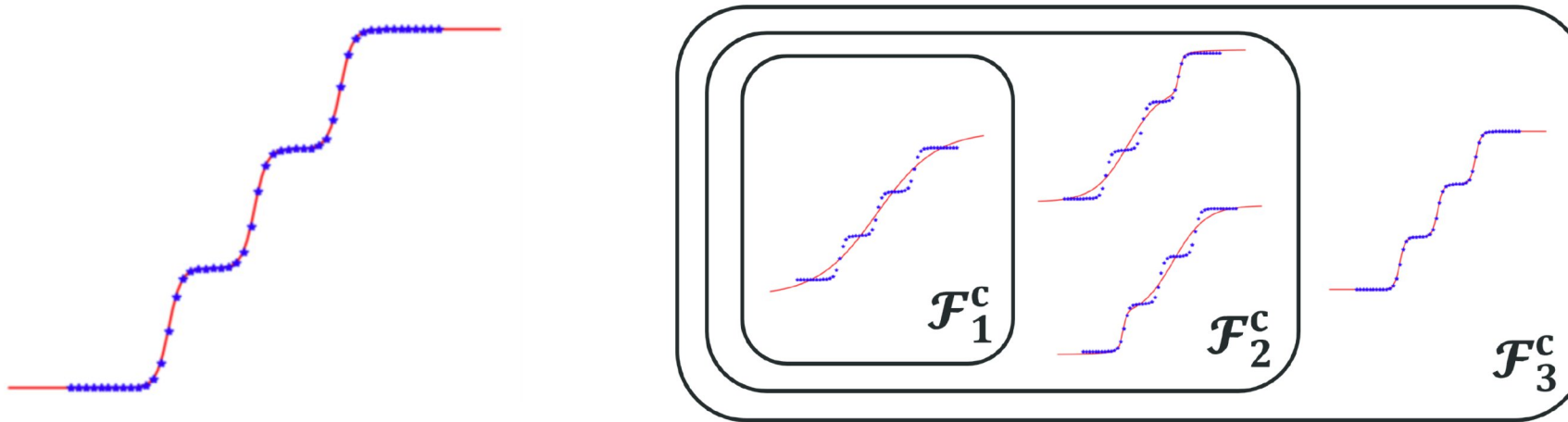
Objective: $f^*(x) = \tanh(x-7) + \tanh(x) + \tanh(x+7)$

Data: $S = \{(x_i, f^*(x_i))\}_{i=1}^{50}$

Model: $f_\theta(x) = \sum_{j=1}^3 a_j \tanh(w_j x + b_j)$ ($\theta = [a_j, w_j, b_j]_{j=1}^3$)

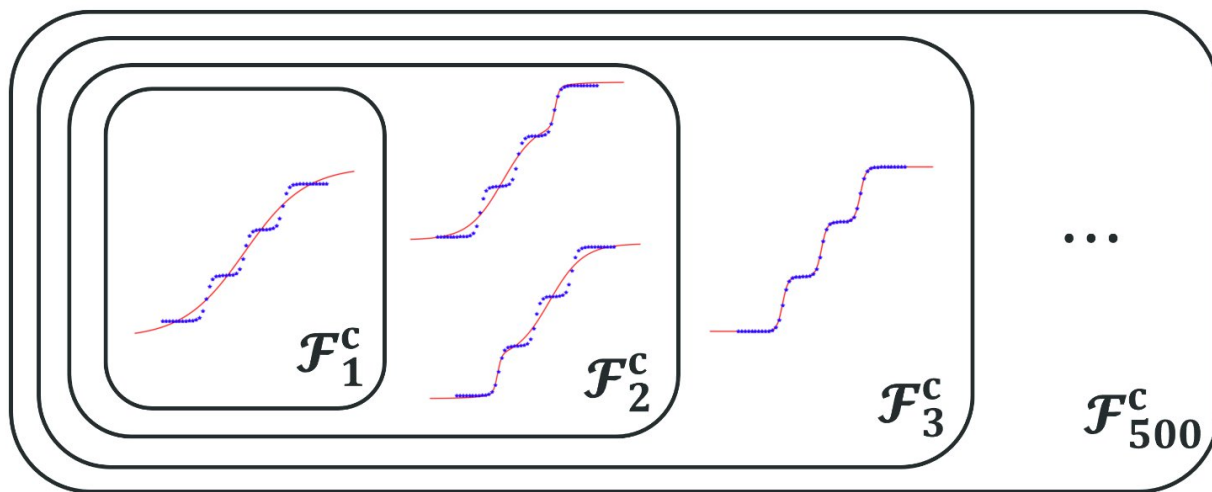
Loss landscape: $R_S(\theta) = \frac{1}{50} \sum_{i=1}^{50} |f_\theta(x_i) - f^*(x_i)|^2$

Critical functions $\mathcal{F}^c := \{f_\theta(\cdot) | \nabla R_S(\theta) = 0\}$?

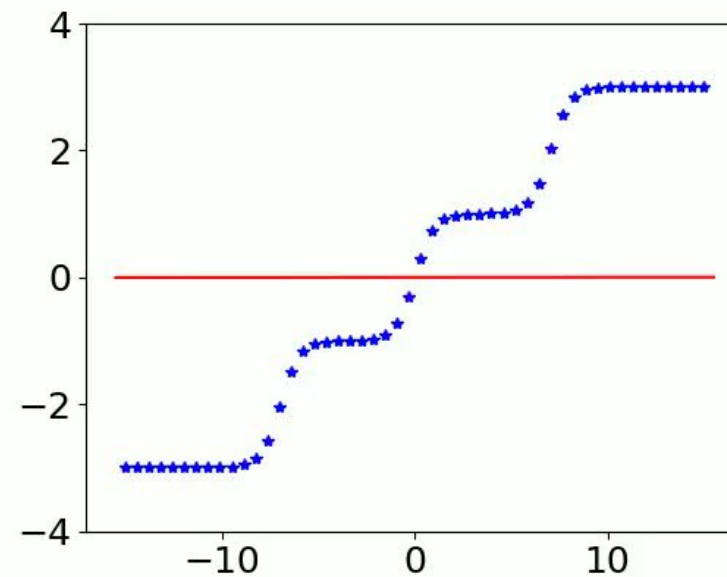




Example: analysis of hierarchical structure



500 tanh neuron





Key to the proof of Embedding Principle

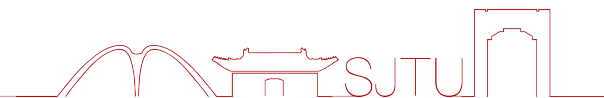
Key: discovering critical embedding

Discover **embedding** $\mathcal{T} : \mathbb{R}^{M_{\text{narr}}} \rightarrow \mathbb{R}^{M_{\text{wide}}}$ such that for $\theta_{\text{wide}} = \mathcal{T}(\theta_{\text{narr}})$

(i) **output preserving**: $f_{\theta_{\text{narr}}} = f_{\theta_{\text{wide}}}$;

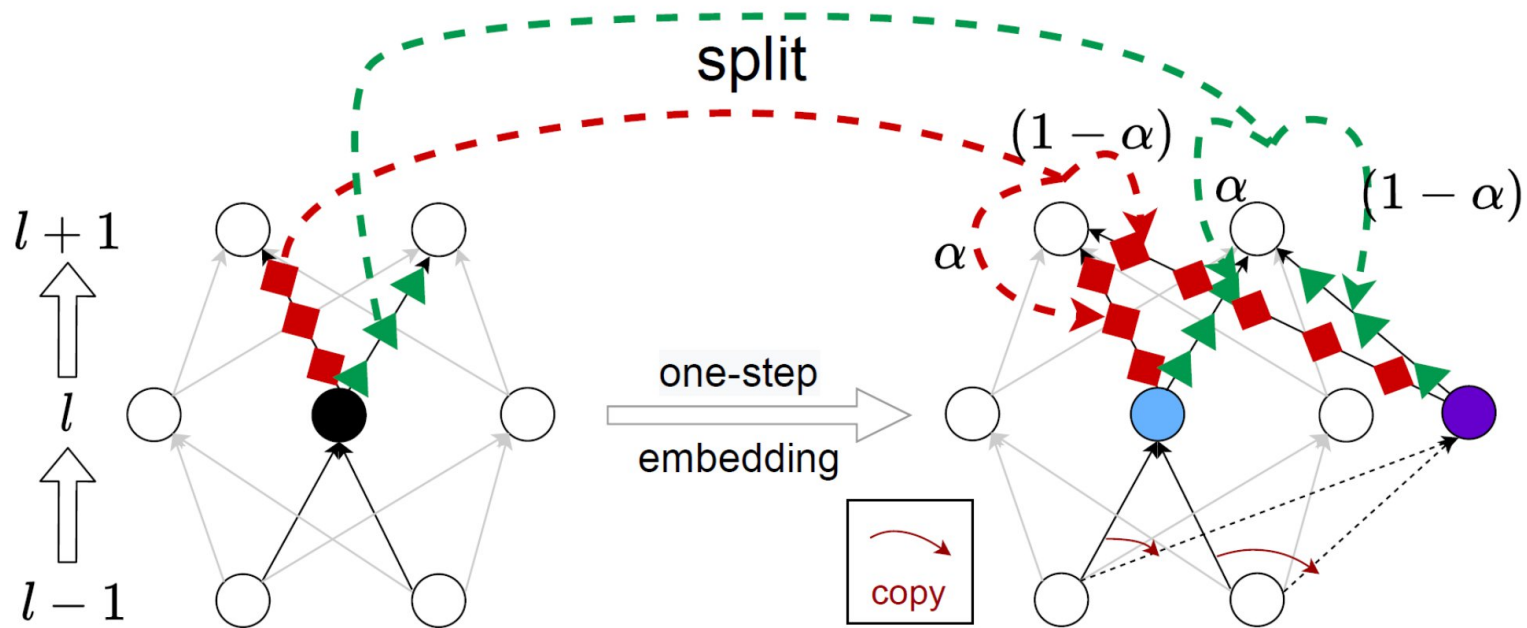
(ii) **criticality preserving**: if θ_{narr} is a critical point, then θ_{wide} is also a critical point.

critical embedding exists \Rightarrow Embedding Principle

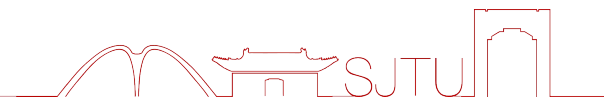




One-step splitting embedding



One-step embedding $\mathcal{T}_{l,s}^{\alpha}$.





Splitting embedding is critical embedding

Proposition (**output and representation preserving**)

For any point θ_{narr} of a DNN, a point θ_{wide} of a wider DNN obtained from θ_{narr} by **one-step embedding** satisfies

$$\mathbf{f}_{\theta_{\text{narr}}}(\mathbf{x}) = \mathbf{f}_{\theta_{\text{wide}}}(\mathbf{x}) \text{ for any } \mathbf{x}.$$

Theorem (**criticality preserving**)

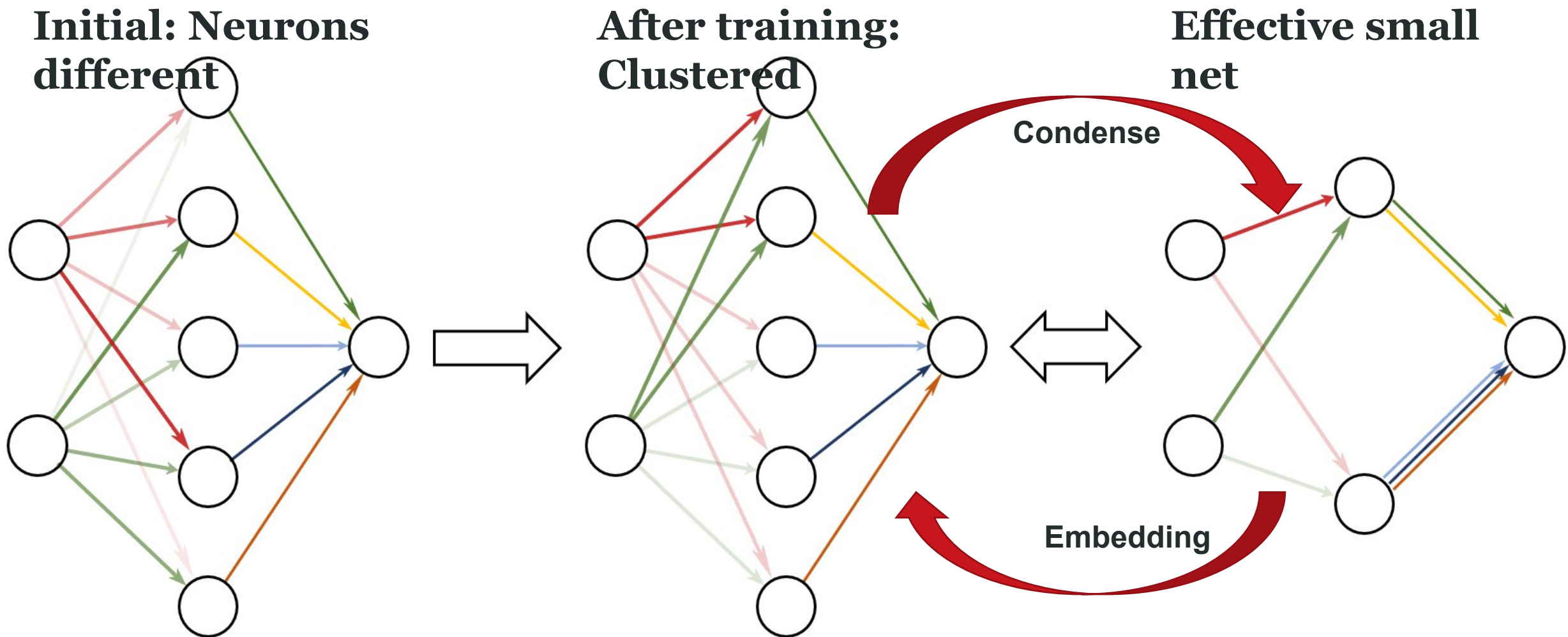
*For any critical point θ_{narr} of a DNN, a point θ_{wide} of a wider DNN obtained from θ_{narr} by **one-step embedding** is a critical point.*

Remark: Obviously, **multi-step embedding**, i.e., composition of **one-step embedding**, is also critical embedding.





Embedding as “inverse” of condensation

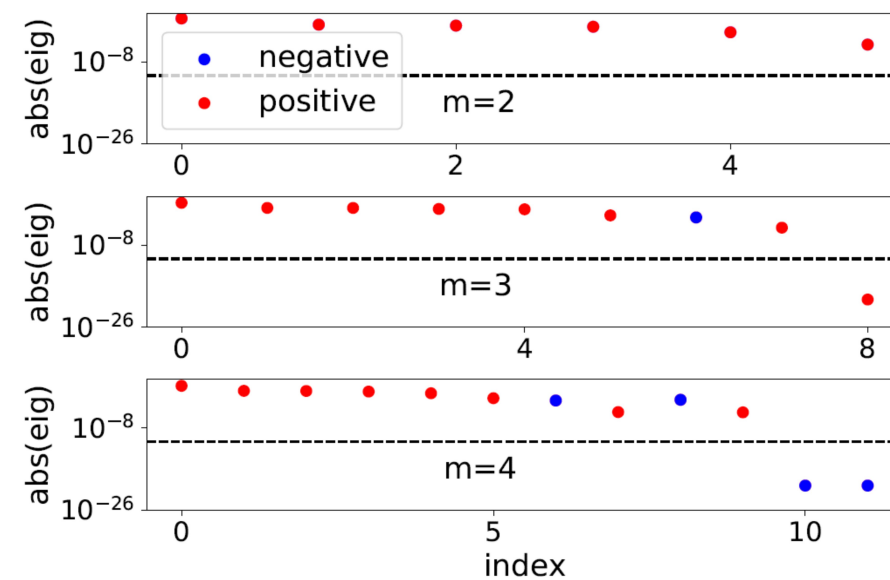
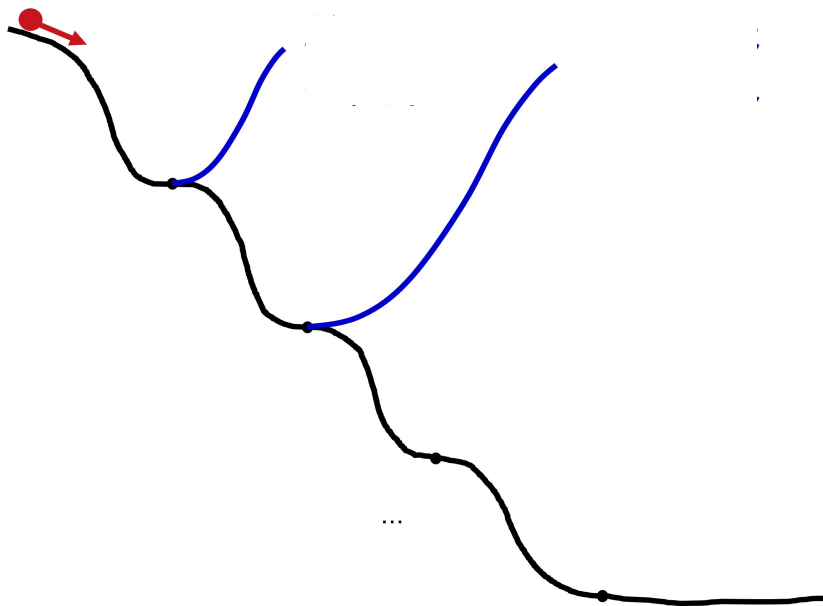


Implication of Embedding Principle



Implication—optimization

local-min of narrow NN \rightarrow saddles of wide NN



(a) synthetic data



Transition to strict saddle points is Irreversible.

Theorem

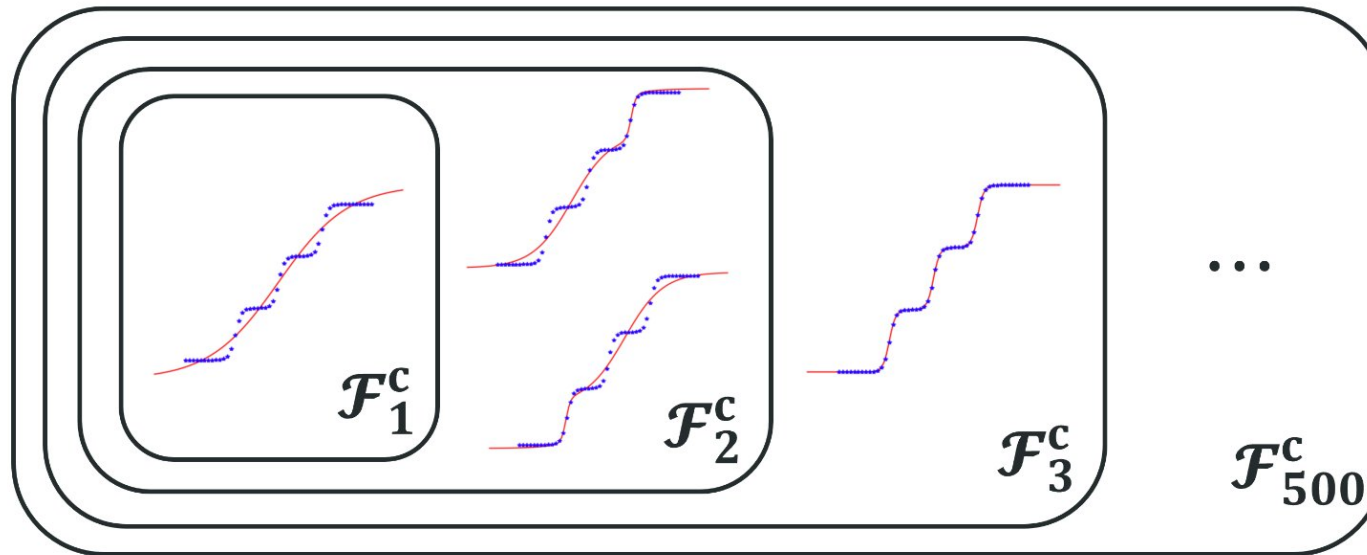
Given an $\text{NN}(\{m_l\}_{l=0}^L)$ and any of its parameters $\theta \in \mathbb{R}^M$, for any critical embedding $\mathcal{T} : \mathbb{R}^M \rightarrow \mathbb{R}^{M'}$ to any wider $\text{NN}(\{m'_l\}_{l=0}^L)$, the number of positive, zero, negative eigenvalues of $\mathbf{H}_S(\mathcal{T}(\theta))$ is no less than the counterparts of $\mathbf{H}_S(\theta)$.



Training and generalization

Observation

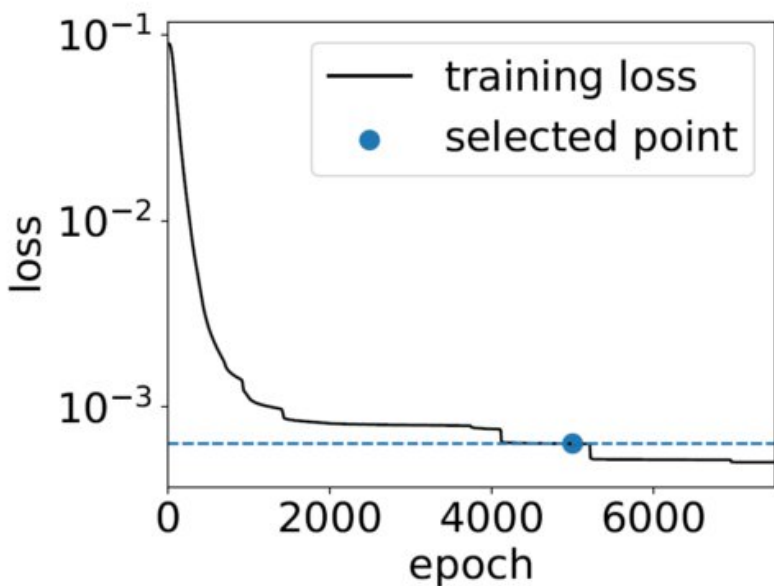
Nonlinear training of DNNs tends to learn simple critical functions.



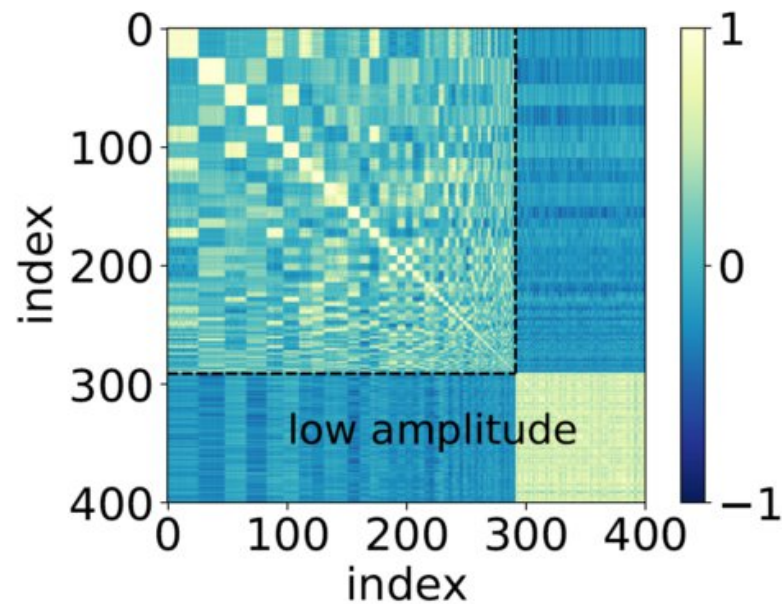


Implication—pruning

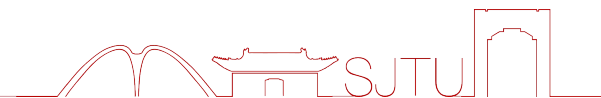
“simple” critical points has huge pruning potential



(a) initial loss



(b) orientation similarity

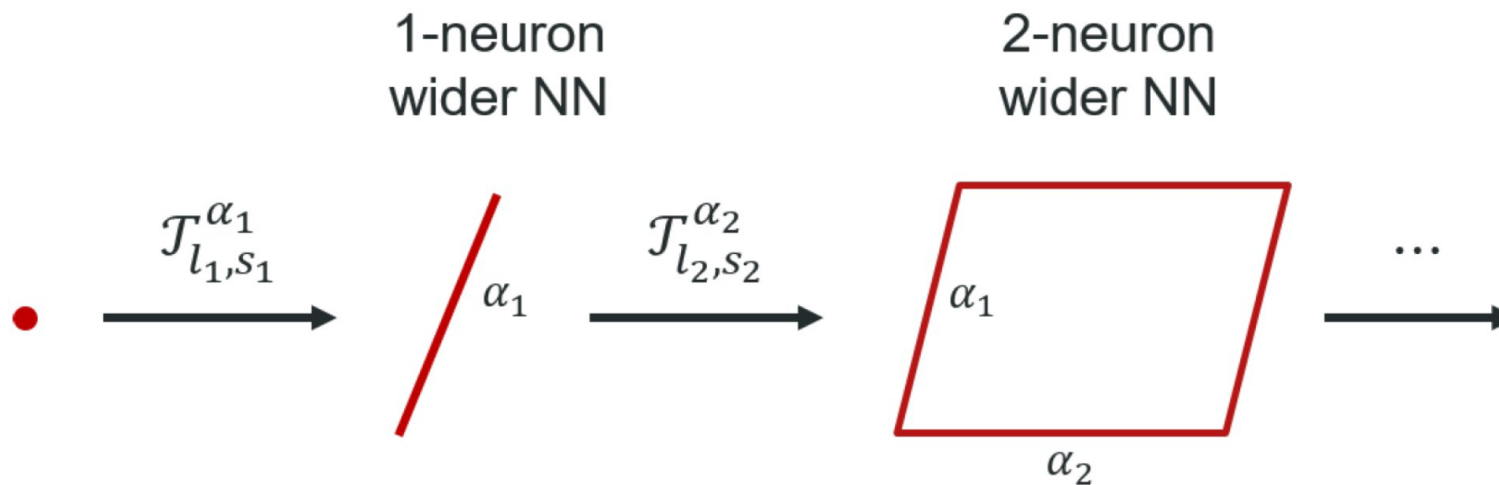




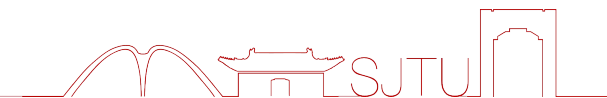
Dimension of critical submanifolds

Theorem (informal)

(Under mild assumption) Any critical point θ^c of a DNN can be embedded to *K -dimensional critical affine subspaces* of a *K -neuron wider DNN*.



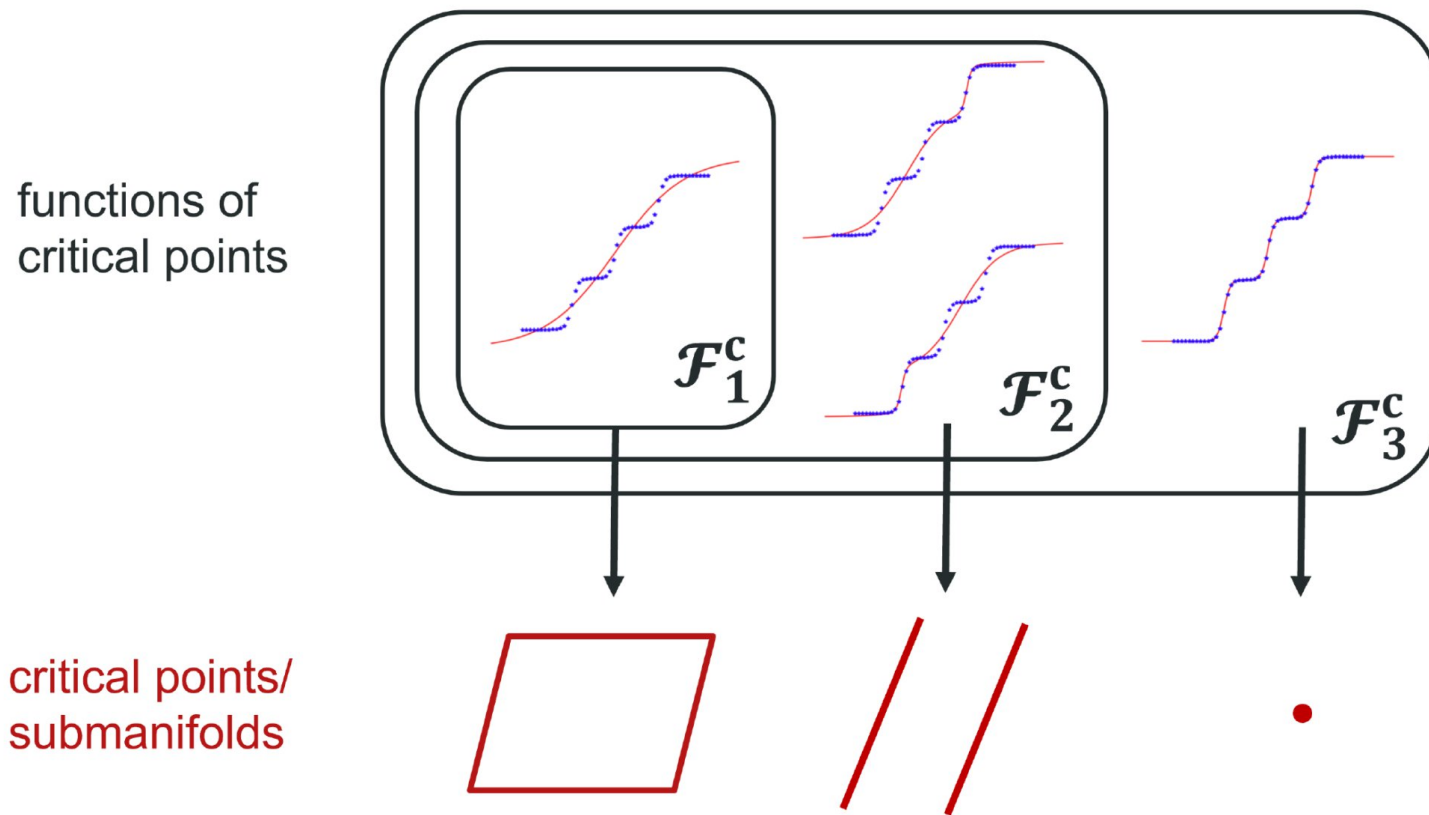
“Simple” critical functions possess high dimensional critical submanifolds.





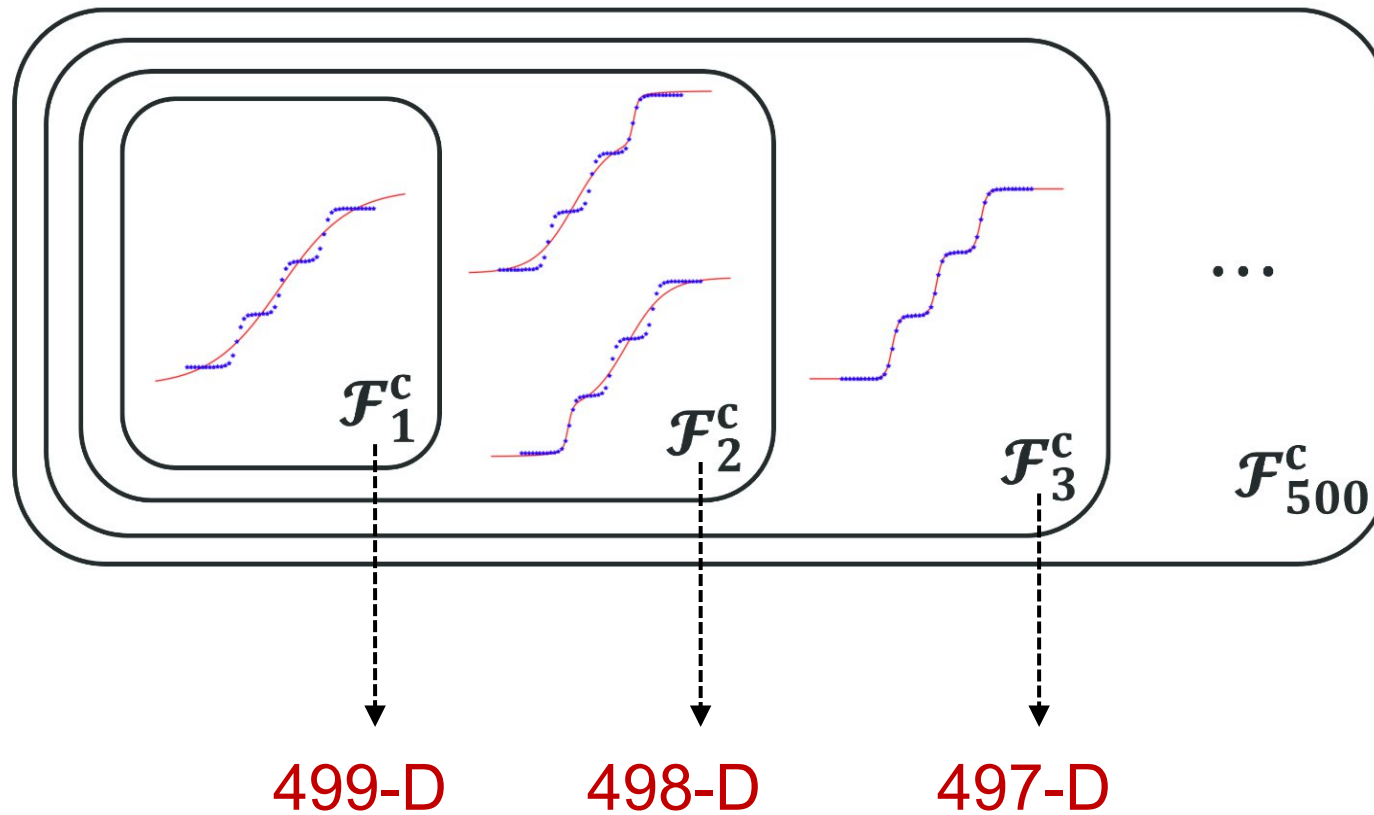
Loss landscape analysis of width-3 tanh-NN

$$R_S(\theta) = \frac{1}{50} \sum_{i=1}^{50} (f_{\theta}(x_i) - y_i)^2, \quad f_{\theta}(x) = \sum_{j=1}^3 a_j \tanh(w_j x + b_j)$$





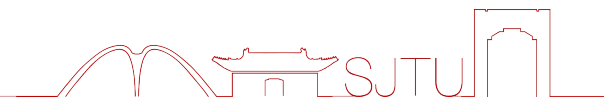
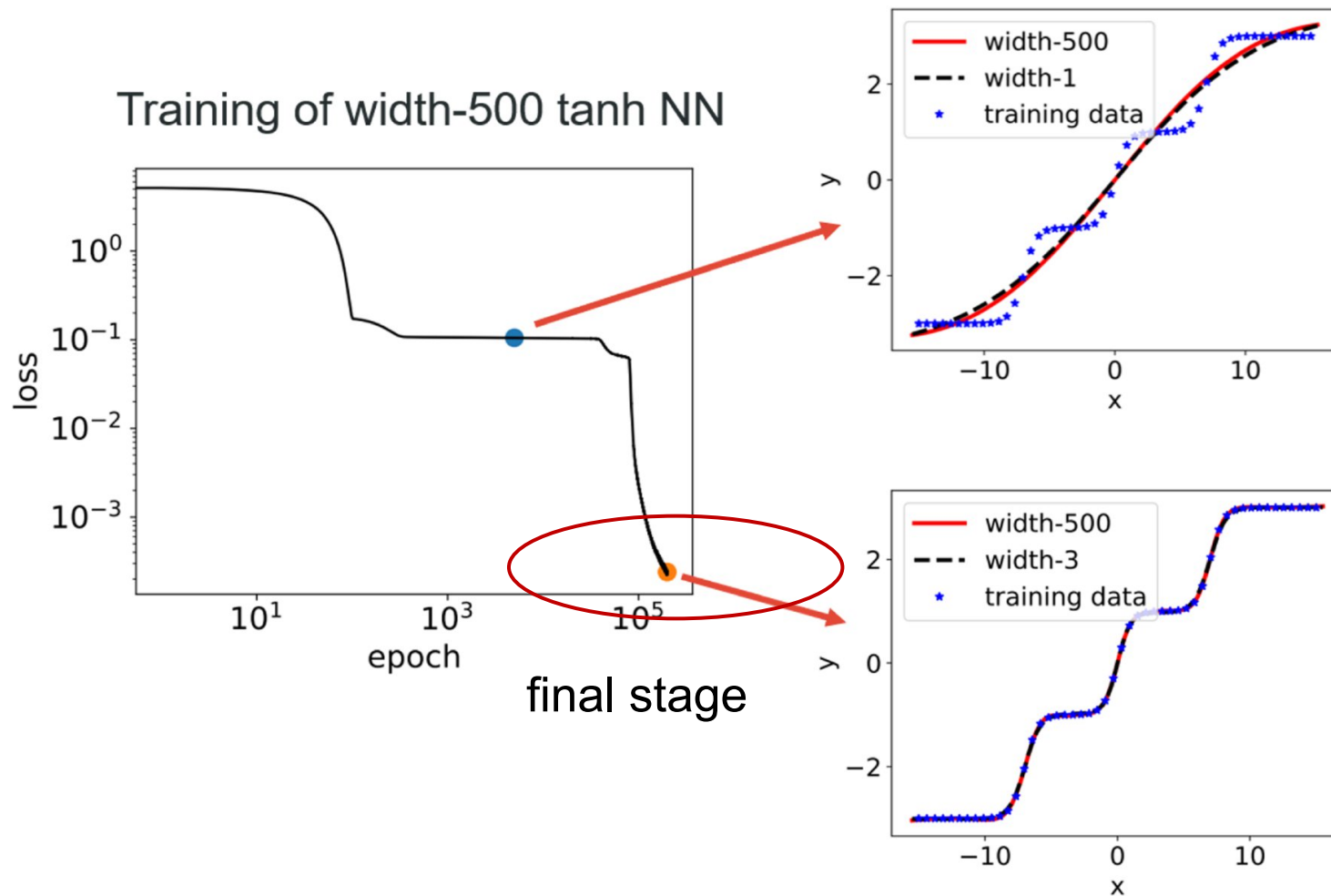
Hierarchical structure for two-layer NN



Final condensation



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: $f^* = \bar{a} \sigma(\bar{w}^T x)$
- Target Set $Q^* = F^{-1}(f^*)$ generally consists of three “branches” (sets)
 - (a) $Q_1 = \{(a_k, w_k)_{k=1}^2 : w_1 = w_2 = \bar{w}, a_1 + a_2 = \bar{a}\}$.
 - (b) $Q_2 = \{(a_k, w_k)_{k=1}^2 : w_1 = \bar{w}, a_1 = \bar{a}, a_2 = 0\}$.
 - (c) $Q_3 = \{(a_k, w_k)_{k=1}^2 : w_2 = \bar{w}, a_2 = \bar{a}, a_1 = 0\}$.

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

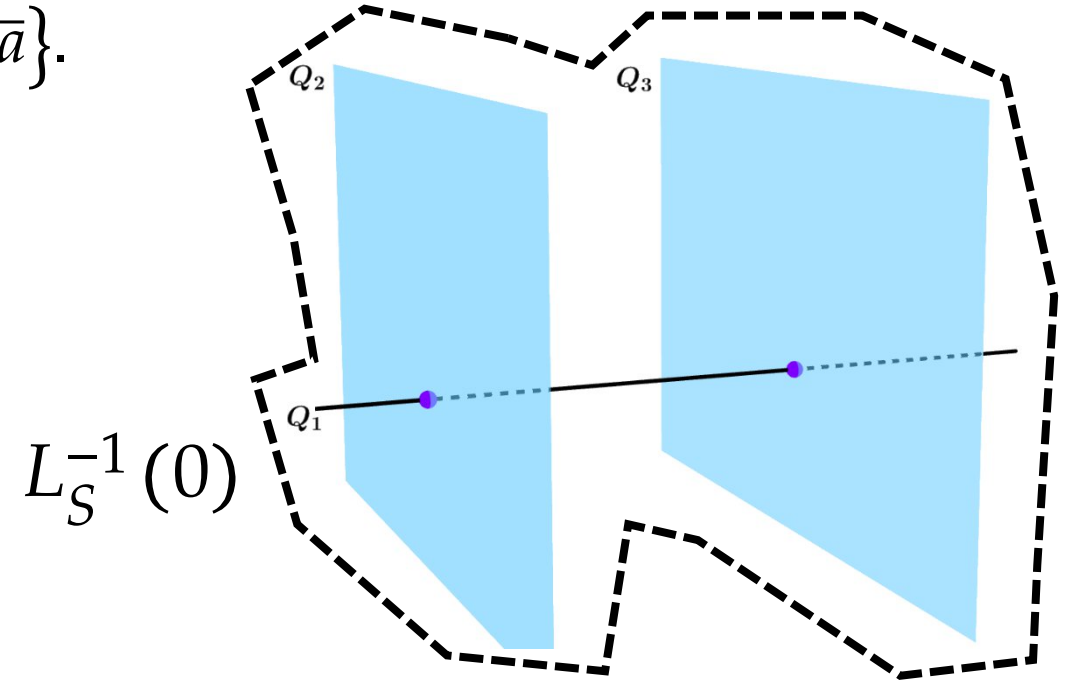
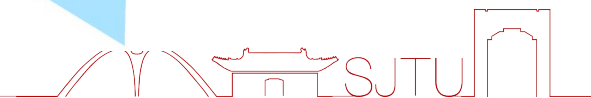
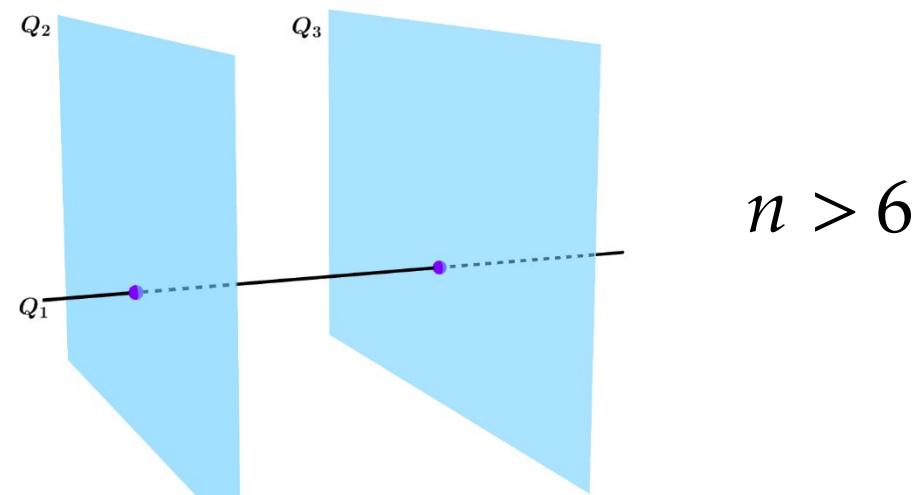
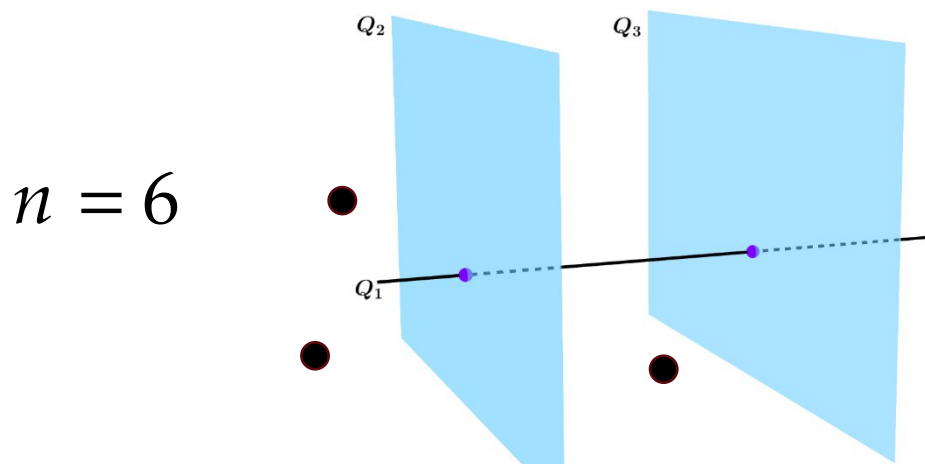
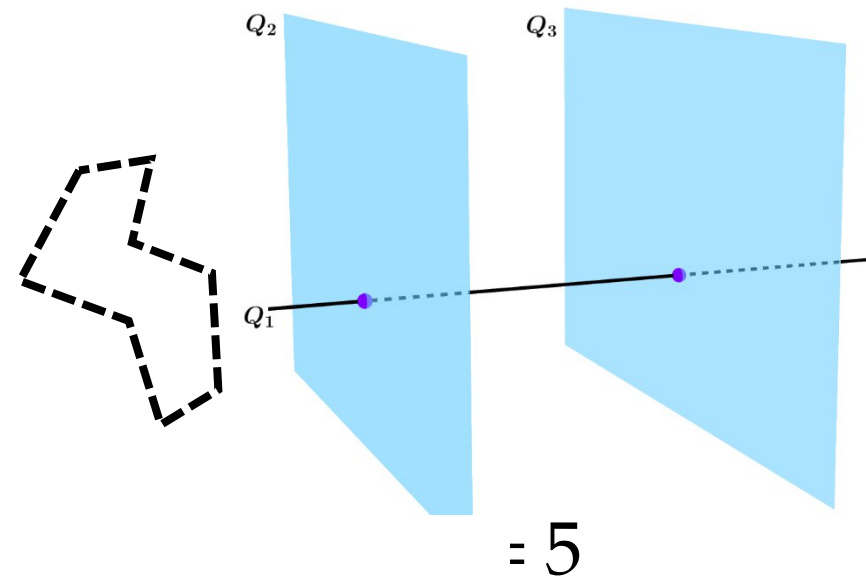
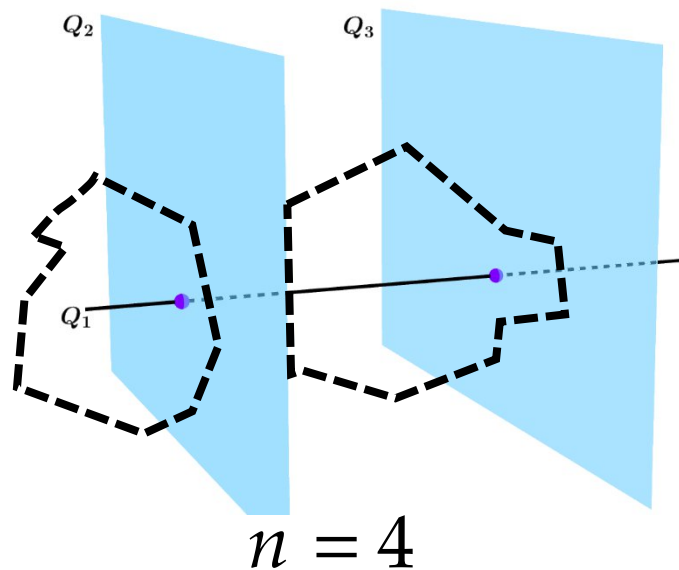
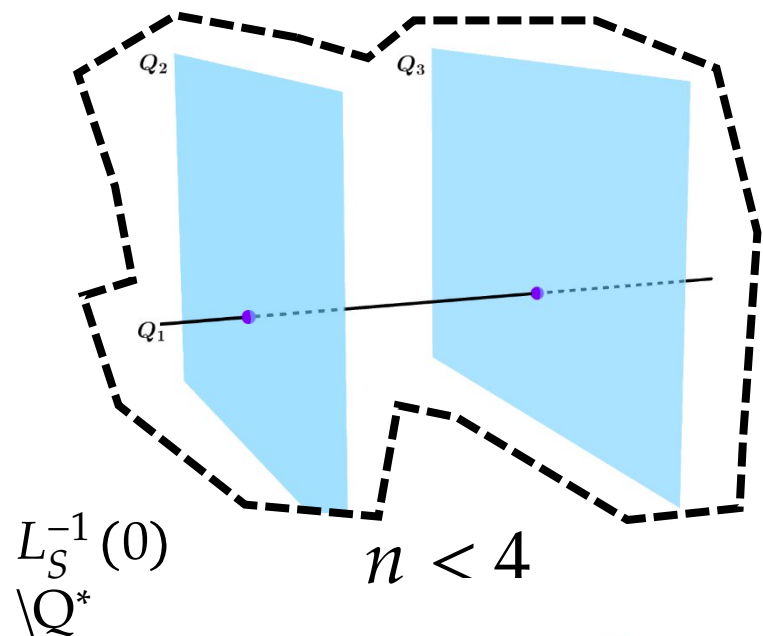


Illustration of Q^1, Q^2, Q^3





Geometry of global minima for final condensation





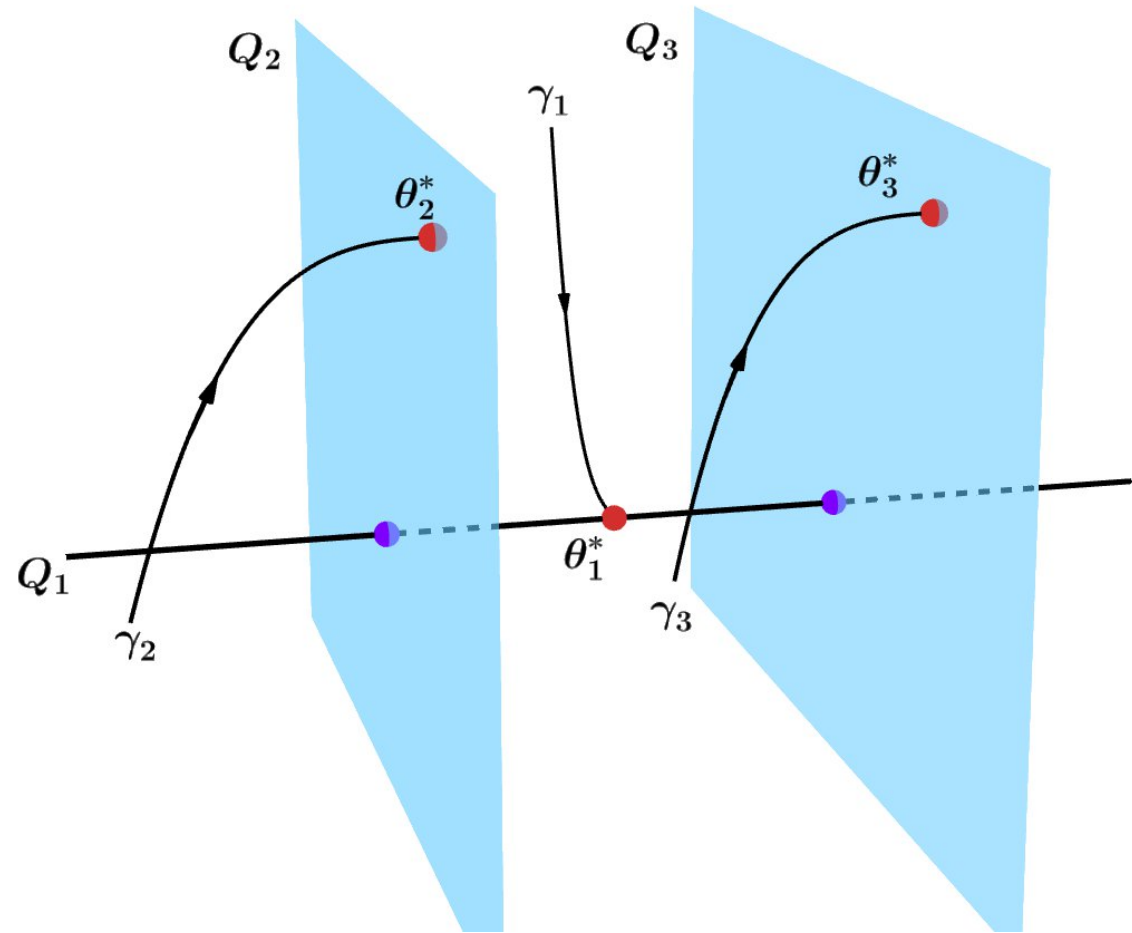
Typical convergence rate for final condensation



Gradient flows near Q^* :

γ_1 : sublinear rate;

γ_2, γ_3 : linear rate.

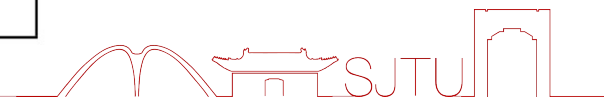




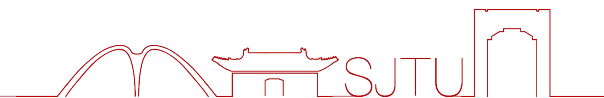
Stability of target branches underlies final condensation

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
$\geq m+m_0d$					✓
\vdots				...	\vdots
$\geq r+(m+m_0-r)d$			✓	...	✓
\vdots		...	\vdots	...	\vdots
$\geq m_0+md$	✓	...	✓	...	✓
$> (d+1)m$	✓*				
✓*: any point in Q^* is recovery stable					

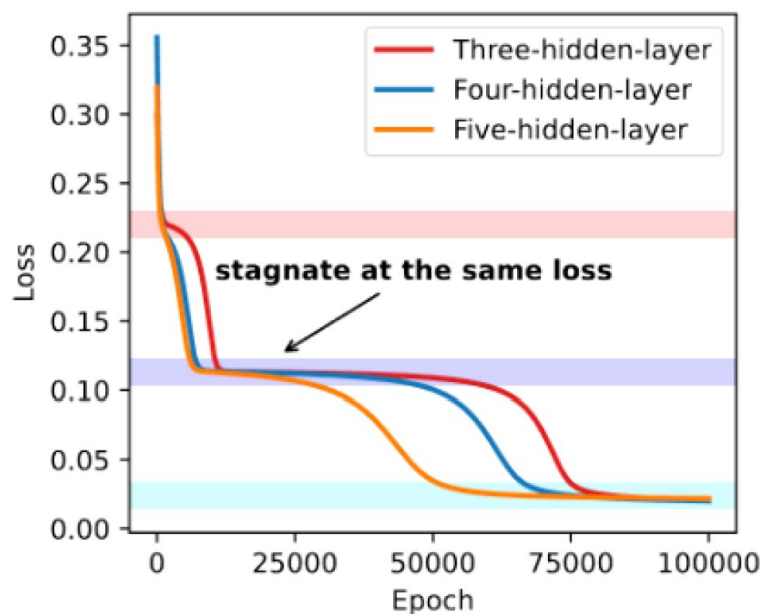


Embedding principle in depth

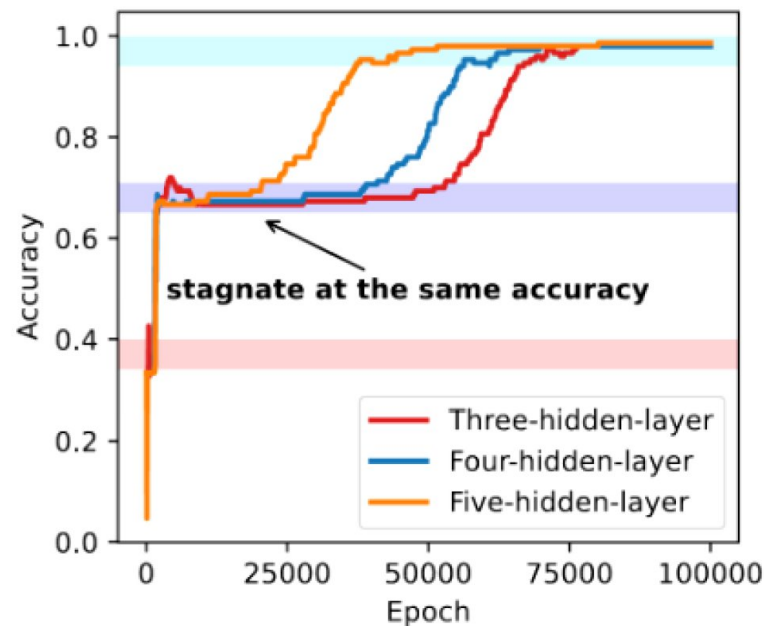




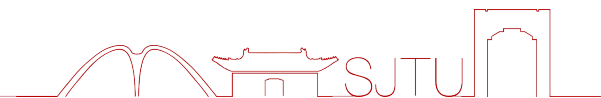
Phenomenon: training similarity in depth



(a) Loss (Iris)



(b) Accuracy (Iris)

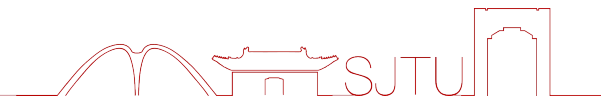




Embedding Principle in depth

Embedding principle in depth

the loss landscape of any network “contains” all critical points of all shallower networks.





Key to the Proof

 Goal: discover an embedding operator

Discover a mapping $\mathcal{T} : \mathbb{R}^{M_{\text{shal}}} \rightarrow \mathbb{R}^{M_{\text{deep}}}$ such that for any $\theta_{\text{deep}} \in \mathcal{T}(\theta_{\text{shal}})$, we have

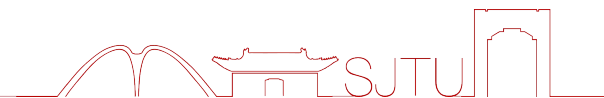
(1) **output preserving:**

$$f_{\theta_{\text{deep}}}(x) = f_{\theta_{\text{shal}}}(x), \forall x \in S_x.$$

(2) **criticality preserving:**

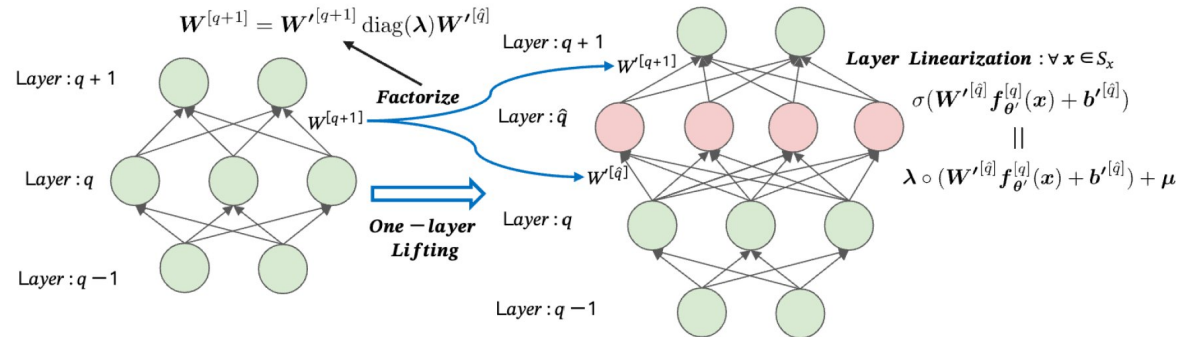
if θ_{shal} is a critical point, then θ_{deep} is also a critical point.

$$\nabla R_S(\theta_{\text{shal}}) = \mathbf{0} \implies \nabla R_S(\theta_{\text{deep}}) = \mathbf{0}.$$



One-layer lifting

Assumption 1. Activation function σ has at least non-constant **linear piece**, e.g., ReLU, Leaky ReLU, ELU, etc.



1. Layer linearization condition

$$\sigma(W'^{[\hat{q}]} f_{\theta'}^{[q]}(x) + b'^{[\hat{q}]}) = \lambda \circ (W'^{[\hat{q}]} f_{\theta'}^{[q]}(x) + b'^{[\hat{q}]}) + \mu, \forall x \in S_x$$

2. Output preserving condition

$$\begin{cases} W'^{[q+1]} \text{diag}(\lambda) W'^{[\hat{q}]} = W^{[q+1]}, \\ W'^{[q+1]} \text{diag}(\lambda) b'^{[\hat{q}]} + W'^{[q+1]} \mu + b'^{[q+1]} = \underline{b^{[q+1]}}. \end{cases}$$



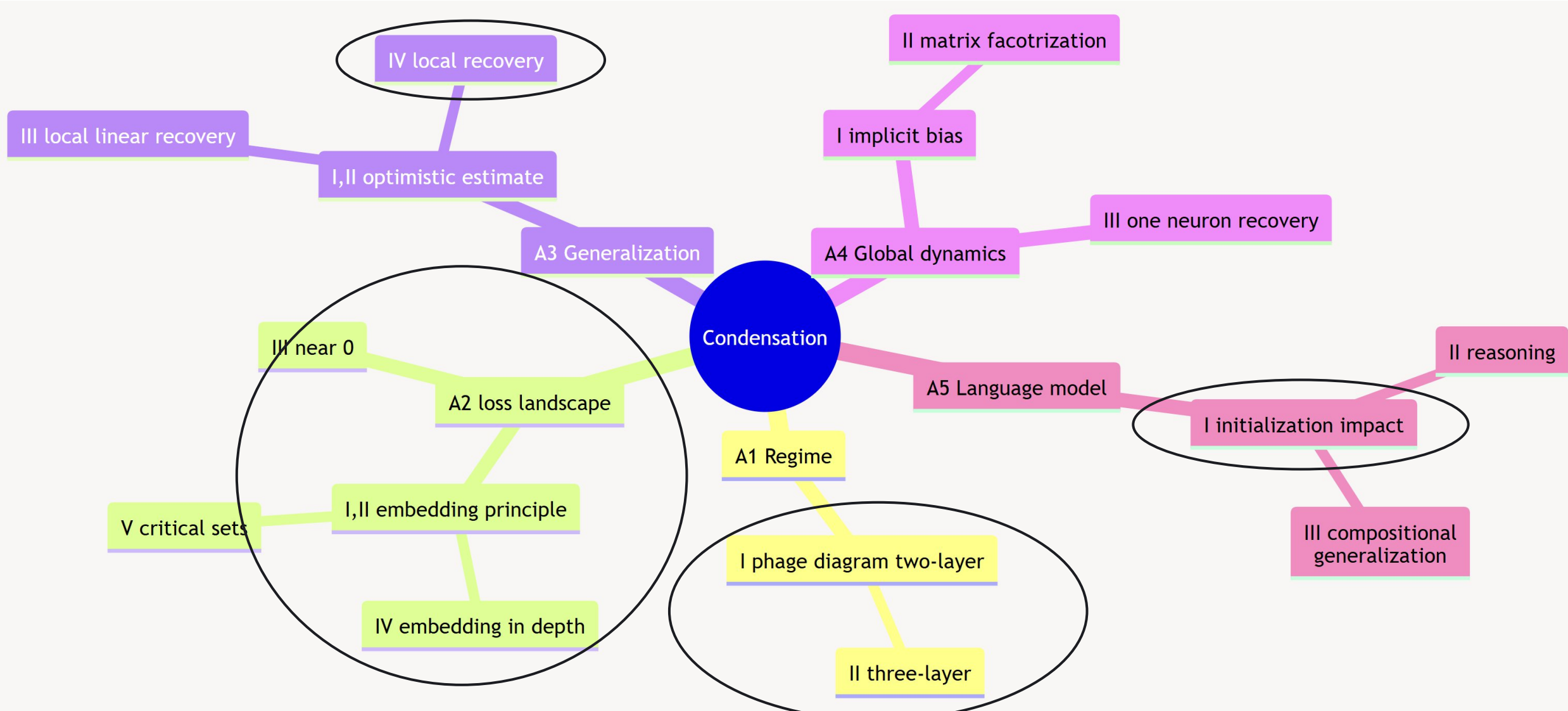
Embedding principle in depth [Bai et al.]

Theorem 1 (embedding principle in depth). Given any NN $\left(\{m_l\}_{l=0}^L\right)$ and data S , for any θ'_c of any shallower NN' $\left(\{m'_l\}_{l=0}^{L'}\right)$ satisfying $\nabla_{\theta} R_S(\theta'_c) = 0$, there exists parameter θ_c in the loss landscape of NN $\left(\{m_l\}_{l=0}^L\right)$ satisfying the following conditions:

- (i) **Output Preserving:** $f_{\theta_c}(x) = f_{\theta'_c}(x)$ for $x \in S_x$;
- (ii) **Criticality Preserving:** $\nabla_{\theta} R_S(\theta_c) = 0$.



Condensation





- ① Why does small initialization lead to (initial) condensation in neural networks?
- ② How does the empirical risk landscape transform as the width of a neural network increases?
- ③ What are the reasons why wider neural networks are often easier to optimize than narrower ones?
- ④ Is it possible for neural networks to achieve zero generalization error for a target function under overparameterization?



Thanks!

饮水思源 爱国荣校