# Appendix of "The Orthogonality of Weight Vectors: the Secret Key to Normalization and Residual Connections"

## 1 Proof of Theorem 2

*Proof.* When using the SGD to update parameters,

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

$\eta_t$ 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

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

$\mathbf{w}_{i,t+1}^l$ 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,

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

Where $\odot$ is the multiplication of corresponding elements. 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 is also a conformal mapping.

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

For a fixed $\kappa$, the degree of shift difference between parameter vectors $\mathbf{w}_{i,t}^l$ and $\mathbf{w}_{j,t}^l$ is $\eta_t \kappa$ when the optimizer is SGD. 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 $\beta_2 \mathbf{v}_{j,t}$, and $\left\| \nabla \mathbf{w}_{i,t}^l \odot \nabla \mathbf{w}_{i,t}^l - \nabla \mathbf{w}_{i,t}^l \odot \nabla \mathbf{w}_{i,t}^l \right\|_2 \geq \kappa$. 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:

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

After adding residual connections and normalization (Here, taking batch normalization as an example, the result of layer normalization are similar.),

$$\mathbf{x}^l = BN\left(\mathbf{x}^{l-1} + f\left(\mathbf{z}^l\right)\right) \tag{5}$$

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

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

Then, $x_i^l = \frac{\hat{x}_i^l - \mu}{\sigma}$.. Next, we calculate the relevant partial derivatives.

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

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

$$\frac{\partial x_i^l}{\partial \hat{x}_j^l} = \frac{\left(\delta_{ij} - \frac{\partial \mu}{\partial \hat{x}_j^l}\right)\sigma - \left(\hat{x}_j^l - \mu\right)\frac{\partial \sigma}{\partial \hat{x}_j^l}}{\sigma^2}$$

$$= \frac{\left(\delta_{ij} - \frac{1}{N}\right) - x_i^l \frac{\partial \sigma}{\partial \hat{x}_j^l}}{\sigma} = \frac{\left(\delta_{ij} - \frac{1}{N}\right) - \frac{x_i^l x_j^l}{N}}{\sigma} \quad (9)$$

$$= \frac{(N\delta_{ij} - 1) - x_i^l x_j^l}{N\sigma}.$$

$$\frac{\partial L}{\partial \hat{x}_j^l} = \sum_i \frac{\partial L}{\partial x_i^l}\frac{\partial x_i^l}{\partial \hat{x}_j^l}$$

$$= \frac{\sum_i \left((x_i^l)'\left[(N\delta_{ij} - 1) - x_i^l x_j^l\right]\right)}{N\sigma}$$

$$= \frac{N\left(x_j^l\right)' - \sum_i \left(x_i^l\right)' - x_j^l \sum_i \left(\left(x_i^l\right)' x_i^l\right)}{N\sigma} \quad (10)$$

$$= \frac{\left(x_j^l\right)'}{\sigma} - \frac{\sum_i \left(x_i^l\right)'}{N\sigma} - \frac{x_j^l \sum_i \left(\left(x_i^l\right)' x_i^l\right)}{N\sigma}$$

$$= \frac{1}{\sigma}\left[\left(x_j^l\right)' - \frac{\sum_i \left(x_i^l\right)'}{N} - \frac{x_j^l \sum_i \left(\left(x_i^l\right)' x_i^l\right)}{N}\right]$$

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

$$\left(\hat{\mathbf{x}}^l\right)' = \frac{\partial L}{\partial \hat{\mathbf{x}}^l}$$

$$= \frac{1}{\sigma}\left[\left(\mathbf{x}^l\right)' - \frac{\left(\left(\mathbf{x}^l\right)', \mathbf{1}\right)}{N}\mathbf{1} - \frac{\left(\left(\mathbf{x}^l\right)', \mathbf{x}^l\right)}{N}\mathbf{x}^l\right] \quad (11)$$

where $(\cdot, \cdot)$ is the inner product in $N$ dimensions. Science

$$\left\|\vec{1}\right\|^2 = N$$

$$\left\|\mathbf{x}^l\right\|^2 = \sum_i \left(x_i^l\right)^2 = \sum_i \frac{N\left(\hat{x}_i^l - \mu\right)^2}{\sum_j \left(\hat{x}_j^l - \mu\right)^2} = N \quad (12)$$

$$\left(\hat{\mathbf{x}}^l\right)' = \frac{\partial L}{\partial \hat{\mathbf{x}}^l}$$

$$= \frac{1}{\sigma}\left[\left(\mathbf{x}^l\right)' - \frac{\left(\left(\mathbf{x}^l\right)', \mathbf{1}\right)}{(\mathbf{1}, \mathbf{1})}\mathbf{1} - \frac{\left(\left(\mathbf{x}^l\right)', \mathbf{x}^l\right)}{(\mathbf{x}^l, \mathbf{x}^l)}\mathbf{x}^l\right] \quad (13)$$

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

$$\nabla W^l = \frac{\partial L}{\partial \mathbf{z}^l}\frac{\partial \mathbf{z}^l}{\partial \mathbf{w}^l} = \frac{\partial L}{\partial \hat{\mathbf{x}}^l}\frac{\partial \hat{\mathbf{x}}^l}{\partial \mathbf{z}^l}\frac{\partial \mathbf{z}^l}{\partial \mathbf{w}^l}$$

$$= \frac{1}{\sigma}f'\left(\mathbf{z}^l\right) \odot \left[\left(\mathbf{x}^l\right)' - \frac{\left(\left(\mathbf{x}^l\right)', \mathbf{1}\right)}{(\mathbf{1}, \mathbf{1})}\mathbf{1} - \frac{\left(\left(\mathbf{x}^l\right)', \mathbf{x}^l\right)}{(\mathbf{x}^l, \mathbf{x}^l)}\mathbf{x}^l\right]\left(\mathbf{x}^{l-1}\right)^T$$

$$(14)$$

$\odot$ is the element-wise multiplication. $\left(\mathbf{x}^l\right)' - \frac{\left(\left(\mathbf{x}^l\right)', \mathbf{1}\right)}{(\mathbf{1}, \mathbf{1})}\mathbf{1} - \frac{\left(\left(\mathbf{x}^l\right)', \mathbf{x}^l\right)}{(\mathbf{x}^l, \mathbf{x}^l)}\mathbf{x}^l$ is the projection transformation of vector $\left(\mathbf{x}^l\right)'$ in the directions of vectors $\mathbf{1}$ and $\mathbf{x}^l$, abbreviated as $P_{\mathbf{1}\perp}P_{\mathbf{x}^l\perp}\left(\mathbf{x}^l\right)'$. Let $\delta_{NR}^l = f'\left(\mathbf{z}^l\right) \odot \left(W^{l+1}\delta^{l+1}\right)$, $NR$ represents the addition of normalization and residual connections.

When the batch normalization and residual connections are also introduced to the neurons in the $l+1$ layer,

$$\hat{\mathbf{x}}^{l+1} = \mathbf{x}^l + f\left(W^{l+1}\mathbf{x}^l + b^{l+1}\right), \mathbf{x}^{l+1} = BN\left(\hat{\mathbf{x}}^{l+1}\right) \quad (15)$$

For the equation (14),

$$\left(\mathbf{x}^l\right)' = \frac{\partial L}{\partial \mathbf{x}^l} = \frac{\partial \hat{\mathbf{x}}^{l+1}}{\partial \mathbf{x}^l}\frac{\partial L}{\partial \hat{\mathbf{x}}^{l+1}}$$

$$= \frac{1}{\sigma}\left(\mathbf{1} + f'\left(\mathbf{z}^{l+1}\right)\right)$$

$$\odot W^{l+1}\left[\left(\mathbf{x}^{l+1}\right)' - \frac{\left(\left(\mathbf{x}^{l+1}\right)', \mathbf{1}\right)}{(\mathbf{1}, \mathbf{1})}\mathbf{1} - \frac{\left(\left(\mathbf{x}^{l+1}\right)', \mathbf{x}^{l+1}\right)}{(\mathbf{x}^{l+1}, \mathbf{x}^{l+1})}\mathbf{x}^{l+1}\right]$$

$$= \frac{1}{\sigma}\left(\mathbf{1} + f'\left(\mathbf{z}^{l+1}\right)\right) \odot W^{l+1}P_{\mathbf{1}\perp}P_{\mathbf{x}^{l+1}\perp}\left(\mathbf{x}^{l+1}\right)'$$

$$(16)$$

So,

$$\nabla W^l = \frac{1}{\sigma^l}f'\left(\mathbf{z}^l\right) \odot P_{\mathbf{1}\perp}P_{\mathbf{x}^l\perp}$$

$$\left[\frac{1}{\sigma^{l+1}}\left(\mathbf{1} + f'\left(\mathbf{z}^{l+1}\right)\right) \odot W^{l+1}P_{\mathbf{1}\perp}P_{\mathbf{x}^{l+1}\perp}\left(\mathbf{x}^{l+1}\right)'\right]\left(\mathbf{x}^{l-1}\right)^T$$

$$(17)$$

At the same time,

$$\delta_{NR}^l = \frac{1}{\sigma^l}f'\left(\mathbf{z}^l\right) \odot P_{\mathbf{1}\perp}P_{\mathbf{x}^l\perp}$$

$$\left[\frac{1}{\sigma^{l+1}}\left(\mathbf{1} + f'\left(\mathbf{z}^{l+1}\right)\right) \odot W^{l+1}P_{\mathbf{1}\perp}P_{\mathbf{x}^{l+1}\perp}\left(\mathbf{x}^{l+1}\right)'\right] \quad (18)$$

If the neural network does not incorporate the normalization and residual connections,

$$\mathbf{x}^l = \hat{\mathbf{x}}^l = f\left(\mathbf{z}^l\right) = f\left(W^l\mathbf{x}^{l-1} + \mathbf{b}^l\right)$$

$$\nabla W^l = \delta^l \left(\mathbf{x}^{l-1}\right)^T \quad (19)$$

$$\delta^l = f'\left(\mathbf{z}^l\right) \odot \left(W^{l+1}\delta^{l+1}\right)$$

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

$$\nabla \mathbf{w}_i^{c,l} = x_i^{l-1}\delta^l$$

$$\nabla \mathbf{w}_j^{c,l} = x_j^{l-1}\delta^l \quad (20)$$

After adding normalization and residual connections,

$$\nabla \mathbf{w}_{NR,i}^{c,l} = x_{NR,i}^{l-1}\delta_{NR}^l$$

$$\nabla \mathbf{w}_{NR,j}^{c,l} = x_{NR,j}^{l-1}\delta_{NR}^l \quad (21)$$

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

$$
\begin{aligned}
\nabla\mathbf{w}^{r,l}_{NR,i} = \delta^l_{NR,i} \left(\mathbf{x}^{l-1}_{NR}\right)^T \\
\nabla\mathbf{w}^{r,l}_{NR,j} = \delta^l_{NR,j} \left(\mathbf{x}^{l-1}_{NR}\right)^T
\end{aligned}
\tag{22}
$$

$$
\begin{aligned}
\nabla\mathbf{w}^{r,l}_i = \delta^l_i \left(\mathbf{x}^{l-1}\right)^T \\
\nabla\mathbf{w}^{r,l}_j = \delta^l_j \left(\mathbf{x}^{l-1}\right)^T
\end{aligned}
\tag{23}
$$

Formula (18) indicates that, under the influence of vectors **1** and variance $\sigma$, the values of $\delta^l_{NR,i}$ are closer to $\delta^l_{NR,j}$ compared to $\delta^l_i$ and $\delta^l_j$. Therefore, the optimizer exhibits better conformal capabilities. In conclusion, normalization and residual links generally enhance the optimizer conformal mapping ability ☐

## 3 The explanation for LoRA maintaining orthogonality of parameter vectors.

For a pre-trained weight matrix $W_0 \in R^{d\times k}$, the core improvement of LoRA is that constrain the update of $W_0$ by representing the latter with a low-rank decomposition

$$
W_0 + \Delta W = W_0 + BA
\tag{24}
$$

where $B \in R^{d\times r}$, $A \in R^{r\times k}$, and the rank $r \ll \min(d, k)$. For any two vectors in the parameter matrix $W_0$, the essence of Equation (24) remains a translation transformation of vectors. In Equation (24), $\Delta W$ is decomposed into the product of two low-rank matrices, $B$ and $A$, yet matrices $B$ and $A$ are still computed from the gradients of parameters. The normalization and residual links reduce the differences between gradient vectors, hence diminishing the distinctions between any vectors within matrices $B$ and $A$. Consequently, the results of Equation (24) for any two vectors in $W_0$ tend to undergo translation transformations along the same direction, ensuring that the orthogonality between parameter vectors is not significantly altered in the network fine-tuning with LoRA.

## 4 Notations Introduced Previously and Some 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 by the activation funtion |
| $I_k$ | $R^{k \times k}$ | Identity matrix of size $k$ |
| $\mathbf{1}_k$ | $R^n$ | All-ones vector |

Table 1: Notations introduced previously.

| Optimizers | PII Layer_1↑ | PII Layer_2↓ | PII Layer_3↓ |
|---|---|---|---|
| SGD | 1.0/0.858 | 0.0.070/0.121 | 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 mapping capability of the SGD optimizer is the best, while the conformal mapping capabilities of Adam and AdamW are comparable

| Functions/ Optimizers | PII Original | PII ↓ Norm_Res | ACC % Original | ACC % ↑ 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, normalization and residual connections can significantly reduce PII.Therefore, it enhances the orthogonality of parameter vectors, simultaneously improving the network's capability.

| Optimizers | Hidden layer_1 | Hidden layer_2 | Hidden layer_3 |
|---|---|---|---|
| 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 both 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.