

Appendix of "The Orthogonality of Weight Vectors: The Key Characteristics of Normalization and Residual Connections"

1 Proof of Theorem 2

Proof. When using the SGD optimizer to update parameters,

$$\mathbf{w}_{i,t+1}^l = \mathbf{w}_{i,t}^l - \eta_t \nabla \mathbf{w}_{i,t}^l, \quad (1)$$

η_t is the learning rate. Let $\mathbf{w}_{i,t}^l$ and $\mathbf{w}_{j,t}^l$ be any two vectors of the weight W_t^l then

$$\begin{aligned} \mathbf{w}_{i,t+1}^l &= \mathbf{w}_{i,t}^l - \eta_t \nabla \mathbf{w}_{i,t}^l \\ \mathbf{w}_{j,t+1}^l &= \mathbf{w}_{j,t}^l - \eta_t \nabla \mathbf{w}_{j,t}^l, \end{aligned} \quad (2)$$

$\mathbf{w}_{i,t+1}^l$ and $\mathbf{w}_{j,t+1}^l$ are the translation transformations of $\mathbf{w}_{i,t}^l$ and $\mathbf{w}_{j,t}^l$ along the $\eta_t \nabla \mathbf{w}_{i,t}^l$ and $\eta_t \nabla \mathbf{w}_{j,t}^l$ directions, in particular, when $\nabla \mathbf{w}_{i,t}^l = \nabla \mathbf{w}_{j,t}^l$, $\mathbf{w}_{i,t}^l$ and $\mathbf{w}_{j,t}^l$ move in the same direction. Therefore, the angle between vectors remains unchanged under translation, and in this case, the translation transformation is a conformal mapping. In the same way, when updating parameters using the Adam optimizer,

$$\begin{aligned} \mathbf{w}_{i,t+1}^l &= \mathbf{w}_{i,t}^l - \eta_t \frac{\mathbf{m}_{t+1}}{\sqrt{\mathbf{v}_{t+1}} + \epsilon} b_{t+1} \\ \mathbf{m}_{t+1} &= \beta_1 \mathbf{m}_t + (1 - \beta_1) \nabla \mathbf{w}_t^l \\ \mathbf{v}_{t+1} &= \beta_2 \mathbf{v}_t + (1 - \beta_2) \nabla \mathbf{w}_t^l \odot \nabla \mathbf{w}_t^l \\ b_{t+1} &= \frac{\sqrt{1 - \beta_2^{t+1}}}{1 - \beta_1^{t+1}}. \end{aligned} \quad (3)$$

Where \odot is the element-wise multiplication. $\beta_1, \beta_2 \in [0, 1)$ are hyperparameters. β_1^{t+1} and β_2^{t+1} denote β_1 and β_2 to the power $t + 1$. For any two vectors $\mathbf{w}_{i,t}^l$ and $\mathbf{w}_{j,t}^l$, $\mathbf{m}_{i,t+1} = \mathbf{m}_{j,t+1}$ and $\mathbf{v}_{i,t+1} = \mathbf{v}_{j,t+1}$ when $\nabla \mathbf{w}_{i,t}^l = \nabla \mathbf{w}_{j,t}^l$, so the Adam optimizer is also a conformal mapping at this time.

However, $\nabla \mathbf{w}_{i,t}^l \neq \nabla \mathbf{w}_{j,t}^l$ in general. Let $\kappa = \|\nabla \mathbf{w}_{i,t}^l - \nabla \mathbf{w}_{j,t}^l\|_2$ be the degree of difference between $\nabla \mathbf{w}_{i,t}^l$ and $\nabla \mathbf{w}_{j,t}^l$. Obviously, $\kappa \geq 0$. As κ increases, the degree of difference between $\nabla \mathbf{w}_{i,t}^l$ and $\nabla \mathbf{w}_{j,t}^l$ becomes larger. As κ decreases, the degree of difference becomes smaller. When κ approaches 0, the optimizer approaches a conformal mapping.

If the optimizer is SGD, the degree of shift difference between parameter vectors $\mathbf{w}_{i,t}^l$ and $\mathbf{w}_{j,t}^l$ is $\eta_t \kappa$. If the optimizer is the Adam, The degree of difference in translation transformations between $\mathbf{w}_{i,t}^l$ and $\mathbf{w}_{j,t}^l$ is determined

by $\mathbf{m}_{i,t+1}$, $\mathbf{m}_{j,t+1}$, $\mathbf{v}_{i,t+1}$, and $\mathbf{v}_{j,t+1}$. Where $\mathbf{m}_{i,t+1}$ and $\mathbf{m}_{j,t+1}$ represent translation transformations of $\beta_1 \mathbf{m}_{i,t}$ and $\beta_1 \mathbf{m}_{j,t}$, and $\mathbf{v}_{i,t+1}$ and $\mathbf{v}_{j,t+1}$ represent translation transformations of $\beta_2 \mathbf{v}_{i,t}$ and $\beta_2 \mathbf{v}_{j,t}$. For any parameter vector, the Adam optimizer needs to undergo three translation transformations, and $\|\nabla \mathbf{w}_{i,t}^l \odot \nabla \mathbf{w}_{i,t}^l - \nabla \mathbf{w}_{j,t}^l \odot \nabla \mathbf{w}_{j,t}^l\|_2 \geq \kappa$. Therefore, the degree of difference in translation vectors between $\mathbf{w}_{i,t}^l$ and $\mathbf{w}_{j,t}^l$ will be significantly greater than $\eta_t \kappa$. The ability of the SGD optimizer for conformal mapping surpasses that of the Adam optimizer. \square

2 Proof of Theorem 3

Proof. For a FNN, the forward propagation process of the neural network is as follows:

$$\mathbf{z}^l = W^l \mathbf{x}^{l-1} + \mathbf{b}^l, \mathbf{x}^l = f(\mathbf{z}^l). \quad (4)$$

After adding residual connections and normalization,

$$\mathbf{x}^l = \text{BN}(\mathbf{x}^{l-1} + f(\mathbf{z}^l)), \quad (5)$$

where $\hat{\mathbf{x}}^l = \mathbf{x}^{l-1} + f(\mathbf{z}^l)$, $\mathbf{x}^l = \text{BN}(\hat{\mathbf{x}}^l)$. Let

$$\mu = \frac{1}{N} \sum_i \hat{x}_i^l, \sigma^2 = \frac{1}{N} \sum_i (\hat{x}_i^l - \mu)^2. \quad (6)$$

Then, $x_i^l = \frac{\hat{x}_i^l - \mu}{\sigma}$. We take the batch normalization as an example, the result of layer normalization are similar. Next, we calculate the relevant partial derivatives. For any $j \in \{1, 2, \dots, N\}$,

$$\frac{\partial \mu}{\partial \hat{x}_j^l} = \frac{1}{N} \quad (7)$$

$$\begin{aligned} \frac{\partial \sigma}{\partial \hat{x}_j^l} &= \frac{1}{2\sigma N} \sum_i \left[2(\hat{x}_i^l - \mu) \left(\Phi_{ij} - \frac{1}{N} \right) \right] \\ &= \frac{1}{N\sigma} \sum_i [(\hat{x}_i^l - \mu) \Phi_{ij}] - \frac{1}{N^2\sigma} \sum_i (\hat{x}_i^l - \mu) \\ &= \frac{\hat{x}_j^l - \mu}{N\sigma} - 0 = \frac{x_j^l}{N}, \end{aligned} \quad (8)$$

where

$$\Phi_{ij} = \begin{cases} 1 & i = j \\ 0 & i \neq j \end{cases}. \quad (9)$$

50 Then

$$\begin{aligned}\frac{\partial x_i^l}{\partial \hat{x}_j^l} &= \frac{\left(\Phi_{ij} - \frac{\partial \mu}{\partial \hat{x}_j^l}\right) \sigma - (\hat{x}_j^l - \mu) \frac{\partial \sigma}{\partial \hat{x}_j^l}}{\sigma^2} \\ &= \frac{\left(\Phi_{ij} - \frac{1}{N}\right) - x_i^l \frac{\partial \sigma}{\partial \hat{x}_j^l}}{\sigma} = \frac{\left(\Phi_{ij} - \frac{1}{N}\right) - \frac{x_i^l x_j^l}{N}}{\sigma} \quad (10) \\ &= \frac{(N\Phi_{ij} - 1) - x_i^l x_j^l}{N\sigma}.\end{aligned}$$

$$\begin{aligned}\frac{\partial L}{\partial \hat{x}_j^l} &= \sum_i \frac{\partial L}{\partial x_i^l} \frac{\partial x_i^l}{\partial \hat{x}_j^l} \\ &= \frac{\sum_i \left((x_i^l)' [(N\Phi_{ij} - 1) - x_i^l x_j^l] \right)}{N\sigma} \\ &= \frac{N (x_j^l)' - \sum_i (x_i^l)' - x_j^l \sum_i \left((x_i^l)' x_i^l \right)}{N\sigma} \\ &= \frac{(x_j^l)'}{\sigma} - \frac{\sum_i (x_i^l)'}{N\sigma} - \frac{x_j^l \sum_i \left((x_i^l)' x_i^l \right)}{N\sigma} \\ &= \frac{1}{\sigma} \left[(x_j^l)' - \frac{\sum_i (x_i^l)'}{N} - \frac{x_j^l \sum_i \left((x_i^l)' x_i^l \right)}{N} \right]. \quad (11)\end{aligned}$$

51 Let $\mathbf{x}^l = \{x_1^l, x_2^l, \dots, x_N^l\}$, then

$$\begin{aligned}(\hat{\mathbf{x}}^l)' &= \frac{\partial L}{\partial \hat{\mathbf{x}}^l} \\ &= \frac{1}{\sigma} \left[(\mathbf{x}^l)' - \frac{((\mathbf{x}^l)', \mathbf{1})}{N} \mathbf{1} - \frac{((\mathbf{x}^l)', \mathbf{x}^l)}{N} \mathbf{x}^l \right], \quad (12)\end{aligned}$$

52 where (\cdot, \cdot) is the inner product in N dimensions. Since

$$\begin{aligned}\|\vec{1}\|^2 &= N \\ \|\mathbf{x}^l\|^2 &= \sum_i (x_i^l)^2 = \sum_i \frac{N (\hat{x}_i^l - \mu)^2}{\sum_j (\hat{x}_j^l - \mu)^2} = N, \quad (13)\end{aligned}$$

53 we can obtain

$$\begin{aligned}(\hat{\mathbf{x}}^l)' &= \frac{\partial L}{\partial \hat{\mathbf{x}}^l} \\ &= \frac{1}{\sigma} \left[(\mathbf{x}^l)' - \frac{((\mathbf{x}^l)', \mathbf{1})}{(\mathbf{1}, \mathbf{1})} \mathbf{1} - \frac{((\mathbf{x}^l)', \mathbf{x}^l)}{(\mathbf{x}^l, \mathbf{x}^l)} \mathbf{x}^l \right]. \quad (14)\end{aligned}$$

54 For neurons in layer l , $\mathbf{z}^l = W^l \mathbf{x}^{l-1} + \mathbf{b}^l$, then

$$\begin{aligned}\nabla W^l &= \frac{\partial L}{\partial \mathbf{z}^l} \frac{\partial \mathbf{z}^l}{\partial \mathbf{w}^l} = \frac{\partial L}{\partial \hat{\mathbf{x}}^l} \frac{\partial \hat{\mathbf{x}}^l}{\partial \mathbf{z}^l} \frac{\partial \mathbf{z}^l}{\partial \mathbf{w}^l} \\ &= \frac{1}{\sigma} f'(\mathbf{z}^l) \odot \left[(\mathbf{x}^l)' - \frac{((\mathbf{x}^l)', \mathbf{1})}{(\mathbf{1}, \mathbf{1})} \mathbf{1} - \frac{((\mathbf{x}^l)', \mathbf{x}^l)}{(\mathbf{x}^l, \mathbf{x}^l)} \mathbf{x}^l \right] (\mathbf{x}^{l-1})^T, \quad (15)\end{aligned}$$

\odot is the element-wise multiplication. $(\mathbf{x}^l)' - \frac{((\mathbf{x}^l)', \mathbf{1})}{(\mathbf{1}, \mathbf{1})} \mathbf{1} - \frac{((\mathbf{x}^l)', \mathbf{x}^l)}{(\mathbf{x}^l, \mathbf{x}^l)} \mathbf{x}^l$ is the projection transformation of vector $(\mathbf{x}^l)'$ in the directions of vectors $\mathbf{1}$ and \mathbf{x}^l , abbreviated as $P_{\mathbf{1}^\perp P_{\mathbf{x}^l}^\perp}(\mathbf{x}^l)'$. When the batch normalization and residual connections are also introduced to the neurons in the $l+1$ layer,

$$\hat{\mathbf{x}}^{l+1} = \mathbf{x}^l + f(W^{l+1} \mathbf{x}^l + \mathbf{b}^{l+1}), \mathbf{x}^{l+1} = \text{BN}(\hat{\mathbf{x}}^{l+1}). \quad (16)$$

For the equation (15),

$$\begin{aligned}(\mathbf{x}^l)' &= \frac{\partial L}{\partial \mathbf{x}^l} = \frac{\partial \hat{\mathbf{x}}^{l+1}}{\partial \mathbf{x}^l} \frac{\partial L}{\partial \hat{\mathbf{x}}^{l+1}} \\ &= \frac{1}{\sigma} (\mathbf{1} + f'(\mathbf{z}^{l+1})) \\ &\odot W^{l+1} \left[(\mathbf{x}^{l+1})' - \frac{((\mathbf{x}^{l+1})', \mathbf{1})}{(\mathbf{1}, \mathbf{1})} \mathbf{1} - \frac{((\mathbf{x}^{l+1})', \mathbf{x}^{l+1})}{(\mathbf{x}^{l+1}, \mathbf{x}^{l+1})} \mathbf{x}^{l+1} \right] \\ &= \frac{1}{\sigma} (\mathbf{1} + f'(\mathbf{z}^{l+1})) \odot W^{l+1} P_{\mathbf{1}^\perp P_{\mathbf{x}^{l+1}}^\perp}(\mathbf{x}^{l+1})', \quad (17)\end{aligned}$$

so,

$$\begin{aligned}\nabla W^l &= \frac{1}{\sigma^l} f'(\mathbf{z}^l) \odot P_{\mathbf{1}^\perp P_{\mathbf{x}^l}^\perp} \\ &\left[\frac{1}{\sigma^{l+1}} (\mathbf{1} + f'(\mathbf{z}^{l+1})) \odot W^{l+1} P_{\mathbf{1}^\perp P_{\mathbf{x}^{l+1}}^\perp}(\mathbf{x}^{l+1})' \right] (\mathbf{x}^{l-1})^T. \quad (18)\end{aligned}$$

At the same time,

$$\begin{aligned}\delta_{NR}^l &= \frac{1}{\sigma^l} f'(\mathbf{z}^l) \odot P_{\mathbf{1}^\perp P_{\mathbf{x}^l}^\perp} \\ &\left[\frac{1}{\sigma^{l+1}} (\mathbf{1} + f'(\mathbf{z}^{l+1})) \odot W^{l+1} P_{\mathbf{1}^\perp P_{\mathbf{x}^{l+1}}^\perp}(\mathbf{x}^{l+1})' \right], \quad (19)\end{aligned}$$

where σ^l and σ^{l+1} are variance of neurons in layer l and $l+1$. If the neural network does not incorporate the normalization and residual connections,

$$\begin{aligned}\mathbf{x}^l &= \hat{\mathbf{x}}^l = f(\mathbf{z}^l) = f(W^l \mathbf{x}^{l-1} + \mathbf{b}^l) \\ \nabla W^l &= \delta^l (\mathbf{x}^{l-1})^T \\ \delta^l &= f'(\mathbf{z}^l) \odot (W^{l+1} \delta^{l+1}). \quad (20)\end{aligned}$$

If $\nabla \mathbf{w}_i^l$ and $\nabla \mathbf{w}_j^l$ are any two row vectors in ∇W^l , so without adding normalization and residual connections,

$$\begin{aligned}\nabla \mathbf{w}_i^{c,l} &= x_i^{l-1} \delta^l \\ \nabla \mathbf{w}_j^{c,l} &= x_j^{l-1} \delta^l. \quad (21)\end{aligned}$$

69 After adding normalization and residual connections,

$$\begin{aligned}\nabla \mathbf{w}_{NR,i}^{c,l} &= x_{NR,i}^{l-1} \delta_{NR}^l \\ \nabla \mathbf{w}_{NR,j}^{c,l} &= x_{NR,j}^{l-1} \delta_{NR}^l.\end{aligned}\quad (22)$$

70 Normalization makes $x_{NR,i}^{l-1}$ and $x_{NR,j}^{l-1}$ belong to a distribu-
71 tion with an expected value of 0 and a variance of 1. Compared to x_i^{l-1} and x_j^{l-1} , $x_{NR,i}^{l-1}$ and $x_{NR,j}^{l-1}$ are closer, resulting
72 in smaller differences between $\nabla \mathbf{w}_{NR,i}^{c,l}$ and $\nabla \mathbf{w}_{NR,j}^{c,l}$. This
73 makes the parameter updates closer to translational transformations, thereby improving the optimizer ability for conformal mappings. In the same way, if $\nabla \mathbf{w}_i^l$ and $\nabla \mathbf{w}_j^l$ are any two
74 column vectors in ∇W^l , then

$$\begin{aligned}\nabla \mathbf{w}_{NR,i}^{r,l} &= \delta_{NR,i}^l (\mathbf{x}_{NR}^{l-1})^T \\ \nabla \mathbf{w}_{NR,j}^{r,l} &= \delta_{NR,j}^l (\mathbf{x}_{NR}^{l-1})^T,\end{aligned}\quad (23)$$

$$\begin{aligned}\nabla \mathbf{w}_i^{r,l} &= \delta_i^l (\mathbf{x}^{l-1})^T \\ \nabla \mathbf{w}_j^{r,l} &= \delta_j^l (\mathbf{x}^{l-1})^T.\end{aligned}\quad (24)$$

75 Equation (19) indicates that, under the influence of vectors
76 $\mathbf{1}$ and variance σ , the values of $\delta_{NR,i}^l$ are closer to $\delta_{NR,j}^l$
77 compared to δ_i^l and δ_j^l . Therefore, the optimizer exhibits better conformal capabilities. In conclusion, normalization and residual connections generally enhance the optimizer conformal capability. \square

84 3 The Explanation for LoRA Maintaining 85 Orthogonality of Parameter Vectors

86 For a pre-trained weight matrix $W_0 \in R^{d \times k}$, the core im-
87 provement of LoRA is that it constrains the update of W_0 by
88 representing the latter with a low-rank decomposition

$$W_0 + \Delta W = W_0 + BA, \quad (25)$$

89 where $B \in R^{d \times r}$, $A \in R^{r \times k}$, and the rank $r \ll \min(d, k)$.
90 For any two vectors in the parameter matrix W_0 , the essence
91 of Equation (25) remains a translation transformation of vec-
92 tors. In Equation (25), ΔW is decomposed into the prod-
93 uct of two low-rank matrices, B and A , yet matrices B and
94 A are still computed from the gradients of parameters. The
95 normalization and residual connections reduce the differences
96 between gradient vectors, hence diminishing the distinctions
97 between any vectors within matrices B and A . Consequently,
98 the results of Equation (25) for any two vectors in W_0 tend
99 to undergo translation transformations along the same direc-
100 tion, ensuring that the orthogonality between parameter vec-
101 tors is not significantly altered in the network fine-tuning with
102 LoRA.

103 4 Notations Introduced Previously and Some 104 Additional Experiments

Notation	Type	Definition
l	Integer	Number of network layers
N	Integer	Batch size
M_l	Integer	Number of neurons in layer l
k	Integer	Output dimension
W^l	$R^{M_l \times M_{l-1}}$	Weight from layer $l-1$ to l
\mathbf{b}^l	R^{M_l}	Bias from layer $l-1$ to l
\mathbf{x}^l	R^{M_l}	Neurons of layer l
\mathbf{z}^l	R^{M_l}	Neurons of layer l not activated
I_k	$R^{k \times k}$	Identity matrix of size k
$\mathbf{1}_k$	R^k	All-ones vector

Table 1: Notations introduced previously.

Optimizers	PII Layer_1 \uparrow	PII Layer_2 \downarrow	PII Layer_3 \downarrow
SGD	1.0/0.858	0.070/0.121	0.050/0.062
Adam	1.0/0.333	0.070/0.324	0.050/0.152
AdamW	1.0/0.352	0.070/0.311	0.050/0.166
RMSprop	1.0/0.384	0.070/0.309	0.050/0.154

Table 2: Conformal capabilities of optimizers in CNNs. The datasets is CIFAR-100. Similar to the results on CIFAR-10, the conformal capability of the SGD optimizer is the best, while the conformal capabilities of Adam and AdamW optimizers are comparable.

Functions/ Optimizers	PII Original	PII \downarrow Norm_Res	ACC % Original	ACC % \uparrow Norm_Res
GELU/Adam	0.212	0.154	39.38	49.43
GELU/SGD	0.194	0.166	51.59	53.34

Table 3: The influence of the combination of normalization and residual connections on the orthogonality of parameter vectors and testing accuracy in CNNs. The datasets is CIFAR-100. For Adam and SGD optimizers, normalization and residual connections can significantly reduce PII. Therefore, it enhances the orthogonality of parameter vectors, simultaneously improving the network’s capability.

Optimizers	layer_1 \downarrow	layer_2 \downarrow	layer_3 \downarrow
Original	0.131	0.148	0.153
Normalization	0.100	0.107	0.128
Residual	0.111	0.137	0.147
Norm_Residual	0.079	0.090	0.107

Table 4: Changes in PII after separately adding normalization and residual connections. "Original" represents the results without the addition of normalization and residual connections; "Normalization" represents the results with only normalization added; "Residual" represents the results with only residual connections added, and "Norm_Residual" represents the results with normalization and residual connections added simultaneously. Experimental results indicate that adding only normalization or residual connections can also reduce PII, but the most significant reduction occurs when both normalization and residual connections are added simultaneously, resulting in the lowest PII.

105 Deep learning models are trained with two Intel Core i9
106 Processors (5.80 GHz), and 2 NVIDIA GeForce RTX 4090
107 GPU (24GB). For the MNIST datasets, we train an MLP
108 with three hidden layers of size 500. For the CIFAR-10 and
109 CIFAR-100 datasets, we trained a CNN with 4 convolution
110 layers and a FNN. For the WikiText-2 datasets, we trained
111 a network consisting of six Transformers. For the PLM and
112 LLM, we directly invoke the pre-trained models, eliminating
113 the need for retraining. All experimental parameters are in-
114 cluded in the source code we provide.