

Supplement to: Embedding Principle of Loss Landscape of Deep Neural Networks

Yaoyu Zhang^{1,2,*}, Zhongwang Zhang¹ †, Tao Luo¹, Zhi-Qin John Xu^{1‡}

¹ School of Mathematical Sciences, Institute of Natural Sciences, MOE-LSC and
Qing Yuan Research Institute, Shanghai Jiao Tong University

² Shanghai Center for Brain Science and Brain-Inspired Technology
{zhyy.sjtu, 0123zzw666, luotao41, xuzhiqin}@sjtu.edu.cn.

A Appendix

Lemma (Lemma 1). *Given a L -layer ($L \geq 2$) fully-connected neural network with width (m_0, \dots, m_L) , for any network parameters $\theta = (\mathbf{W}^{[1]}, \mathbf{b}^{[1]}, \dots, \mathbf{W}^{[L]}, \mathbf{b}^{[L]})$ and for any $l \in [L-1]$, $s \in [m_l]$, we have the expressions for $\theta' := \mathcal{T}_{l,s}^\alpha(\theta)$ (see Fig. S2 for an illustration)*

(i) feature vectors in $\mathbf{F}_{\theta'}$: $\mathbf{f}_{\theta'}^{[l']} = \mathbf{f}_\theta^{[l']}$, $l' \neq l$ and $\mathbf{f}_{\theta'}^{[l]} = \left[(\mathbf{f}_\theta^{[l]})^\top, (\mathbf{f}_\theta^{[l]})_s \right]^\top$;

(ii) feature gradients in $\mathbf{G}_{\theta'}$: $\mathbf{g}_{\theta'}^{[l']} = \mathbf{g}_\theta^{[l']}$, $l' \neq l$ and $\mathbf{g}_{\theta'}^{[l]} = \left[(\mathbf{g}_\theta^{[l]})^\top, (\mathbf{g}_\theta^{[l]})_s \right]^\top$;

(iii) error vectors in $\mathbf{Z}_{\theta'}$: $\mathbf{z}_{\theta'}^{[l']} = \mathbf{z}_\theta^{[l']}$, $l' \neq l$

and $\mathbf{z}_{\theta'}^{[l]} = \left[(\mathbf{z}_\theta^{[l]})^\top_{[1:s-1]}, (1-\alpha)(\mathbf{z}_\theta^{[l]})_s, (\mathbf{z}_\theta^{[l]})^\top_{[s+1:m_l]}, \alpha(\mathbf{z}_\theta^{[l]})_s \right]^\top$.

Proof. (i) By the construction of θ' , it is clear that $\mathbf{f}_{\theta'}^{[l']} = \mathbf{f}_\theta^{[l']}$ for any $l' \in [l-1]$. Then

$$\mathbf{f}_{\theta'}^{[l]} = \sigma \circ \left(\begin{bmatrix} \mathbf{W}^{[l]} \\ \mathbf{W}_{s,[1:m_{l-1}]}^{[l]} \end{bmatrix} \mathbf{f}_\theta^{[l-1]} + \begin{bmatrix} \mathbf{b}^{[l]} \\ \mathbf{b}_s^{[l]} \end{bmatrix} \right) = \begin{bmatrix} \mathbf{f}_\theta^{[l]} \\ (\mathbf{f}_\theta^{[l]})_s \end{bmatrix}. \quad (1)$$

Note that

$$\alpha \left[\mathbf{0}_{m_{l+1} \times (s-1)}, -\mathbf{W}_{[1:m_{l+1}],s}^{[l+1]}, \mathbf{0}_{m_{l+1} \times (m_l - s)}, \mathbf{W}_{[1:m_{l+1}],s}^{[l+1]} \right] \begin{bmatrix} \mathbf{f}_\theta^{[l]} \\ (\mathbf{f}_\theta^{[l]})_s \end{bmatrix} = \mathbf{0}_{m_{l+1} \times 1}.$$

Thus

$$\mathbf{f}_{\theta'}^{[l+1]} = \sigma \circ \left(\begin{bmatrix} \mathbf{W}^{[l+1]} \\ \mathbf{W}_{m_{l+1} \times 1}^{[l+1]} \end{bmatrix} \begin{bmatrix} \mathbf{f}_\theta^{[l]} \\ (\mathbf{f}_\theta^{[l]})_s \end{bmatrix} + \mathbf{0}_{m_{l+1} \times 1} + \mathbf{b}^{[l+1]} \right) = \mathbf{f}_\theta^{[l+1]}. \quad (2)$$

Next, by the construction of θ' again, it is clear that $\mathbf{f}_{\theta'}^{[l']} = \mathbf{f}_\theta^{[l']}$ for any $l' \in [l+1:L]$.

(ii) The results for feature gradients $\mathbf{g}_{\theta'}^{[l']}$, $l' \in [L]$ can be calculated in a similar way.

(iii) By the backpropagation and the above facts in (i), we have $\mathbf{z}_{\theta'}^{[L]} = \nabla \ell(\mathbf{f}_{\theta'}^{[L]}, \mathbf{y}) = \nabla \ell(\mathbf{f}_\theta^{[L]}, \mathbf{y}) = \mathbf{z}_\theta^{[L]}$.

*Corresponding author: zhyy.sjtu@sjtu.edu.cn.

†Part of this work is done when ZZ was an undergraduate student of Zhiyuan Honors Program at Shanghai Jiao Tong University.

‡Corresponding author: xuzhiqin@sjtu.edu.cn.

Recalling the recurrence relation for $l' \in [l + 1 : L - 1]$, then we recursively obtain the following equality for l' from $L - 1$ down to $l + 1$:

$$\mathbf{z}_{\theta'}^{[l']} = (\mathbf{W}^{[l'+1]})^\top \mathbf{z}_{\theta'}^{[l'+1]} \circ \mathbf{g}_{\theta'}^{[l'+1]} = (\mathbf{W}^{[l'+1]})^\top \mathbf{z}_{\theta}^{[l'+1]} \circ \mathbf{g}_{\theta}^{[l'+1]} = \mathbf{z}_{\theta}^{[l']}. \quad (3)$$

Next,

$$\begin{aligned} \mathbf{z}_{\theta'}^{[l]} &= \left(\left[\mathbf{W}^{[m_{l+1}]} \right], \mathbf{0}_{m_{l+1} \times 1} \right) + \alpha \left[\mathbf{0}_{m_{l+1} \times (s-1)}, -\mathbf{W}_{[1:m_{l+1}],s}^{[l+1]}, \mathbf{0}_{m_{l+1} \times (m_l-s)}, \mathbf{W}_{[1:m_{l+1}],s}^{[l+1]} \right]^\top \mathbf{z}_{\theta}^{[l+1]} \circ \mathbf{g}_{\theta}^{[l+1]} \\ &= \begin{bmatrix} \mathbf{z}_{\theta}^{[l]} \\ 0 \end{bmatrix} + \begin{bmatrix} \mathbf{0}_{m_{l+1} \times (s-1)} \\ -\alpha(\mathbf{z}_{\theta}^{[l]})_s \\ \mathbf{0}_{m_{l+1} \times (m_l-s)} \\ \alpha(\mathbf{z}_{\theta}^{[l]})_s \end{bmatrix} \\ &= \left[(\mathbf{z}_{\theta}^{[l]})_{[1:s-1]}^\top, (1-\alpha)(\mathbf{z}_{\theta}^{[l]})_s, (\mathbf{z}_{\theta}^{[l]})_{[s+1:m_l]}^\top, \alpha(\mathbf{z}_{\theta}^{[l]})_s \right]^\top. \end{aligned} \quad (4)$$

Finally,

$$\begin{aligned} \mathbf{z}_{\theta'}^{[l-1]} &= \left[(\mathbf{W}^{[l]})^\top, (\mathbf{W}^{[l]})_{s,[1:m_{l-1}]}^\top \right] \left(\begin{bmatrix} \mathbf{z}_{\theta}^{[l]} \\ 0 \end{bmatrix} + \begin{bmatrix} \mathbf{0}_{m_{l+1} \times (s-1)} \\ -\alpha(\mathbf{z}_{\theta}^{[l]})_s \\ \mathbf{0}_{m_{l+1} \times (m_l-s)} \\ \alpha(\mathbf{z}_{\theta}^{[l]})_s \end{bmatrix} \right) \circ \begin{bmatrix} \mathbf{g}_{\theta}^{[l]} \\ (\mathbf{g}_{\theta}^{[l]})_s \end{bmatrix} \\ &= (\mathbf{W}^{[l]})^\top \mathbf{z}_{\theta}^{[l]} \circ \mathbf{g}_{\theta}^{[l]} + \mathbf{0}_{m_{l-1} \times 1} \\ &= \mathbf{z}_{\theta}^{[l-1]}. \end{aligned} \quad (5)$$

This with the recurrence relation again leads to $\mathbf{z}_{\theta'}^{[l']} = \mathbf{z}_{\theta}^{[l']}$ for all $l' \in [1 : l - 1]$. \square

Proposition (Proposition 1: one-step embedding preserves network properties). *Given a L -layer ($L \geq 2$) fully-connected neural network with width (m_0, \dots, m_L) , for any network parameters $\theta = (\mathbf{W}^{[1]}, \mathbf{b}^{[1]}, \dots, \mathbf{W}^{[L]}, \mathbf{b}^{[L]})$ and for any $l \in [L - 1]$, $s \in [m_l]$, the following network properties are preserved for $\theta' = \mathcal{T}_{l,s}^\alpha(\theta)$:*

(i) *output function is preserved: $f_{\theta'}(\mathbf{x}) = f_{\theta}(\mathbf{x})$ for all \mathbf{x} ;*

(ii) *empirical risk is preserved: $R_S(\theta') = R_S(\theta)$;*

(iii) *the sets of features are preserved: $\left\{ (\mathbf{f}_{\theta'}^{[l]})_i \right\}_{i \in [m_{l+1}]} = \left\{ (\mathbf{f}_{\theta}^{[l]})_i \right\}_{i \in [m_l]}$ and $\left\{ (\mathbf{f}_{\theta'}^{[l']})_i \right\}_{i \in [m_{l'}]} = \left\{ (\mathbf{f}_{\theta}^{[l']})_i \right\}_{i \in [m_{l'}]}$ for $l' \in [L] \setminus \{l\}$;*

Proof. The properties (i)–(iii) are direct consequences of Lemma 1. \square

Theorem (Theorem 1: criticality preserving). *Given a L -layer ($L \geq 2$) fully-connected neural network with width (m_0, \dots, m_L) , for any network parameters $\theta = (\mathbf{W}^{[1]}, \mathbf{b}^{[1]}, \dots, \mathbf{W}^{[L]}, \mathbf{b}^{[L]})$ and for any $l \in [L - 1]$, $s \in [m_l]$, if $\nabla_{\theta} R_S(\theta) = \mathbf{0}$, then $\nabla_{\theta} R_S(\theta') = \mathbf{0}$.*

Proof. Gradient of loss with respect to network parameters of each layer can be computed from \mathbf{F} , \mathbf{G} , and \mathbf{Z} as follows

$$\begin{aligned} \nabla_{\mathbf{W}^{[l']}} R_S(\theta) &= \nabla_{\mathbf{W}^{[l']}} \mathbb{E}_S \ell(\mathbf{f}_{\theta}(\mathbf{x}), \mathbf{y}) = \mathbb{E}_S \left(\mathbf{z}_{\theta}^{[l']} \circ \mathbf{g}_{\theta}^{[l']} (\mathbf{f}_{\theta}^{[l'-1]})^\top \right), \\ \nabla_{\mathbf{b}^{[l']}} R_S(\theta) &= \nabla_{\mathbf{b}^{[l']}} \mathbb{E}_S \ell(\mathbf{f}_{\theta}(\mathbf{x}), \mathbf{y}) = \mathbb{E}_S (\mathbf{z}_{\theta}^{[l']} \circ \mathbf{g}_{\theta}^{[l']}). \end{aligned}$$

Then, by Lemma 1, we have $\nabla_{\mathbf{W}^{[l']}} R_S(\theta') = \nabla_{\mathbf{W}^{[l']}} R_S(\theta) = \mathbf{0}$ for $l' \neq l, l+1$ and $\nabla_{\mathbf{b}^{[l']}} R_S(\theta') = \nabla_{\mathbf{b}^{[l']}} R_S(\theta) = \mathbf{0}$ for $l' \neq l$. Also, for any $j \in [m_{l+1}], k \in [m_l]$, since $(\mathbf{z}_{\theta'}^{[l+1]})_j = (\mathbf{z}_{\theta}^{[l+1]})_j$, $(\mathbf{g}_{\theta'}^{[l+1]})_j = (\mathbf{g}_{\theta}^{[l+1]})_j$, and $(\mathbf{f}_{\theta'}^{[l]})_k = (\mathbf{f}_{\theta}^{[l]})_k$, $(\mathbf{f}_{\theta'}^{[l]})_{m_l+1} = (\mathbf{f}_{\theta}^{[l]})_s$, we obtain

$$\begin{aligned} \nabla_{\mathbf{W}_{j,k}^{[l+1]}} R_S(\theta') &= \nabla_{\mathbf{W}_{j,k}^{[l+1]}} R_S(\theta) = 0, \\ \nabla_{\mathbf{W}_{j,m_l+1}^{[l+1]}} R_S(\theta') &= \nabla_{\mathbf{W}_{j,s}^{[l+1]}} R_S(\theta) = 0. \end{aligned}$$

Similarly, for any $j \in [m_l] \setminus \{s\}, k \in [m_{l-1}]$, we have

$$\begin{aligned}
\nabla_{\mathbf{W}_{j,k}^{[l]}} R_S(\boldsymbol{\theta}') &= \nabla_{\mathbf{W}_{j,k}^{[l]}} R_S(\boldsymbol{\theta}) = 0, \\
\nabla_{\mathbf{b}_j^{[l]}} R_S(\boldsymbol{\theta}') &= \nabla_{\mathbf{b}_j^{[l]}} R_S(\boldsymbol{\theta}) = 0, \\
\nabla_{\mathbf{W}_{s,k}^{[l]}} R_S(\boldsymbol{\theta}') &= (1 - \alpha) \nabla_{\mathbf{W}_{s,k}^{[l]}} R_S(\boldsymbol{\theta}) = 0, \\
\nabla_{\mathbf{W}_{m_l+1,k}^{[l]}} R_S(\boldsymbol{\theta}') &= \alpha \nabla_{\mathbf{W}_{s,k}^{[l]}} R_S(\boldsymbol{\theta}) = 0, \\
\nabla_{\mathbf{b}_s^{[l]}} R_S(\boldsymbol{\theta}') &= (1 - \alpha) \nabla_{\mathbf{b}_s^{[l]}} R_S(\boldsymbol{\theta}) = 0, \\
\nabla_{\mathbf{b}_{m_l+1}^{[l]}} R_S(\boldsymbol{\theta}') &= \alpha \nabla_{\mathbf{b}_s^{[l]}} R_S(\boldsymbol{\theta}) = 0.
\end{aligned}$$

Collecting all the above equalities, we have $\nabla_{\boldsymbol{\theta}} R_S(\boldsymbol{\theta}') = \mathbf{0}$. \square

Lemma (Lemma 2: increment of the degree of degeneracy). *Given a L -layer ($L \geq 2$) fully-connected neural network with width (m_0, \dots, m_L) , if there exists $l \in [L - 1]$, $s \in [m_l]$, and a q -dimensional differential manifold \mathcal{M} consisting of critical points of R_S such that for any $\boldsymbol{\theta} \in \mathcal{M}$, $\mathbf{W}_{[1:m_{l+1}],s}^{[l+1]} \neq \mathbf{0}$, then $\mathcal{M}' := \{\mathcal{T}_{l,s}^\alpha(\boldsymbol{\theta}) | \boldsymbol{\theta} \in \mathcal{M}, \alpha \in \mathbb{R}\}$ is a $(d + 1)$ -dimensional differential manifold consisting of critical points for the corresponding L -layer fully-connected neural network with width $(m_0, \dots, m_{l-1}, m_l + 1, m_{l+1}, \dots, m_L)$.*

Proof. For any $\boldsymbol{\theta} \in \mathcal{M}$, let $\{\mathbf{e}_i(\boldsymbol{\theta})\}_{i=1}^d$ be a basis of its tangent space $T_{\boldsymbol{\theta}}\mathcal{M}$. Then for any $\alpha \in \mathbb{R}$, the tangent space of $\boldsymbol{\theta}' = \mathcal{T}_{l,s}^\alpha(\boldsymbol{\theta}) \in \mathcal{M}'$ is spanned by $\{\mathcal{T}_{l,s}(\mathbf{e}_1(\boldsymbol{\theta})), \dots, \mathcal{T}_{l,s}(\mathbf{e}_d(\boldsymbol{\theta})), \mathcal{V}_{l,s}(\boldsymbol{\theta})\}$. Since $\mathcal{T}_{l,s}$ is linear and injective, $\{\mathcal{T}_{l,s}(\mathbf{e}_i(\boldsymbol{\theta}))\}_{i=1}^d$ is also a linearly independent set. Moreover, since $\mathbf{W}_{[1:m_{l+1}],m_{l+1}}^{[l+1]} = \mathbf{0}$ for any vector in parameter space applied with $\mathcal{T}_{l,s}$, then, $\{\mathcal{T}_{l,s}(\mathbf{e}_1(\boldsymbol{\theta})), \dots, \mathcal{T}_{l,s}(\mathbf{e}_d(\boldsymbol{\theta})), \mathcal{V}_{l,s}(\boldsymbol{\theta})\}$ are independent if and only if $\mathcal{V}_{l,s}(\boldsymbol{\theta}) \neq \mathbf{0}$, i.e., $\mathbf{W}_{[1:m_{l+1}],s}^{[l+1]} \neq \mathbf{0}$. \square

Remark 1. *The requirement that \mathcal{M} is a q -dimensional differential manifold can be relaxed to that \mathcal{M} is a q -dimensional topological manifold. In the latter case, \mathcal{M}' is a $(d + 1)$ -dimensional topological manifold.*

Theorem (Theorem 2: degeneracy of embedded critical points). *Consider two L -layer ($L \geq 2$) fully-connected neural networks $\text{NN}_A(\{m_i\}_{i=0}^L)$ and $\text{NN}_B(\{m'_i\}_{i=0}^L)$ which is K -neuron wider than NN_A . Suppose that the critical point $\boldsymbol{\theta}_A = (\mathbf{W}^{[1]}, \mathbf{b}^{[1]}, \dots, \mathbf{W}^{[L]}, \mathbf{b}^{[L]})$ satisfy $\mathbf{W}^{[l]} \neq \mathbf{0}$ for each layer $l \in [L]$. Then the parameters $\boldsymbol{\theta}_A$ of NN_A can be critically embedded to a K -dimensional critical affine subspace $\mathcal{M}_B = \{\boldsymbol{\theta}_B + \sum_{i=1}^K \alpha_i \mathbf{v}_i | \alpha_i \in \mathbb{R}\}$ of loss landscape of NN_B . Here $\boldsymbol{\theta}_B = (\prod_{i=1}^K \mathcal{T}_{l_i, s_i})(\boldsymbol{\theta}_A)$ and $\mathbf{v}_i = \mathcal{T}_{l_K, s_K} \cdots \mathcal{V}_{l_i, s_i} \cdots \mathcal{T}_{l_1, s_1} \boldsymbol{\theta}_A$.*

Proof. The assumption $\mathbf{W}^{[l]} \neq \mathbf{0}, l \in [L]$ implies the existence of non-silent neurons, i.e., existing $s \in [m_l]$ such that $\mathbf{W}_{[1:m_{l+1}],s}^{[l+1]} \neq \mathbf{0}$, for any $l \in [L - 1]$ with $m'_l > m_l$.

In this proof, we misuse notation and denote $m_l = m_l(\boldsymbol{\theta})$ for the width of the l -th layer for any fully-connected neural network with parameters $\boldsymbol{\theta}$. For such a general network with parameters $\boldsymbol{\theta}$, we introduce the following operator. Given an index set J and for any $l \in [L]$, $s \in [m_l]$, we define

$$\mathcal{V}_{l,s,J}(\boldsymbol{\theta}) = \left(\mathbf{0}_{m_0 \times m_1}, \dots, \left[\mathbf{0}_{m_{l+1} \times (s-1)}, - \sum_{j \in J} \mathbf{W}_{[1:m_{l+1}],j}^{[l+1]}, \mathbf{0}_{m_{l+1} \times (m_l - s)}, \sum_{j \in J} \mathbf{W}_{[1:m_{l+1}],j}^{[l+1]} \right], \dots \right).$$

Clearly, $\mathcal{V}_{l,s,J} = \sum_{j \in J} \mathcal{V}_{l,s,\{j\}}$. If for all $j \in J$, $\mathbf{W}_{j,[1:m_l-1]}^{[l]} = \mathbf{W}_{s,[1:m_l-1]}^{[l]}$, $\mathbf{W}_{[1:m_l+1],j}^{[l+1]} = \beta_j \sum_{j' \in J} \mathbf{W}_{[1:m_l+1],j'}^{[l+1]}$ and $\sum_{j' \in J} \mathbf{W}_{[1:m_l+1],j'}^{[l+1]} \neq \mathbf{0}$, then for $\beta_s \neq 0$, we have

$$\begin{aligned} \mathcal{T}_{l,s,J}^\alpha \boldsymbol{\theta} &= (\mathcal{T}_{l,s} + \alpha \mathcal{V}_{l,s,J}) \boldsymbol{\theta} \\ &= (\mathcal{T}_{l,s} + \alpha \sum_{j \in J} \mathcal{V}_{l,s,j}) \boldsymbol{\theta} \\ &= (\mathcal{T}_{l,s} + \alpha \sum_{j \in J} \frac{\beta_j}{\beta_s} \mathcal{V}_{l,s}) \boldsymbol{\theta} \\ &= \mathcal{T}_{l,s}^{\alpha'} \boldsymbol{\theta}, \end{aligned}$$

where $\alpha' = \alpha \sum_{j \in J} \frac{\beta_j}{1 - \beta_j}$, is simply the one-step critical embedding. We can extend this to the case of $\beta_s = 0$.

Then, for $J' = J \cup \{m_l + 1\}$, $l' = l$, $s' \in J$,

$$\begin{aligned} \mathcal{V}_{l',s',J'} \mathcal{V}_{l,s,J} &= (\mathcal{V}_{l',s',J} + \mathcal{V}_{l',s',m_l+1}) \mathcal{V}_{l,s,J} \\ &= \mathcal{V}_{l',s',J} \mathcal{V}_{l,s,J} + \mathcal{V}_{l',s',m_l+1} \mathcal{V}_{l,s,J} \\ &= \mathcal{V}_{l',s',s} \mathcal{V}_{l,s,J} + \mathcal{V}_{l',s',m_l+1} \mathcal{V}_{l,s,J} \\ &= \mathbf{0}. \end{aligned}$$

In general, we have

$$\begin{aligned} \mathcal{V}_{l',s',J'} \prod_{i=1}^N \mathcal{T}_{l'_i,s'_i} \mathcal{V}_{l,s,J} &= (\mathcal{V}_{l',s',J} + \mathcal{V}_{l',s',m_l+1}) \prod_{i=1}^N \mathcal{T}_{l'_i,s'_i} \mathcal{V}_{l,s,J} \\ &= \mathcal{V}_{l',s',J} \prod_{i=1}^N \mathcal{T}_{l'_i,s'_i} \mathcal{V}_{l,s,J} + \mathcal{V}_{l',s',m_l+1} \prod_{i=1}^N \mathcal{T}_{l'_i,s'_i} \mathcal{V}_{l,s,J} \\ &= \mathcal{V}_{l',s',s} \prod_{i=1}^N \mathcal{T}_{l'_i,s'_i} \mathcal{V}_{l,s,J} + \mathcal{V}_{l',s',m_l+1} \prod_{i=1}^N \mathcal{T}_{l'_i,s'_i} \mathcal{V}_{l,s,J} \\ &= \mathbf{0}. \end{aligned}$$

For $l' \neq l$ or $s' \notin J$, obviously we have $\mathcal{V}_{l',s',J'} \mathcal{V}_{l,s,J} = \mathbf{0}$ and $\mathcal{V}_{l',s',J'} \prod_{i=1}^N \mathcal{T}_{l'_i,s'_i} \mathcal{V}_{l,s,J} = \mathbf{0}$.

Now we are ready to prove the lemma. Let $J_i = \{s_i\} \cup \{m_l + \#\{i | l_i = l, i \in [j]\} | l_j = l, s_j = s, j \in [K]\}$, where $\#$ indicates number of elements in a set. Then

$$\begin{aligned} \prod_{i=1}^K \mathcal{T}_{l_i,s_i,J_i}^{\alpha_i} &= \prod_{i=1}^K (\mathcal{T}_{l_i,s_i} + \alpha_i \mathcal{V}_{l_i,s_i,J_i}) \\ &= \prod_{i=1}^K \mathcal{T}_{l_i,s_i} + \sum_{i=1}^K \alpha_i \mathcal{T}_{l_K,s_K} \cdots \mathcal{V}_{l_i,s_i,J_i} \cdots \mathcal{T}_{l_1,s_1}, \\ &= \prod_{i=1}^K \mathcal{T}_{l_i,s_i} + \sum_{i=1}^K \alpha_i \mathcal{T}_{l_K,s_K} \cdots \mathcal{V}_{l_i,s_i} \cdots \mathcal{T}_{l_1,s_1}, \end{aligned}$$

which is a critical embedding for any $[\alpha_i]_{i=1}^K \in \mathbb{R}^K$. This completes the proof. \square

A.1 Trivial critical transforms

In general, neuron-index permutation among the same layer is a trivial criticality invariant transform because of the layer-wise intrinsic symmetry of DNN models. Therefore, any critical point/manifold

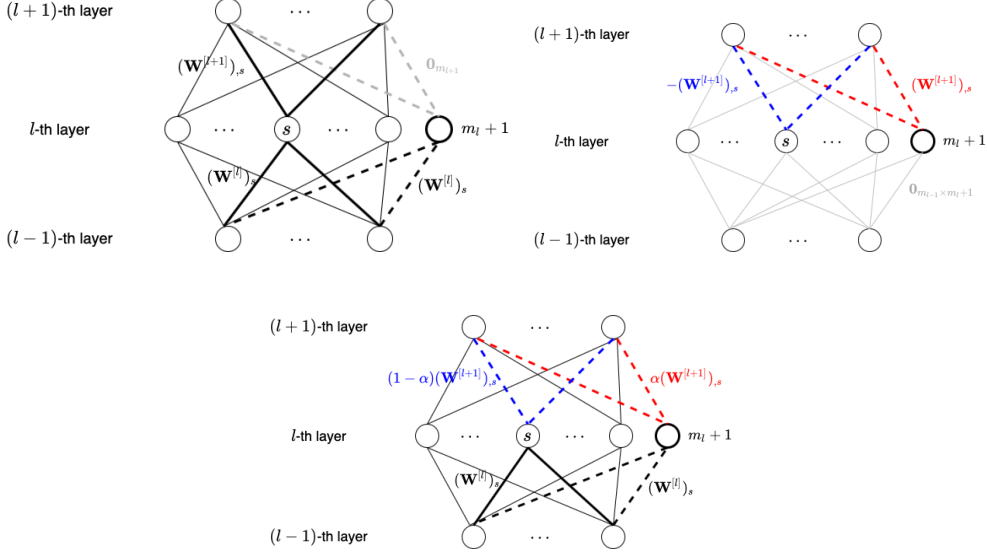


Figure S1: Illustration of $\mathcal{T}_{l,s}$, $\mathcal{V}_{l,s}$, and $\mathcal{T}_{l,s}^\alpha$.

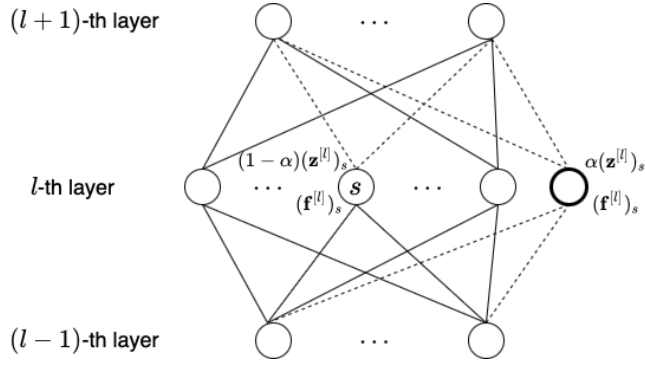


Figure S2: Illustration of \mathbf{F} and \mathbf{Z}

may result in multiple "mirror" critical points/manifolds of the loss landscape through all possible permutations. However, this transform does not inform about the degeneracy of critical points/manifolds.

For any p -homogeneous activation function σ i.e., $\sigma(\beta x) = \beta^p \sigma(x)$ for any $\beta > 0$ and $x \in \mathbb{R}$, we define for any $l \in [L-1]$, $s \in [m_l]$ the following scaling transform $\theta' = \mathcal{S}_{l,s}^\beta(\theta)$ ($\beta \neq 0$) such that $\mathbf{W}_{[1:m_{l+1}],s}^{\prime[l+1]} = \frac{1}{\beta^p} \mathbf{W}_{[1:m_{l+1}],s}^{[l+1]}$ and $\mathbf{W}_{s,[1:m_{l-1}]}^{\prime[l]} = \beta \mathbf{W}_{s,[1:m_{l-1}]}^{[l]}$, $\mathbf{b}_s^{\prime[l]} = \beta \mathbf{b}_s^{[l]}$, and all the other entries remain the same. Clearly, this transform is also a critical transform. Moreover, it informs about one more degenerate dimension for each neuron with $\left\| \mathbf{W}_{[1:m_{l+1}],s}^{\prime[l+1]} \right\|_2 \left\| \left(\mathbf{W}_{s,[1:m_{l-1}]}^{[l]\top}, \mathbf{b}_s^{[l]\top} \right)^\top \right\|_2 \neq 0$. This critical scaling transform is trivial in a sense that it is an obvious result of the cross-layer scaling preserving intrinsic to each DNN of homogeneous activation function, not relevant to cross-width landscape similarity between DNNs we focus on.

B Details of experiments

For the 1D fitting experiments (Figs. 1, 3(a), 4), we use tanh as the activation function, mean squared error (MSE) as the loss function. We use the full-batch gradient descent with learning rate 0.005 to train NNs for 300000 epochs. The initial distribution of all parameters follows a normal distribution with a mean of 0 and a variance of $\frac{1}{m^3}$.

For the iris classification experiment (Fig. 3(b)), we use sigmoid as the activation function, MSE as the loss function. We use the default Adam optimizer of full batch with learning rate 0.02 to train for 500000 epochs. The initial distribution of all parameters follows a normal distribution with mean 0 and variance $\frac{1}{m^6}$.

For the experiment of MNIST classification (Fig. 5), we use ReLU as the activation function, MSE as the loss function. We also use the default Adam optimizer of full batch with learning rate 0.00003 to train for 100000 epochs. The initial distribution of all parameters follows a normal distribution with mean 0 and variance $\frac{1}{m^6}$.

To obtain the empirical diagram in Fig. 4, we run 200 trials each for width-1, width-2 and width-3 tanh NNs with variance of initial parameters $\frac{1}{m^3}$ ($m = 1, 2, 3$) for 300000 epochs. Then we find all parameters with gradient less than 10^{-10} , which we define as empirical critical points, throughout the training in total 600 trajectories. Next, we cluster them based on their loss values, output functions, input parameters of neurons and only 4 different cases arises after excluding the trivial case of constant zero output. Their output functions are shown in the figure.

Remark that, although Figs. 1 and 5 are case studies each based on a random trial, similar phenomenon can be easily observed as long as the initialization variance is properly small, i.e., far from the linear/kernel/NTK regime.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
 - (b) Did you describe the limitations of your work? [Yes] See Section 6.
 - (c) Did you discuss any potential negative societal impacts of your work? [N/A]
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [Yes] See Section 5.
 - (b) Did you include complete proofs of all theoretical results? [Yes] See Appendix.
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] In the material.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] See Section 3.2
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [N/A]
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] the provider information Will be shown in Acknowledgement.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [Yes]
 - (b) Did you mention the license of the assets? [No] The MNIST dataset is well known.
 - (c) Did you include any new assets either in the supplemental material or as a URL? [No]
 - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]

- (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
- (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]