# Training Networks in Null Space of Feature Covariance for Continual Learning (Supplementary Material)

Anonymous CVPR 2021 submission

Paper ID 656

In this supplementary material, we first introduce additional notations, and give the proof of Lemma 1 in Appendix A. Then we discuss the reason why the parameter update satisfying Condition 2 is the descent direction in Appendix B. Finally, in Appendix C, we prove that $\langle \Delta \mathbf{w}_{t,s}, \mathbf{g}_{t,s} \rangle \geq 0$ claimed in Sec. 4.1 of the manuscript.

## Notations

We first introduce additional notations here. When feeding data $X_p$ from task $\mathcal{T}_p$ ($p \leq t$) to $f$ with parameters $\mathbf{w}_{t,s}$, the input feature and output feature at the $l$-th linear layer are denoted as $X_{p,t,s}^l$ and $O_{p,t,s}^l$ respectively, then

$$O_{p,t,s}^l = X_{p,t,s}^l w_{t,s}^l, \quad X_{p,t,s}^{l+1} = \sigma_l(O_{p,t,s}^l)$$

with $X_{p,t,s}^1 = X_p$. In addition, by denoting the learning rate as $\alpha$, we have

$$w_{t,s}^l = w_{t,s-1}^l - \alpha \Delta w_{t,s-1}^l, l = 1, \ldots, L.$$

## Appendix A

In this appendix, we show the proof of Lemma 1 in the manuscript. Lemma 1 tells us that, when we train network on task $\mathcal{T}_t$, the network retains its training loss on data $X_p$ in the training process, if the network parameter update satisfies Eqn. (1) at each training step. We first recall Lemma 1 as follows, then give the proof.

**Lemma 1.** *Given the data $X_p$ from task $\mathcal{T}_p$, and the network $f$ with $L$ linear layers is trained on task $\mathcal{T}_t$ ($t > p$). If network parameter update $\Delta w_{t,s}^l$ lies in the null space of $X_{p,t-1}^l$, i.e.,*

$$X_{p,t-1}^l \Delta w_{t,s}^l = 0, \tag{1}$$

*at each training step $s$, for the $l$-th layer of $f$ ($l = 1, \ldots, L$), we have $X_{p,t}^l = X_{p,t-1}^l$ and $f(X_p, \tilde{\mathbf{w}}_{t-1}) = f(X_p, \tilde{\mathbf{w}}_t)$.*

*Proof.* The proof is based on the recursive structure of network and iterative training process. We first prove that $X_{p,t,1}^l = X_{p,t-1}^l$ and $f(X_p, \mathbf{w}_{t,1}) = f(X_p, \tilde{\mathbf{w}}_{t-1})$ hold for $s = 1$, and then illustrate that $X_{p,t,s}^l = X_{p,t-1}^l$ and $f(X_p, \mathbf{w}_{t,s}) = f(X_p, \tilde{\mathbf{w}}_{t-1})$ hold for each $s > 1$, which suggests that Lemma 1 holds.

When $s = 1$, considering that we initialize parameters $\mathbf{w}_{t,0} = \tilde{\mathbf{w}}_{t-1}$, we have

$$X_{p,t,0}^l = X_{p,t-1}^l, \quad O_{p,t,0}^l = O_{p,t-1}^l. \tag{2}$$

Therefore, at the first layer ($l = 1$) where $X_{p,t,1}^1 = X_{p,t,0}^1 = X_{p,t-1}^1$ (all of them equal to $X_p$ when $l = 1$),

$$\begin{aligned}
O_{p,t,1}^1 &= X_{p,t,1}^1 w_{t,1}^1 \\
&= X_{p,t,0}^1 (w_{t,0}^1 - \alpha \Delta w_{t,0}^1) \\
&= X_{p,t,0}^1 w_{t,0}^1 - \alpha X_{p,t-1}^1 \Delta w_{t,0}^1 \\
&= X_{p,t,0}^1 w_{t,0}^1 \\
&= O_{p,t,0}^1, \tag{3}
\end{aligned}$$

where the fourth equation holds due to Eqn. (1). Furthermore, we have

$$X_{p,t,1}^2 = \sigma_1(O_{p,t,1}^1) = \sigma_1(O_{p,t,0}^1) = X_{p,t,0}^2 = X_{p,t-1}^2, \tag{4}$$

*i.e.*, the input feature $X_{p,t,1}^2$ equals to $X_{p,t-1}^2$ at the second linear layer, based on which, we can recursively prove that

$$O_{p,t,1}^l = O_{p,t,0}^l = O_{p,t-1}^l$$

and

$$X_{p,t,1}^l = X_{p,t,0}^l = X_{p,t-1}^l$$

for $l = 3, \ldots, L$ by replacing $l = 1$ with $l = 2, \ldots, L$ in Eqns. (3) and (4), then we have $f(X_p, \mathbf{w}_{t,1}) = f(X_p, \tilde{\mathbf{w}}_{t-1})$.

We now have proved that $X_{p,t,s}^l = X_{p,t-1}^l$, $O_{p,t,s}^l = O_{p,t-1}^l$ ($l = 1, \ldots, L$) and $f(X_p, \mathbf{w}_{t,s}) = f(X_p, \tilde{\mathbf{w}}_{t-1})$ hold for $s = 1$. Considering the iterative training process, we can prove that

$$X_{p,t,s}^l = X_{p,t-1}^l, \quad O_{p,t,s}^l = O_{p,t-1}^l \quad (l = 1, \ldots, L)$$

and

$$f(X_p, \mathbf{w}_{t,s}) = f(X_p, \tilde{\mathbf{w}}_{t-1})$$

hold for $s = 2, ...$, by repeating the above process with $s = 2, ....$

Finally, we have $X_{p,t}^l = X_{p,t-1}^l$ and $f(X_p, \tilde{\mathbf{w}}_{t-1}) = f(X_p, \tilde{\mathbf{w}}_t)$, since Lemma 1 holds for each $s \geq 1$. □

## Appendix B

We first recall the Condition 2 in the manuscript as follows, then prove that parameter update $\Delta\mathbf{w}_{t,s}$ satisfying condition 2 is the descent direction, *i.e.*, the training loss after updating parameters using $\Delta\mathbf{w}_{t,s}$ will decrease.

**Condition 2** (plasticity)**.** *Assume that the network $f$ is being trained on task $\mathcal{T}_t$, and $\mathbf{g}_{t,s} = \{g_{t,s}^1, \ldots, g_{t,s}^L\}$ denotes the parameter update generated by a gradient-descent training algorithm for training $f$ at training step $s$. $\langle\Delta\mathbf{w}_{t,s}, \mathbf{g}_{t,s}\rangle > 0$ should hold where $\langle\cdot,\cdot\rangle$ represents inner product.*

We now discuss the reason why $\Delta\mathbf{w}_{t,s}$ is the descent direction, if it satisfies condition 2. For clarity, we denote the loss for training network $f$ as $\mathcal{L}(\mathbf{w})$ which ignores the data term with no effect. The discussion can also be found in Lemma 2 of the lecture[1].

By denoting the learning rate as $\alpha$, and $h(\alpha) \triangleq \mathcal{L}(\mathbf{w}_{t,s} - \alpha\Delta\mathbf{w}_{t,s})$, according to Taylor's theorem, we have

$$h(\alpha) = h(0) + \nabla_\alpha h(0) + o(\alpha),$$

*i.e.*,

$$\mathcal{L}(\mathbf{w}_{t,s} - \alpha\Delta\mathbf{w}_{t,s}) = \mathcal{L}(\mathbf{w}_{t,s}) - \alpha\langle\Delta\mathbf{w}_{t,s}, \mathbf{g}_{t,s}\rangle + o(\alpha),$$

where $\frac{|o(\alpha)|}{\alpha} \to 0$ when $\alpha \to 0$. Therefore, there exists $\bar{\alpha} > 0$ such that

$$|o(\alpha)| < \alpha|\langle\Delta\mathbf{w}_{t,s}, \mathbf{g}_{t,s}\rangle|, \ \forall\alpha \in (0, \bar{\alpha}).$$

Together with the condition $\langle\Delta\mathbf{w}_{t,s}, \mathbf{g}_{t,s}\rangle > 0$, we can conclude that $\mathcal{L}(\mathbf{w}_{t,s} - \alpha\Delta\mathbf{w}_{t,s}) < \mathcal{L}(\mathbf{w}_{t,s})$ for all $\alpha \in (0, \bar{\alpha})$. Therefore, parameter update $\Delta\mathbf{w}_{t,s}$ satisfying condition 2 is the descent direction.

## Appendix C

Here, we give the proof of $\langle\Delta\mathbf{w}_{t,s}, \mathbf{g}_{t,s}\rangle \geq 0$ with $\Delta w_{t,s}^l = U_2^l(U_2^l)^\top g_{t,s}^l$, which is claimed in Sec 4.1 of the manuscript. The proof mainly utilizes the properties of Kro-

---

[1] http://www.princeton.edu/~aaa/Public/Teaching/ORF363_COS323/F14/ORF363_COS323_F14_Lec8.pdf

necker product [1, Eqns. (2.10) and (2.13)].

$$\langle\Delta\mathbf{w}_{t,s}, \mathbf{g}_{t,s}\rangle = \sum_{l=1}^L \langle U_2^l(U_2^l)^\top g_{t,s}^l, g_{t,s}^l\rangle$$

$$= \sum_{l=1}^L \mathrm{vec}(U_2^l(U_2^l)^\top g_{t,s}^l I)^\top \mathrm{vec}(g_{t,s}^l)$$

$$= \sum_{l=1}^L \mathrm{vec}((U_2^l)^\top g_{t,s}^l)^\top (I \otimes (U_2^l)^\top)\mathrm{vec}(g_{t,s}^l)$$

$$= \sum_{l=1}^L \mathrm{vec}((U_2^l)^\top g_{t,s}^l)^\top \mathrm{vec}((U_2^l)^\top g_{t,s}^l)$$

$$\geq 0, \tag{5}$$

where $\mathrm{vec}(\cdot)$ is the vectorization of $\cdot$, $I$ is the identity matrix and $\otimes$ is the Kronecker product.

## Appendix D

We now discuss the difference between our algorithm and OWM [2] in details as follows. (1) We provide novel theoretical conditions for the stability and plasticity of network based on feature covariance. (2) The null space of ours is defined as the null space of feature covariance matrix which is easy to be accumulated after each task (refer to Q1 & Alg. 2). While the projection matrix in OWM is $\mathbf{P}_l = \mathbf{I}_l - \mathbf{A}_l(\mathbf{A}_l^\top\mathbf{A}_l + \beta_l\mathbf{I}_l)^{-1}\mathbf{A}_l^\top$ where $\mathbf{A}_l$ consists of all previous features of layer $l$. (3) With the coming of new tasks, our covariance matrix is incrementally updated without approximation error, while $\mathbf{P}_l$ of OWM is updated by recursive least square, where the approximation error of matrix inversion (because of the additionally introduced $\beta_l\mathbf{I}$) will be accumulated. (4) Our approach relies on a hyperparameter $a$ in line 14 of Alg. 2, for approximating the null space of covariance, which can balance the stability and plasticity as discussed in lines 572-579 and Fig. 5. It is easy to set the hyperparameter (line 614 and Figs. 4, 5). But we find that it is hard to tune the hyperparameter $\beta_l$ in OWM for each layer to balance the approximation error and computational stability. (5) Experimental comparison with OWM on three benchmarks are shown in Tabs. 1-3. The ACC of ours are 4.88%, 7.48% and 8.3% higher than OWM with comparable BWT. Please refer to Q4 for comparison on ImageNet with deeper networks. We will clarify these differences by extending the discussions in Sect. 2.

## References

[1] Alexander Graham. *Kronecker products and matrix calculus with applications.* Courier Dover Publications, 2018. 2

[2] Guanxiong Zeng, Yang Chen, Bo Cui, and Shan Yu. Continual learning of context-dependent processing in neural networks. *Nature Machine Intelligence*, 1(8):364–372, 2019. 2