

# MetaFGNet: Supplementary Material

Yabin Zhang, Hui Tang and Kui Jia

School of Electronic and Information Engineering

South China University of Technology, Guangzhou, China

{zhang.yabin, eehuitang}@mail.scut.edu.cn, kuijia@scut.edu.cn

This is the supplementary material of [4] to demonstrate how to evaluate gradients of the meta-learning loss with respect to model parameters.

The gradients of meta-learning loss are uniformly formulated as:

$$[\Delta(\theta_b; \mathcal{T}_j), \Delta(\theta_c^t; \mathcal{T}_j)] = \nabla_{\theta^t} \frac{1}{|\mathcal{T}_j|} L(\mathcal{T}_j; \theta^{t'}) = \nabla_{\theta^{t'}} \frac{1}{|\mathcal{T}_j|} L(\mathcal{T}_j; \theta^{t'}) \left[ \mathbf{I} - \eta \frac{1}{|\mathcal{T}_i|} \left( \frac{\partial^2 L(\mathcal{T}_i; \theta^t)}{\partial(\theta^t)^2} \right) \right], \quad (1)$$

which is easily implemented by TensorFlow [1] but a little bit more complicated with PyTorch [3]. We present to implement it using PyTorch in this supplementary material. Firstly, we expand Equ. 1 as:

$$\nabla_{\theta^{t'}} \frac{1}{|\mathcal{T}_j|} L(\mathcal{T}_j; \theta^{t'}) - \eta \nabla_{\theta^{t'}} \frac{1}{|\mathcal{T}_j|} L(\mathcal{T}_j; \theta^{t'}) \frac{1}{|\mathcal{T}_i|} \left( \frac{\partial^2 L(\mathcal{T}_i; \theta^t)}{\partial(\theta^t)^2} \right),$$

where  $\nabla_{\theta^{t'}} \frac{1}{|\mathcal{T}_j|} L(\mathcal{T}_j; \theta^{t'})$  is directly computed via back propagation. Hence, we focus on calculating:

$$\nabla_{\theta^{t'}} L(\mathcal{T}_j; \theta^{t'}) \left( \frac{\partial^2 L(\mathcal{T}_i; \theta^t)}{\partial(\theta^t)^2} \right). \quad (2)$$

The Hessian matrix in Equ. 2 can be transformed into:

$$\frac{\partial^2 L(\mathcal{T}_i; \theta^t)}{\partial(\theta^t)^2} = \frac{\partial^2 L(\mathcal{T}_i; \theta^t)}{\partial \theta^t \partial(\theta^t)^T} = \frac{\partial}{\partial \theta^t} \left[ \frac{\partial L(\mathcal{T}_i; \theta^t)}{\partial(\theta^t)^T} \right] = \frac{\partial \mathbf{g}^t}{\partial \theta^t}, \quad (3)$$

where  $\theta^t \in \mathbb{R}^{1 \times N}$ <sup>1</sup>,  $\mathbf{g}^t = \frac{\partial L(\mathcal{T}_i; \theta^t)}{\partial(\theta^t)^T} \in \mathbb{R}^{N \times 1}$ ,  $\frac{\partial^2 L(\mathcal{T}_i; \theta^t)}{\partial(\theta^t)^2} \in \mathbb{R}^{N \times N}$ . The  $i$ th row of  $\frac{\partial^2 L(\mathcal{T}_i; \theta^t)}{\partial(\theta^t)^2}$  is rewritten as  $\frac{\partial \mathbf{g}_i^t}{\partial \theta^t}$ , which can be obtained on the basis of the chain rule.

To date, we are able to compute gradients of the meta-learning loss. However, it is computationally expensive to calculate the Hessian matrix. Thus, we propose to conduct gradient computation in a faster way.

For the sake of brevity, we abbreviate  $\nabla_{\theta^{t'}} L(\mathcal{T}_j; \theta^{t'})$  and  $\frac{\partial^2 L(\mathcal{T}_i; \theta^t)}{\partial(\theta^t)^2}$  as  $\mathbf{G}$  and  $\mathbf{H}$  respectively. The Equ. 2 can be rewritten as:

$$\nabla_{\theta^{t'}} L(\mathcal{T}_j; \theta^{t'}) \left( \frac{\partial^2 L(\mathcal{T}_i; \theta^t)}{\partial(\theta^t)^2} \right) = \mathbf{G} \mathbf{H} = \sum_{i=1}^N \mathbf{G}_{1i} \cdot \mathbf{H}_{i.} = \sum_{i=1}^N \mathbf{G}_{1i} \cdot \frac{\partial \mathbf{g}_i^t}{\partial \theta^t}, \quad (4)$$

where  $\mathbf{G}_{1i}$  indicates the  $i$ th element in the vector  $\mathbf{G}$  and  $\mathbf{H}_{i.}$  indicates the  $i$ th row of  $\mathbf{H}$ .  $\mathbf{G}_{1i}$  is considered as a scalar, which is not related to  $\theta^t$  and thus implies  $\frac{\partial \mathbf{G}_{1i}}{\partial \theta^t} = 0$ . Then Equ. 4 can be rewritten as:

$$\sum_{i=1}^N \mathbf{G}_{1i} \cdot \frac{\partial \mathbf{g}_i^t}{\partial \theta^t} = \sum_{i=1}^N \frac{\partial \mathbf{G}_{1i} \cdot \mathbf{g}_i^t}{\partial \theta^t} = \frac{\partial \left( \sum_{i=1}^N \mathbf{G}_{1i} \cdot \mathbf{g}_i^t \right)}{\partial \theta^t}. \quad (5)$$

It is much faster to evaluate gradients of meta-learning loss through Equ. 5.

There also exists a simpler and faster first-order approximation [2] of Equ. 1 that drops the Hessian (i.e., assuming zero curvature around  $\theta^t$ ), resulting in:  $\nabla_{\theta^{t'}} \frac{1}{|\mathcal{T}_j|} L(\mathcal{T}_j; \theta^{t'}) \left[ \mathbf{I} - \eta \frac{1}{|\mathcal{T}_i|} \left( \frac{\partial^2 L(\mathcal{T}_i; \theta^t)}{\partial(\theta^t)^2} \right) \right] \approx \nabla_{\theta^{t'}} \frac{1}{|\mathcal{T}_j|} L(\mathcal{T}_j; \theta^{t'})$ .

<sup>1</sup>Model parameters in matrix form can be viewed as vectors.

## References

- [1] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, et al. Tensorflow: a system for large-scale machine learning. In *OSDI*, volume 16, pages 265–283, 2016. [1](#)
- [2] C. Finn, P. Abbeel, and S. Levine. Model-agnostic meta-learning for fast adaptation of deep networks. *arXiv preprint arXiv:1703.03400*, 2017. [1](#)
- [3] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer. Automatic differentiation in pytorch. 2017. [1](#)
- [4] Y. Zhang, H. Tang, and K. Jia. Fine-grained visual categorization using meta-learning optimization with sample selection of auxiliary data. *arXiv preprint arXiv:1807.10916*, 2018. [1](#)