

A Appendix

A.1 Basic definitions

In this study, we first consider the neural network with 2 hidden layers,

A two-layer NN is

$$f_{\theta}(\mathbf{x}) = \sum_{j=1}^m a_j \sigma(\mathbf{w}_j \cdot \mathbf{x}), \quad (11)$$

where $\sigma(\cdot)$ is the activation function, $\mathbf{w}_j = (\bar{\mathbf{w}}_j, \mathbf{b}_j) \in \mathbb{R}^{d+1}$ is the neuron feature including the input weight and bias terms, and $\mathbf{x} = (\bar{\mathbf{x}}, 1) \in \mathbb{R}^{d+1}$ is concatenation of the input sample and scalar 1, θ is the set of all parameters, i.e., $\{a_j, \mathbf{w}_j\}_{j=1}^m$. For simplicity, we call \mathbf{w}_j **input weight or weight** and \mathbf{x} input sample.

Then, we consider the neural network with l hidden layers,

$$\begin{aligned} \mathbf{x}^{[0]} &= (\mathbf{x}, 1), \quad \mathbf{x}^{[1]} = (\sigma(\mathbf{W}^{[1]} \mathbf{x}^{[0]}), 1), \quad \mathbf{x}^{[l]} = (\sigma(\mathbf{W}^{[l]} \mathbf{x}^{[l-1]}), 1), \text{ for } l \in \{2, 3, \dots, L\} \\ f(\theta, \mathbf{x}) &= \mathbf{a}^\top \mathbf{x}^{[L]} \triangleq f_{\theta}(\mathbf{x}), \end{aligned} \quad (12)$$

where $\mathbf{W}^{[l]} = (\bar{\mathbf{W}}^{[l]}, \mathbf{b}^{[l]}) \in \mathbb{R}^{(m_l \times m_{l-1})}$, and m_l represents the dimension of the l -th hidden layer. The initialization of $\mathbf{W}_{k,k'}^{[l]}, l \in \{1, 2, 3, \dots, L\}$ and \mathbf{a}_k obey normal distribution $\mathbf{W}_{k,k'}^{[l]} \sim \mathcal{N}(0, \beta_l^2)$ for $l \in \{1, 2, 3, \dots, L\}$ and $\mathbf{a}_k \sim \mathcal{N}(0, \beta_{L+1}^2)$.

The loss function is mean squared error given below,

$$R(\theta) = \frac{1}{2n} \sum_{i=1}^n (f_{\theta}(\mathbf{x}_i) - y_i)^2. \quad (13)$$

For simplification, we denote $f_{\theta}(\mathbf{x})$ as f in following.

A.2 Derivations for concerned quantities

A.2.1 Neural networks with three hidden layers

In order to better understand the gradient of the parameter matrix of the multi-layer neural network, we first consider the case of the three-layer neural network,

$$f_{\theta}(\mathbf{x}) := \mathbf{a}^\top \sigma(\mathbf{W}^{[2]} \sigma(\mathbf{W}^{[1]} \mathbf{x})), \quad (14)$$

with the mean squared error as the loss function,

$$R_s(\theta) = \frac{1}{2n} \sum_{i=1}^n (f_{\theta}(\mathbf{x}_i) - y_i)^2. \quad (15)$$

We calculate $\frac{df}{d\mathbf{W}^{[2]}}$ and $\frac{df}{d\mathbf{W}^{[1]}}$ respectively, using differential form,

$$df = \text{tr}\left(\left(\frac{\partial f}{\partial \mathbf{x}}\right)^\top d\mathbf{x}\right). \quad (16)$$

We consider $\frac{df}{d\mathbf{W}^{[2]}}$ first,

$$\begin{aligned} df &= \text{tr}\{d(\mathbf{a}^\top \sigma(\mathbf{W}^{[2]} \mathbf{x}^{[1]}))\} \\ &= \text{tr}\{\mathbf{a}^\top d(\sigma(\mathbf{W}^{[2]} \mathbf{x}^{[1]}))\} \\ &= \text{tr}\{\mathbf{a}^\top \sigma'(\mathbf{W}^{[2]} \mathbf{x}^{[1]}) \odot d\mathbf{W}^{[2]} \mathbf{x}^{[1]}\} \\ &= \text{tr}\{(\mathbf{a} \odot \sigma'(\mathbf{W}^{[2]} \mathbf{x}^{[1]}))^\top d\mathbf{W}^{[2]} \mathbf{x}^{[1]}\} \\ &= \text{tr}\{\mathbf{x}^{[1]} (\mathbf{a} \odot \sigma'(\mathbf{W}^{[2]} \mathbf{x}^{[1]}))^\top d\mathbf{W}^{[2]}\} \\ &= \text{tr}\{((\mathbf{a} \odot \sigma'(\mathbf{W}^{[2]} \mathbf{x}^{[1]})) \mathbf{x}^{[1]\top})^\top d\mathbf{W}^{[2]}\}, \end{aligned} \quad (17)$$

where \odot is Hadamard product, and it is the multiplication of matrix elements of the same position. Hence,

$$\begin{aligned}\frac{df}{d\mathbf{W}^{[2]}} &= (\mathbf{a} \odot \sigma'(\mathbf{W}^{[2]}\mathbf{x}^{[1]}))\mathbf{x}^{[1]\top} \\ &= \text{diag}\{\sigma'(\mathbf{W}^{[2]}\mathbf{x}^{[1]})\}\mathbf{a}\mathbf{x}^{[1]\top}.\end{aligned}\quad (18)$$

Then, we consider $\frac{df}{d\mathbf{W}^{[1]}}$,

$$\begin{aligned}df &= \text{tr}\{(\mathbf{a} \odot \sigma'(\mathbf{W}^{[2]}\mathbf{x}^{[1]}))^\top \mathbf{W}^{[2]} d\sigma(\mathbf{W}^{[1]}\mathbf{x})\} \\ &= \text{tr}\{(\mathbf{W}^{[2]\top}(\mathbf{a} \odot \sigma'(\mathbf{W}^{[2]}\mathbf{x}^{[1]})))^\top \sigma'(\mathbf{W}^{[1]}\mathbf{x}) \odot d(\mathbf{W}^{[1]}\mathbf{x})\} \\ &= \text{tr}\{((\mathbf{W}^{[2]\top}(\mathbf{a} \odot \sigma'(\mathbf{W}^{[2]}\mathbf{x}^{[1]})) \odot \sigma'(\mathbf{W}^{[1]}\mathbf{x}))^\top d(\mathbf{W}^{[1]}\mathbf{x})\} \\ &= \text{tr}\{(((\mathbf{W}^{[2]\top}(\mathbf{a} \odot \sigma'(\mathbf{W}^{[2]}\mathbf{x}^{[1]})) \odot \sigma'(\mathbf{W}^{[1]}\mathbf{x}))\mathbf{x}^\top)^\top d(\mathbf{W}^{[1]}\mathbf{x})\}.\end{aligned}\quad (19)$$

Hence, we have,

$$\begin{aligned}\frac{df}{d\mathbf{W}^{[1]}} &= ((\mathbf{W}^{[2]\top}(\mathbf{a} \odot \sigma'(\mathbf{W}^{[2]}\mathbf{x}^{[1]}))) \odot \sigma'(\mathbf{W}^{[1]}\mathbf{x}))\mathbf{x}^\top \\ &= \text{diag}\{\sigma'(\mathbf{W}^{[1]}\mathbf{x})\}\mathbf{W}^{[2]\top} \text{diag}\{\sigma'(\mathbf{W}^{[2]}\mathbf{x}^{[1]})\}\mathbf{a}\mathbf{x}^\top.\end{aligned}\quad (20)$$

Through the chain rule, we can get the evolution equation of $\mathbf{W}^{[1]}$ and $\mathbf{W}^{[2]}$,

$$\begin{aligned}\frac{d\mathbf{W}^{[1]}}{dt} &= -\frac{dR_s(\boldsymbol{\theta})}{d\mathbf{W}^{[1]}} \\ &= -\frac{1}{n} \sum_{i=1}^n (f(\boldsymbol{\theta}, \mathbf{x}_i) - y_i) \frac{df}{d\mathbf{W}^{[1]}} \\ &= -\frac{1}{n} \sum_{i=1}^n (f(\boldsymbol{\theta}, \mathbf{x}_i) - y_i) \text{diag}\{\sigma'(\mathbf{W}^{[1]}\mathbf{x}_i)\} \mathbf{W}^{[2]\top} \text{diag}\{\sigma'(\mathbf{W}^{[2]}\mathbf{x}_i^{[1]})\} \mathbf{a}\mathbf{x}_i^\top,\end{aligned}\quad (21)$$

and

$$\begin{aligned}\frac{d\mathbf{W}^{[2]}}{dt} &= -\frac{dR_s(\boldsymbol{\theta})}{d\mathbf{W}^{[2]}} \\ &= -\frac{1}{n} \sum_{i=1}^n (f(\boldsymbol{\theta}, \mathbf{x}_i) - y_i) \frac{df}{d\mathbf{W}^{[1]}} \\ &= -\frac{1}{n} \sum_{i=1}^n (f(\boldsymbol{\theta}, \mathbf{x}_i) - y_i) \text{diag}\{\sigma'(\mathbf{W}^{[2]}\mathbf{x}_i^{[1]})\} \mathbf{a}\mathbf{x}_i^{[1]\top}.\end{aligned}\quad (22)$$

A.2.2 L hidden layers condition

And, we consider the neural network with L hidden layers,

$$\begin{aligned}df &= \text{tr}\{d\mathbf{a}^\top d\sigma(\mathbf{W}^{[L]}\mathbf{x}^{[L-1]})\} \\ &= \text{tr}\{(\mathbf{a} \odot \sigma'(\mathbf{W}^{[L]}\mathbf{x}^{[L-1]}))^\top d\mathbf{W}^{[L]} \sigma(\mathbf{W}^{[L-1]}\mathbf{x}^{[L-2]})\} \\ &= \text{tr}\{(\mathbf{W}^{[L]\top} \Lambda_L)^\top \sigma'(\mathbf{W}^{[L-1]}\mathbf{x}^{[L-2]}) \odot d\mathbf{W}^{[L-1]} \sigma(\mathbf{W}^{[L-2]}\mathbf{x}^{[L-3]})\} \\ &= \text{tr}\{((\mathbf{W}^{[L]\top} \Lambda_L) \odot \sigma'(\mathbf{W}^{[L-1]}\mathbf{x}^{[L-2]}))^\top \mathbf{W}^{[L-1]} d\sigma(\mathbf{W}^{[L-2]}\mathbf{x}^{[L-3]})\} \\ &= (\mathbf{W}^{[L-1]\top} \Lambda_{L-1})^\top d\sigma(\mathbf{W}^{[L-2]}\mathbf{x}^{[L-3]}) \\ &= \dots \\ &= \text{tr}\{\Lambda_k^\top d\mathbf{W}^{[k]}\mathbf{x}^{[k-1]}\} \\ &= \text{tr}\{(\Lambda_k \mathbf{x}^{[k-1]\top})^\top d\mathbf{W}^{[k]}\},\end{aligned}\quad (23)$$

where $\Lambda_l := (\mathbf{W}^{[l+1]^\top} \Lambda_{l+1}) \odot \sigma'(\mathbf{W}^{[l]} \mathbf{x}^{[l-1]})$ for $l = k, k+1 \dots L-1$ and $\Lambda_L := \mathbf{a} \odot \sigma'(\mathbf{W}^{[L]} \mathbf{x}^{[L-1]})$.

Hence, we get,

$$\frac{df}{d\mathbf{W}^{[k]}} = \Lambda_k \mathbf{x}^{[k-1]^\top}. \quad (24)$$

Through the chain rule, we can get the evolution equation of $\mathbf{W}^{[k]}$,

$$\begin{aligned} \frac{d\mathbf{W}^{[k]}}{dt} &= -\frac{dR_s(\boldsymbol{\theta})}{d\mathbf{W}^{[k]}} \\ &= -\frac{1}{n} \sum_{i=1}^n (f(\boldsymbol{\theta}, \mathbf{x}_i) - y_i) \frac{df}{d\mathbf{W}^{[k]}} \\ &= -\frac{1}{n} \sum_{i=1}^n (f(\boldsymbol{\theta}, \mathbf{x}_i) - y_i) \Lambda_k \mathbf{x}_i^{[k-1]^\top}. \end{aligned} \quad (25)$$

Through $\mathbf{a} \odot \sigma'(\mathbf{W} \mathbf{x}) = \text{diag}\{\sigma'(\mathbf{W} \mathbf{x})\} \mathbf{a}$,

Finally, the dynamic system can be obtained:

$$\begin{aligned} \dot{\mathbf{a}} &= \frac{d\mathbf{a}}{dt} = -\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^{[L]} (f(\boldsymbol{\theta}, \mathbf{x}_i) - y_i), \\ \dot{\mathbf{W}}^{[L]} &= \frac{d\mathbf{W}^{[L]}}{dt} = -\frac{1}{n} \sum_{i=1}^n \text{diag}\{\sigma'(\mathbf{W}^{[L]} \mathbf{x}_i^{[L-1]})\} \mathbf{a} \mathbf{x}_i^{[L-1]^\top} (f(\boldsymbol{\theta}, \mathbf{x}_i) - y_i), \\ \dot{\mathbf{W}}^{[k]} &= \frac{d\mathbf{W}^{[k]}}{dt} = -\frac{1}{n} \sum_{i=1}^n \text{diag}\{\sigma'(\mathbf{W}^{[k]} \mathbf{x}_i^{[k-1]})\} \mathbf{E}^{[k+1:L]} \mathbf{a} \mathbf{x}_i^{[k-1]^\top} (f(\boldsymbol{\theta}, \mathbf{x}_i) - y_i) \quad \forall i \in [1 : L-1], \end{aligned} \quad (26)$$

where we use $\mathbf{E}^l(\mathbf{x}) = \mathbf{W}^{[l]^\top} \text{diag}\{\sigma'(\mathbf{W}^{[l]} \mathbf{x}^{[l-1]})\}$. And $\mathbf{E}^{[q:p]} = \mathbf{E}^q \mathbf{E}^{q+1} \dots \mathbf{E}^p$.

Let $r_{k,j} = \|\mathbf{W}_j^{[k]}\|_2$. We have

$$\frac{d}{dt} |r_{k,j}|^2 = \frac{d}{dt} \|\mathbf{W}_j^{[k]}\|^2. \quad (27)$$

Then we obtain

$$\dot{r}_{k,j} r_{k,j} = \dot{\mathbf{W}}_j^{[k]} \cdot \mathbf{W}_j^{[k]}. \quad (28)$$

Finally, we get

$$\begin{aligned} \dot{r}_{k,j} &= \frac{dr_{k,j}}{dt} = \dot{\mathbf{W}}_j^{[k]} \cdot \mathbf{W}_j^{[k]} / r_{k,j} \\ &= \dot{\mathbf{W}}_j^{[k]} \cdot \mathbf{u}_{k,j}, \end{aligned} \quad (29)$$

where $\mathbf{u}_{k,j} = \frac{\mathbf{W}_j^{[k]}}{r_{k,j}}$ is a unit vector. Then we have,

$$\begin{aligned} \dot{\mathbf{u}}_{k,j} &= \frac{d\mathbf{u}_{k,j}}{dt} = \frac{d}{dt} \left(\frac{\mathbf{W}_j^{[k]}}{r_{k,j}} \right) \\ &= \frac{\dot{\mathbf{W}}_j^{[k]} r_{k,j} - \mathbf{W}_j^{[k]} \dot{r}_{k,j}}{r_{k,j}^2} \\ &= \frac{\dot{\mathbf{W}}_j^{[k]} r_{k,j} - \mathbf{W}_j^{[k]} (\dot{\mathbf{W}}_j^{[k]} \cdot \mathbf{u}_{k,j})}{r_{k,j}^2} \\ &= \frac{\dot{\mathbf{W}}_j^{[k]} - \mathbf{u}_{k,j} (\dot{\mathbf{W}}_j^{[k]} \cdot \mathbf{u}_{k,j})}{r_{k,j}}. \end{aligned} \quad (30)$$

To conclude, the quantities we concern are summarized as follows,

$$\begin{cases} \dot{\mathbf{a}} = -\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^{[L]} (f(\boldsymbol{\theta}, \mathbf{x}_i) - y_i) \end{cases} \quad (31)$$

$$\begin{cases} \dot{\mathbf{W}}^{[L]} = -\frac{1}{n} \sum_{i=1}^n \text{diag}\{\sigma'(\mathbf{W}^{[L]} \mathbf{x}_i^{[L-1]})\} \mathbf{a} \mathbf{x}_i^{[L-1]\top} (f(\boldsymbol{\theta}, \mathbf{x}_i) - y_i), \end{cases} \quad (32)$$

$$\begin{cases} \dot{\mathbf{W}}^{[k]} = -\frac{1}{n} \sum_{i=1}^n \text{diag}\{\sigma'(\mathbf{W}^{[k]} \mathbf{x}_i^{[k-1]})\} \mathbf{E}^{[k+1:L]} \mathbf{a} \mathbf{x}_i^{[k-1]\top} (f(\boldsymbol{\theta}, \mathbf{x}_i) - y_i) \quad \forall k \in [1 : L-1] \end{cases} \quad (33)$$

$$\begin{cases} \dot{r}_{k,j} = \dot{\mathbf{W}}_j^{[k]} \cdot \mathbf{u}_{k,j} \end{cases} \quad (34)$$

$$\begin{cases} \dot{\mathbf{u}}_{k,j} = \frac{\dot{\mathbf{W}}_j^{[k]} - \mathbf{u}_{k,j} (\dot{\mathbf{W}}_j^{[k]} \cdot \mathbf{u}_{k,j})}{r_{k,j}}, \end{cases} \quad (35)$$

where we use $\mathbf{E}^l(\mathbf{x}) = \mathbf{W}^{[l]\top} \text{diag}\{\sigma'(\mathbf{W}^{[l]} \mathbf{x}^{[l-1]})\}$. And $\mathbf{E}^{[q:p]} = \mathbf{E}^q \mathbf{E}^{q+1} \dots \mathbf{E}^p$.

A.2.3 Prove for $\mathcal{P}\mathbf{w}$ in 5

We calculate $\mathcal{P}\mathbf{w} \stackrel{\text{leading order}}{\approx} \mathcal{Q}\mathbf{w}$ as following,

$$\begin{aligned} \mathcal{P}\mathbf{w} \approx \mathcal{Q}\mathbf{w} : &= -\frac{1}{n} \sum_{i=1}^n e_i \mathbf{x}_i^{[k-1]} [\text{diag}\{\sigma'(\mathbf{W}^{[k]} \mathbf{x}_i^{[k-1]})\} (\mathbf{E}^{[k+1:L]} \mathbf{a})]_j \\ &+ \left(\frac{1}{n} \sum_{i=1}^n e_i \mathbf{x}_i^{[k-1]} [\text{diag}\{\sigma'(\mathbf{W}^{[k]} \mathbf{x}_i^{[k-1]})\} (\mathbf{E}^{[k+1:L]} \mathbf{a})]_j \cdot \mathbf{u} \right) \mathbf{u} \\ &= -\frac{1}{n} \sum_{i=1}^n e_i \mathbf{x}_i^{[k-1]} [\text{diag}\left\{ \frac{\sigma^{(p)}(\mathbf{0})}{(p-1)!} \odot (\mathbf{W}^{[k]} \mathbf{x}_i^{[k-1]})^{p-1} \right\} (\mathbf{E}^{[k+1:L]} \mathbf{a})]_j \\ &+ \left(\frac{1}{n} \sum_{i=1}^n e_i \mathbf{x}_i^{[k-1]} [\text{diag}\left\{ \frac{\sigma^{(p)}(\mathbf{0})}{(p-1)!} \odot (\mathbf{W}^{[k]} \mathbf{x}_i^{[k-1]})^{p-1} \right\} (\mathbf{E}^{[k+1:L]} \mathbf{a})]_j \cdot \mathbf{u} \right) \mathbf{u} \\ &= -\frac{1}{n} \sum_{i=1}^n e_i \mathbf{x}_i^{[k-1]} [\text{diag}\{(\mathbf{W}^{[k]} \mathbf{x}_i^{[k-1]})^{p-1}\} \text{diag}\left\{ \frac{\sigma^{(p)}(\mathbf{0})}{(p-1)!} \right\} (\mathbf{E}^{[k+1:L]} \mathbf{a})]_j \\ &+ \left(\frac{1}{n} \sum_{i=1}^n e_i \mathbf{x}_i^{[k-1]} [\text{diag}\{(\mathbf{W}^{[k]} \mathbf{x}_i^{[k-1]})^{p-1}\} \text{diag}\left\{ \frac{\sigma^{(p)}(\mathbf{0})}{(p-1)!} \right\} (\mathbf{E}^{[k+1:L]} \mathbf{a})]_j \cdot \mathbf{u} \right) \mathbf{u} \\ &= -\frac{1}{n} \sum_{i=1}^n e_i \mathbf{x}_i^{[k-1]} [\text{diag}\{(\mathbf{W}^{[k]} \mathbf{x}_i^{[k-1]})^{p-1}\}]_j \text{diag}\left\{ \frac{\sigma^{(p)}(\mathbf{0})}{(p-1)!} \right\} (\mathbf{E}^{[k+1:L]} \mathbf{a}) \\ &+ \left(\frac{1}{n} \sum_{i=1}^n e_i \mathbf{x}_i^{[k-1]} [\text{diag}\{(\mathbf{W}^{[k]} \mathbf{x}_i^{[k-1]})^{p-1}\}]_j \text{diag}\left\{ \frac{\sigma^{(p)}(\mathbf{0})}{(p-1)!} \right\} (\mathbf{E}^{[k+1:L]} \mathbf{a}) \cdot \mathbf{u} \right) \mathbf{u} \\ &= -\left(\frac{1}{n} \sum_{i=1}^n e_i \mathbf{x}_i^{[k-1]} (\mathbf{W}_j^{[k]} \mathbf{x}_i^{[k-1]})^{p-1} \right) [\text{diag}\left\{ \frac{\sigma^{(p)}(\mathbf{0})}{(p-1)!} \right\} (\mathbf{E}^{[k+1:L]} \mathbf{a})]_j \\ &+ \left(\left(\frac{1}{n} \sum_{i=1}^n e_i \mathbf{x}_i^{[k-1]} (\mathbf{W}_j^{[k]} \mathbf{x}_i^{[k-1]})^{p-1} \right) [\text{diag}\left\{ \frac{\sigma^{(p)}(\mathbf{0})}{(p-1)!} \right\} (\mathbf{E}^{[k+1:L]} \mathbf{a})]_j \cdot \mathbf{u} \right) \mathbf{u}, \end{aligned} \quad (36)$$

where $(\cdot)^{p-1}$ and $\sigma^p(\cdot)$ operate on component here.

A.3 The Verification of the initial stage

We put the loss of the experiments in the main text here to show that they are indeed in the initial stage of training by the definition.

As is shown in Fig.7 and Fig.8, at the steps demonstrated in the article, loss satisfies the definition of the initial stage, so we consider that they are in the initial stage of training.

Learning rate is not a sensitive to the appearance of condensation. However, a small learning rate could enable us to observe the condensation process in the initial stage under a gradient flow training, more clearly. For example, when the learning rate is relatively small, the initial stage of training may be relatively long, while when the learning rate is relatively large, the initial stage of training may be relatively small.

We empirically find that to ensure the training process follows a gradient follow, where the loss decays monotonically, we have to select a smaller learning rate for large multiplicity p . Therefore, it looks like we have a longer training in our experiments with large p . Note that for a small learning rate in the experiments of small p , we can observe similar phenomena.

In all subsequent experiments in the Appendix, we will no longer show the loss graph of each experiment one by one, but we make sure that they are indeed in the initial stage of training.

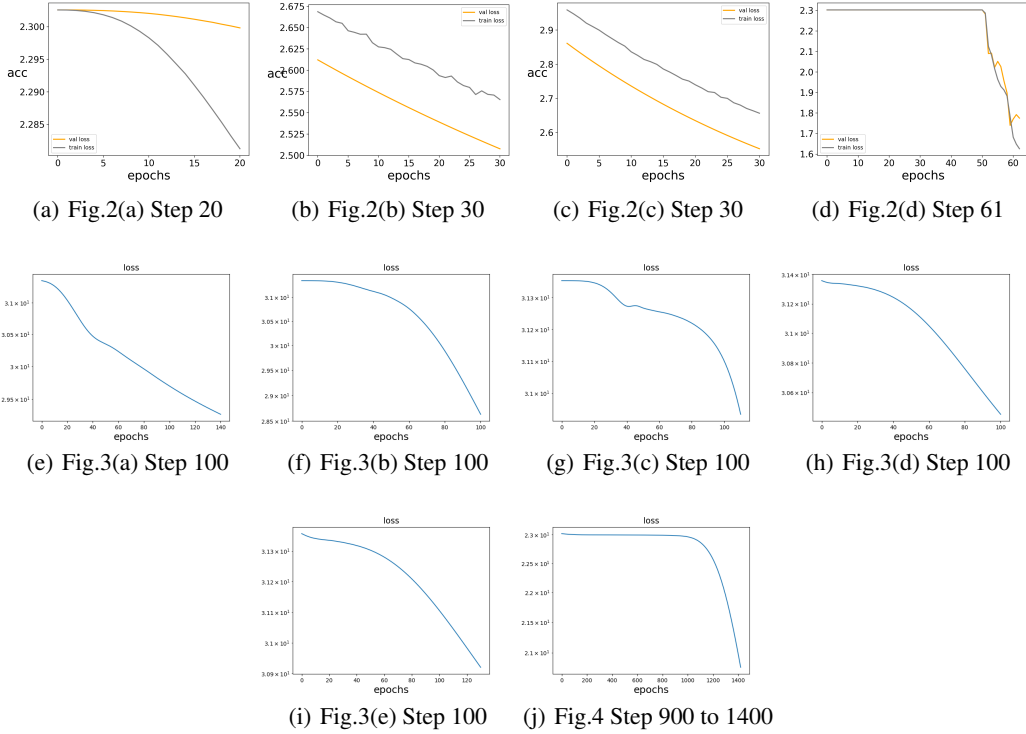


Figure 7: Losses from Fig. 2 to Fig.4. The original figures and the numbers of steps corresponding to each sub-picture are written in the sub-captions.

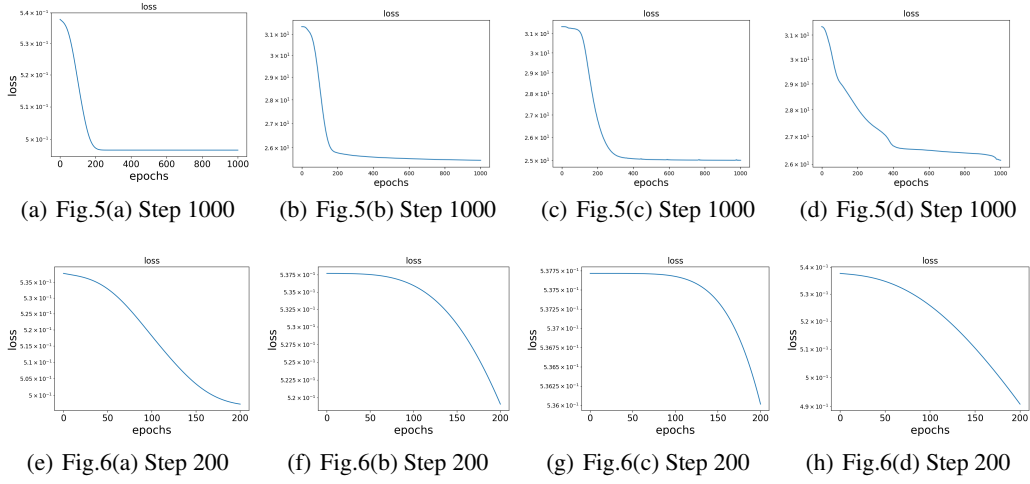
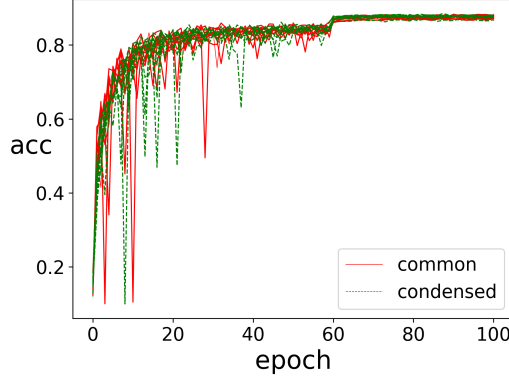


Figure 8: Losses from Fig. 5 to Fig.6. The original figures and the numbers of steps corresponding to each sub-picture are written in the sub-captions.

A.4 Performance of tanh activation function in condensed regime

For practical networks, such as resnet18-like (He et al., 2016) in learning CIFAR10, as shown in Fig. 9 and Table 1, we find that the performance of networks with initialization in the condensed regime is vary similar to the common initialization methods. For both initialization methods, the test accuracy is about 86.5% to 88.5%, where the highest test accuracy and lowest test accuracy for common methods are 88.07% and 86.73%, respectively, while the highest one and lowest one for condensed methods are 88.26% and 86.88%, respectively. This implies that performance of common and condensed initialization is similar.



(a) test accuracy

Figure 9: The test accuracy of Resnet18-like networks with different initialization methods. Each network consists of the convolution part of resnet18 and fully-connected (FC) layers with size 1024-1024-10 and softmax. The convolution part is equipped with ReLU activation and initialized by Glorot normal distribution (Glorot and Bengio, 2010). For FC layers, the activation is $\tanh(x)$ and they are initialized by three common methods (red) and three condensed ones (green) as indicated in Table 1. The learning rate is 10^{-3} for epoch 1-60 and 10^{-4} for epoch 61-100. Adam optimizer with cross-entropy loss and batch size 128 are used for all experiments.

Table 1: Comparison of test accuracy of resnet18 in learning CIFAR10 with common (Glorot and Bengio, 2010) and condensed Gaussian initializations. $\bar{m} = (m_{\text{in}} + m_{\text{out}})/2$. m_{in} : in-layer width. m_{out} : out-layer width. Each line is a trial.

	common			condensed		
	Glorot_uniform	Glorot_normal	$N(0, \frac{1}{\bar{m}})$	$N(0, \frac{1}{m_{\text{out}}^4})$	$N(0, \frac{1}{m_{\text{out}}^3})$	$N(0, (\frac{1}{\bar{m}})^2)$
Test 1	0.8747	0.8759	0.8807	0.8749	0.8744	0.8765
Test 2	0.8715	0.8673	0.8733	0.8763	0.8799	0.8826
Test 3	0.8772	0.8794	0.8788	0.8688	0.8780	0.8771

A.5 multi-layer experimental

The condensation of the six layer without residual connections is shown in 10, whose activation functions for hidden layer 1 to hidden layer 5 are $x^2 \tanh(x)$, $x \tanh(x)$, $\text{sigmoid}(x)$, $\tanh(x)$ and $\text{softplus}(x)$, respectively.

The condensation of the three layer without residual connections is shown in 11, whose activation functions are same for each layer indicated by the corresponding sub-captions.

The condensation of the five layer without residual connections is shown in 12, whose activation functions are same for each layer indicated by the corresponding sub-captions.

The condensation of the five layer with residual connections is shown in 13, whose activation functions are same for each layer indicated by the corresponding sub-captions.

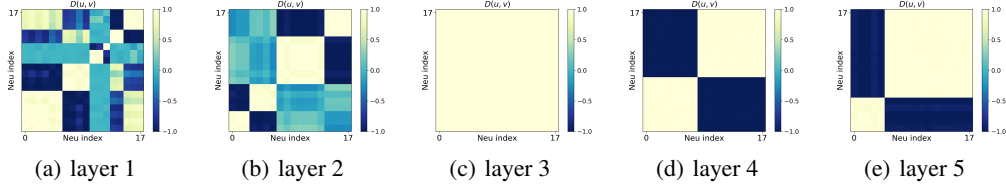


Figure 10: Condensation of six-layer NNs without residual connections. The activation functions for hidden layer 1 to hidden layer 5 are $x^2 \tanh(x)$, $x \tanh(x)$, $\text{sigmoid}(x)$, $\tanh(x)$ and $\text{softplus}(x)$, respectively. The numbers of steps selected in the sub-pictures are epoch 6800, epoch 6800, epoch 6800, epoch 6800 and epoch 6300, respectively, while the NN is only trained once. The color indicates $D(\mathbf{u}, \mathbf{v})$ of two hidden neurons' input weights, whose indexes are indicated by the abscissa and the ordinate, respectively. The training data is 80 points sampled from a 3-dimensional function $\sum_{k=1}^3 4 \sin(12x_k + 1)$, where each x_k is uniformly sampled from $[-4, 2]$. $n = 80$, $d = 3$, $m = 18$, $d_{\text{out}} = 1$, $\text{var} = 0.008^2$, $\text{lr} = 5 \times 10^{-5}$.

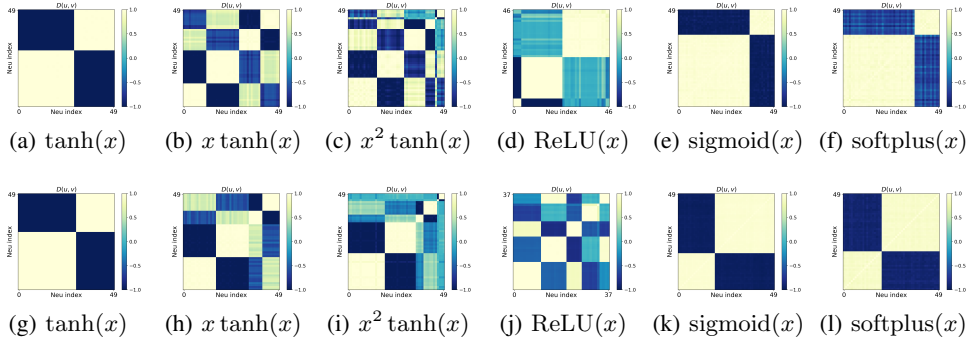


Figure 11: Three-layer NN at epoch 700. (a-f) are for the input weights of the first hidden layer and (g-l) are for the input weights of the second hidden layer. The color indicates $D(\mathbf{u}, \mathbf{v})$ of two hidden neurons' input weights, whose indexes are indicated by the abscissa and the ordinate, respectively. The training data is 80 points sampled from a 5-dimensional function $\sum_{k=1}^5 3 \sin(8x_k + 1)$, where each x_k is uniformly sampled from $[-4, 2]$. $n = 80$, $d = 5$, $m = 50$, $d_{\text{out}} = 1$, $\text{var} = 0.005^2$. $\text{lr} = 10^{-4}, 2 \times 10^{-5}, 1.4 \times 10^{-5}$ for (a-d), (e) and (f), respectively. For (d) and (j), we discard hidden neurons, whose L_2 -norm of its input weight is smaller than 0.1.

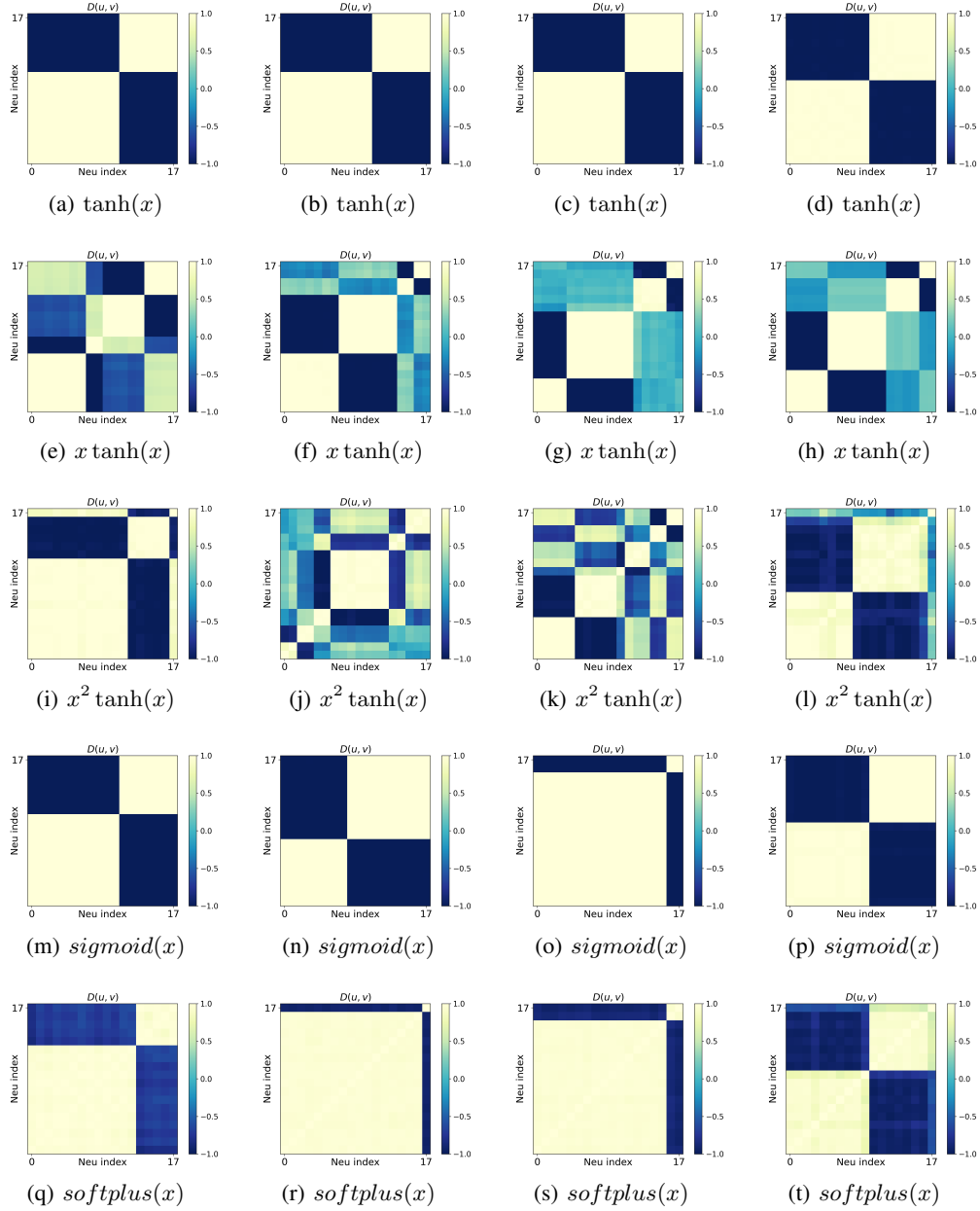


Figure 12: Five-layer NN. The first to fourth columns of each row are for the input weights of neurons from the first to the fourth hidden layers, respectively. The color indicates $D(u, v)$ of two hidden neurons' input weights, whose indexes are indicated by the abscissa and the ordinate, respectively. The training data is 80 points sampled from a 5-dimensional function $\sum_{k=1}^3 3 \sin(10x_k + 1)$, where each x_k is uniformly sampled from $[-4, 2]$. $n = 80$, $d = 5$, $m = 18$, $d_{\text{out}} = 1$, $\text{var} = 0.008^2$. $\text{lr} = 1.5 \times 10^{-5}$, 1.5×10^{-5} , 1.5×10^{-5} , 1.5×10^{-5} , 1.5×10^{-6} and epoch is 10000, 10000, 26000, 10000, 20000 for $\tanh(x)$, $x \tanh(x)$, $x^2 \tanh(x)$, $\text{sigmoid}(x)$, $\text{softplus}(x)$, respectively.

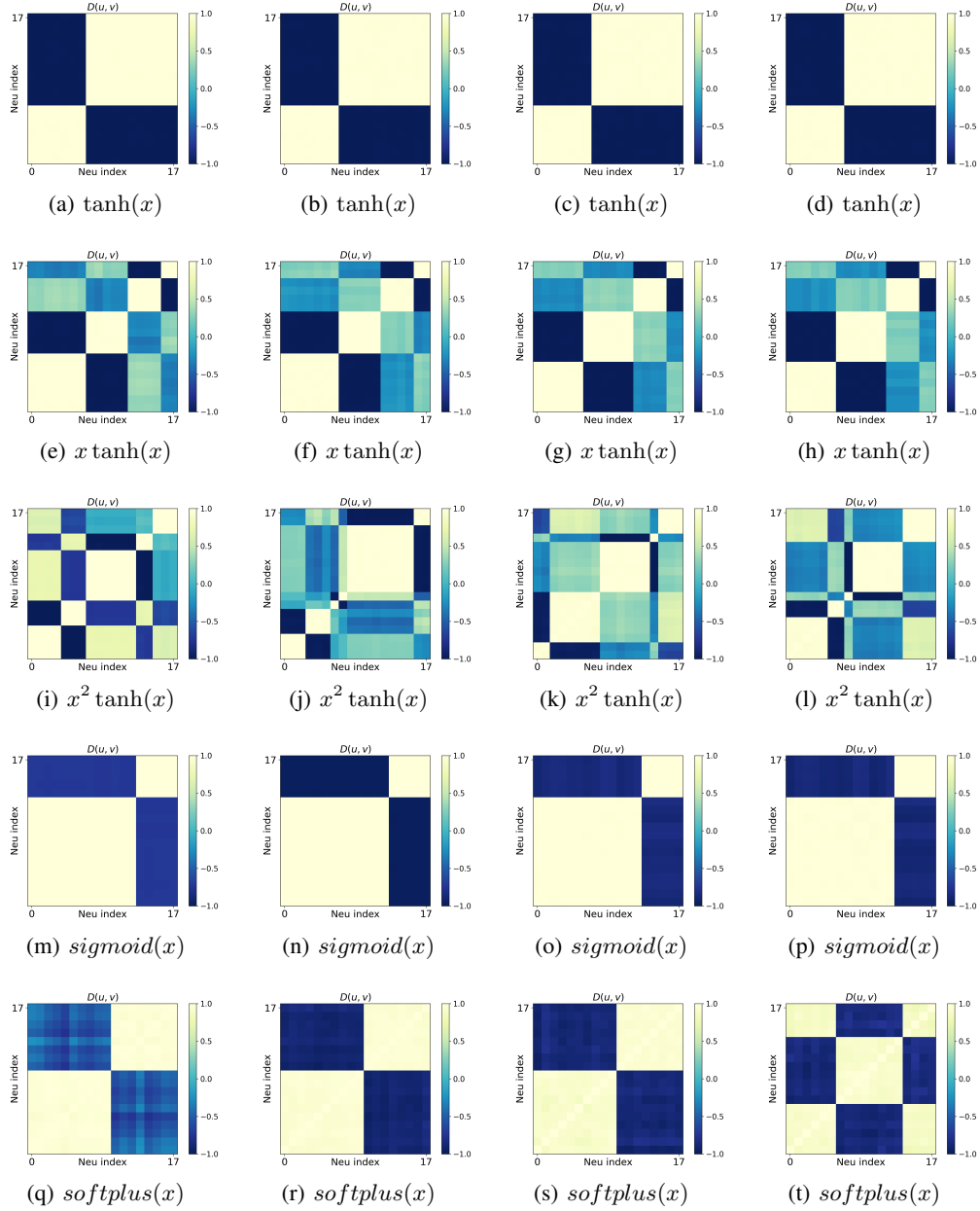


Figure 13: Five-layer NN. The first to fourth columns of each row are for the input weights of neurons from the first to the fourth hidden layers, respectively. The color indicates $D(u, v)$ of two hidden neurons' input weights, whose indexes are indicated by the abscissa and the ordinate, respectively. The training data is 80 points sampled from a 5-dimensional function $\sum_{k=1}^3 3 \sin(10x_k + 1)$, where each x_k is uniformly sampled from $[-4, 2]$. $n = 80$, $d = 5$, $m = 18$, $d_{\text{out}} = 1$, $\text{var} = 0.008^2$. $\text{lr} = 1 \times 10^{-4}$, 1×10^{-4} , 1×10^{-4} , 5×10^{-5} , 5×10^{-5} and epoch is 400, 400, 400, 3000, 360, 400 for $\tanh(x)$, $x \tanh(x)$, $x^2 \tanh(x)$, $x^2 \tanh(x)$, $\text{sigmoid}(x)$, $\text{softplus}(x)$, respectively.

A.6 Several steps during the evolution of condensation at the initial stage

In the article, we only give the results of the last step of each condense, while the details of the evolution of condensation are lacking, which may provide a better understanding. Therefore, we show these details in Fig. 14, Fig. 15, Fig. 16 and Fig. 17, which also further illustrate the rationality of the experimental results and facilitate the understanding of the evolution of condensation in the initial stage.

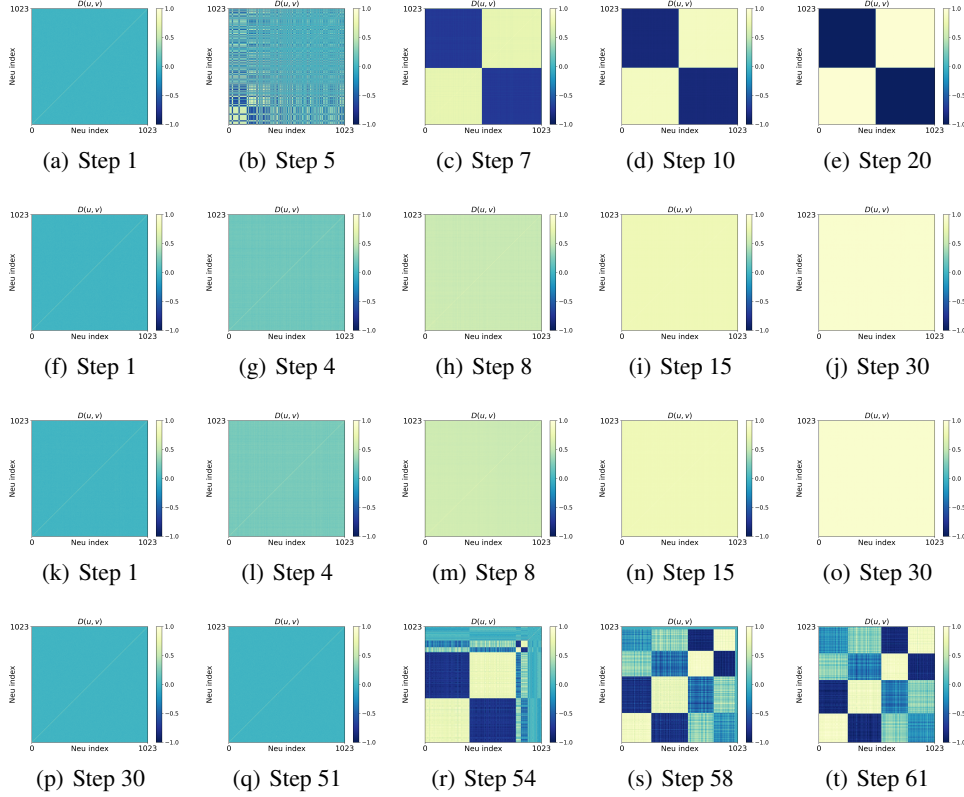


Figure 14: Evolution of condensation of Fig. 2(a), Fig. 2(b), Fig. 2(c), and Fig. 2(d). The evolution from the first row to the fourth row are corresponding to the Fig. 2(a), Fig. 2(b), Fig. 2(c), and Fig. 2(d). The numbers of evolutionary steps are shown in the sub-captions, where sub-figures in the last row are the epochs in the article.

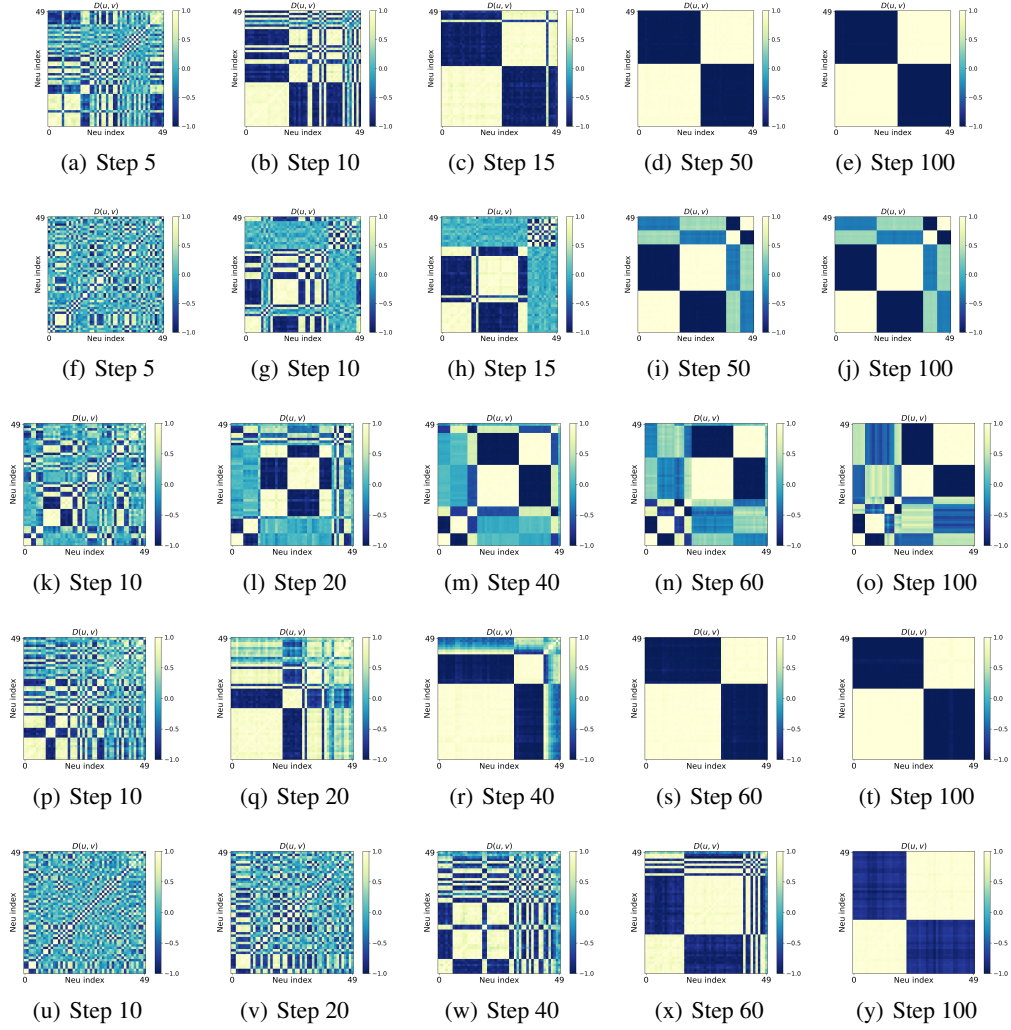


Figure 15: Evolution of condensation from Fig. 3(a) to 3(e). The evolution from the first row to the fifth row are corresponding to the Fig. 3(a), Fig. 3(b), Fig. 3(c), Fig. 3(d), Fig. 3(e). The numbers of evolutionary steps are shown in the sub-captions, where sub-figures in the last row are the epochs in the article.

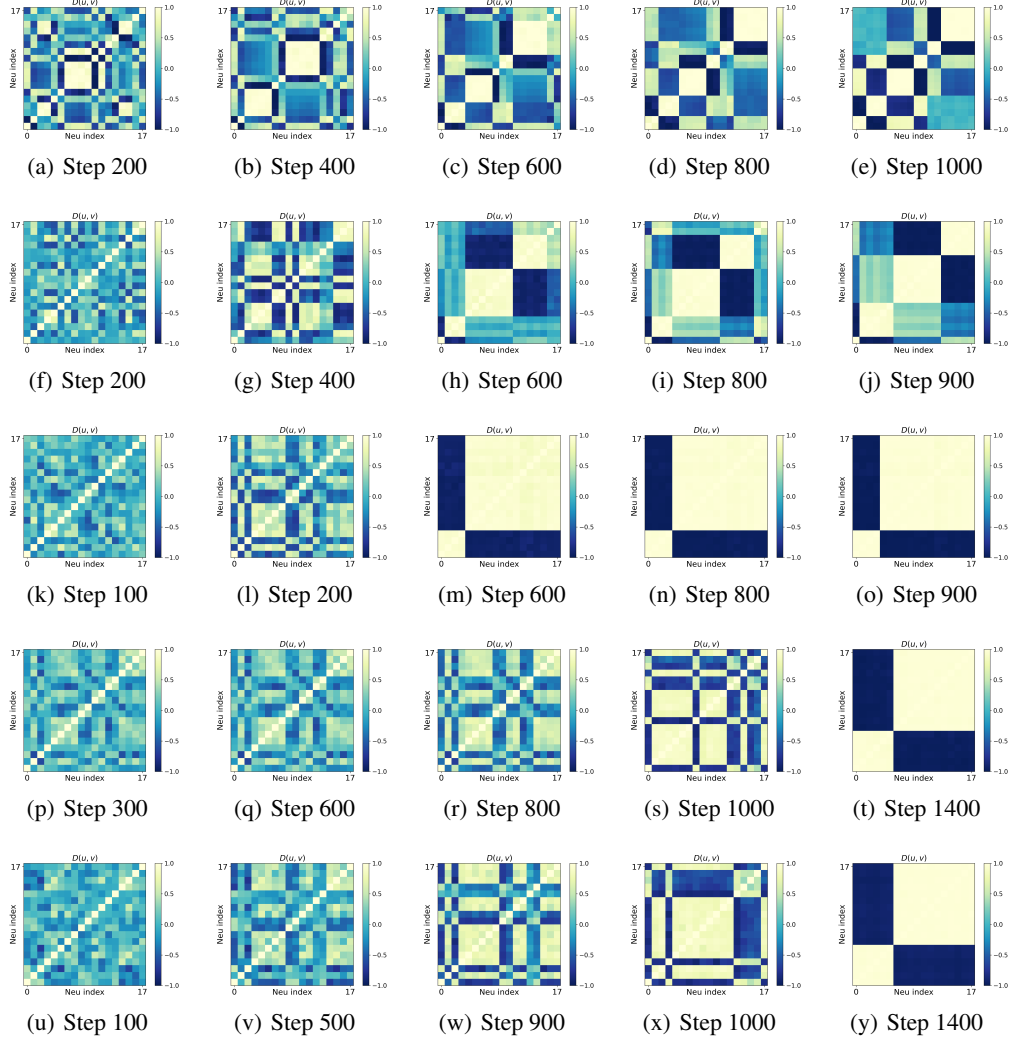


Figure 16: Evolution of condensation from Fig. 4(a) to 4(e). The evolution from the first row to the fifth row are corresponding to the Fig. 4(a), Fig. 4(b), Fig. 4(c), Fig. 4(d), Fig. 4(e). The numbers of evolutionary steps are shown in the sub-captions, where sub-figures in the last row are the epochs in the article.

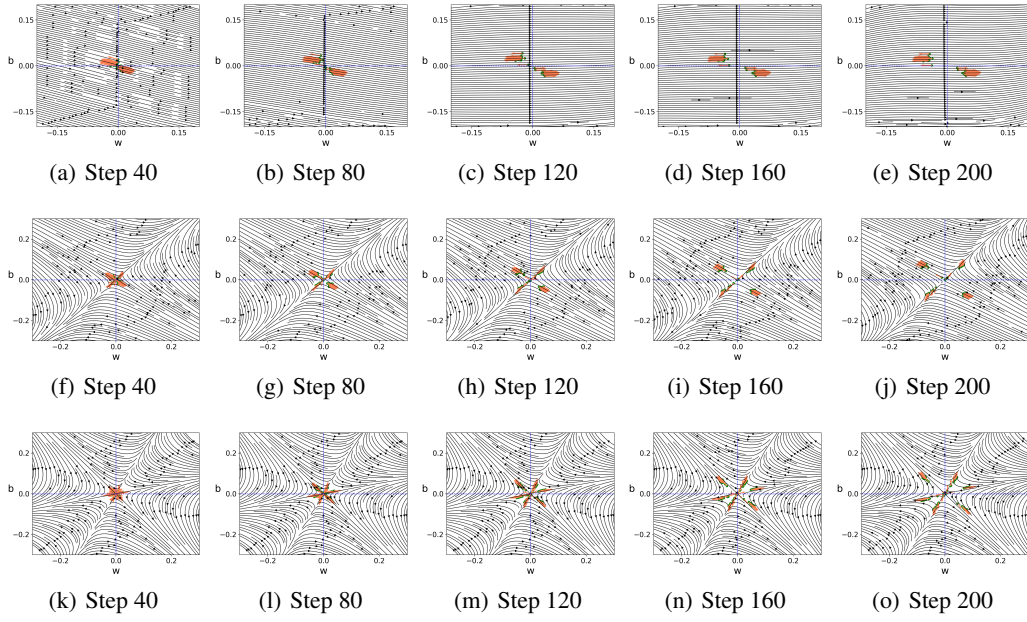


Figure 17: Evolution of condensation from Fig. 6(a) to 6(c). The evolution from the first row to the fifth row are corresponding to the Fig. 6(a), Fig. 6(b), and Fig. 6(c). The numbers of evolutionary steps are shown in the sub-captions, where sub-figures in the last row are the epochs in the article.

A.7 The influence of training data on condensation

We also find that when the training data is less oscillated, the NN may condense at fewer directions. For example, as shown in Fig. 18(a), compared with the high frequency function in Fig. 3, we only change the target function to be a lower-frequency function, i.e., $\sum_{k=1}^5 3.5 \sin(2x_k + 1)$. In this case, the NN with $x^2 \tanh(x)$ only condenses at three directions, in which two are opposite. For MNIST data in Fig. 18(b), we find that, the NN with $x^2 \tanh(x)$ condenses at one line, which may suggest that the function for fitting MNIST dataset is a low-frequency function. For CIFAR100 data in Fig. 18(c), we find that input weights of the first FC layer with $x \tanh(x)$ condense at only one line, which implies that features extracted by the convolution part of the NN may own low complexity.

These experiments does not contradict to our results, which claim that the *maximal* number of condensed orientations in the initial training is twice the multiplicity of the activation function used in general NNs.

For CIFAR100 dataset, we use Resnet18-like neural network, which has been described in Fig. 2. Besides, the input dimension is $d = 32 * 32 * 3$, the output dimension is $d_{\text{out}} = 100$, and all parameters are initialized by a Gaussian distribution $N(0, \text{var})$. The total data size is n . The training method is Adam with batch size 128, learning rate lr and cross-entropy loss.

For MNIST dataset, we use fully-connected neural network with size, $d-m \cdots m-d_{\text{out}}$. The input dimension is $d = 784$, and the output dimension is $d_{\text{out}} = 10$. The number of hidden neurons m is specified in Fig. 18. All parameters are initialized by a Gaussian distribution $N(0, \text{var})$. The total data size is n . The training method is Adam with full batch, learning rate lr and MSE loss.

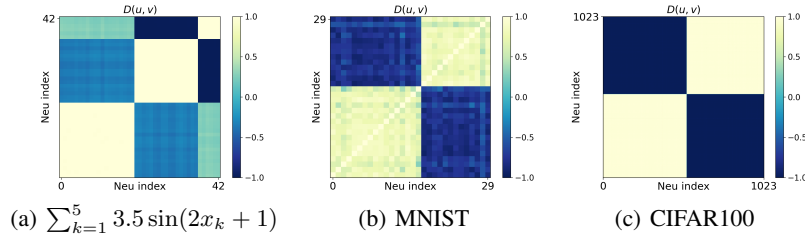


Figure 18: Condensation of low-frequency functions with two-layer NNs in (a,b) and condensation of the first FC layer of the Resnet18-like network on CIFAR100 in (c). The color indicates $D(u, v)$ of two hidden neurons' input weights, whose indexes are indicated by the abscissa and the ordinate. For (a,b), two-layer NN at epoch: 100 with activation function: $x^2 \tanh(x)$. For (a), we discard about 15% of hidden neurons, in which the L_2 -norm of each input weight is smaller than 0.04, while remaining those bigger than 0.4. The mean magnitude here for each parameter is $(0.4^2/785)^{0.5} \sim 0.01$, which should also be quite small. All settings in (a) are the same as Fig. 3, except for the lower frequency target function. Parameters for (b) are $n = 60000$, $d = 784$, $m = 30$, $d_{\text{out}} = 10$, $\text{var} = 0.001^2$, $\text{lr} = 5 \times 10^{-5}$. The structure and parameters of the Resnet18-like neural network for (c) is the same as Fig. 2, except for the data set CIFAR100 and learning rate $\text{lr} = 1 \times 10^{-6}$.