

ROLE OF TWO LEARNING RATES IN CONVERGENCE OF MODEL-AGNOSTIC META-LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

Model-agnostic meta-learning (MAML) is known as a powerful meta-learning method. However, MAML is notorious for being hard to train because of the existence of two learning rates. Therefore, in this paper, we derive the conditions that inner learning rate α and meta-learning rate β must satisfy for MAML to converge to minima with some simplifications. We find that the upper bound of β depends on α , in contrast to the case of using the normal gradient descent method. Moreover, we show that the threshold of β increases as α approaches its own upper bound. This result is verified by experiments on various few-shot tasks and architectures; specifically, we perform sinusoid regression and classification of Omniglot and MiniImagenet datasets with a multilayer perceptron and a convolutional neural network. Based on this outcome, we present a guideline for determining the learning rates: first, search for the largest possible α ; next, tune β based on the chosen value of α .

1 INTRODUCTION

A pillar of human intelligence is the ability to learn and adapt to unseen tasks quickly and based on only a limited quantity of data. Although machine learning has achieved remarkable results, many recent models require massive quantities of data and are designed for solving particular tasks. Meta-learning, one of the ways of tackling this problem, tries to develop a model that can adapt to new tasks quickly by learning to learn new concepts from few data points (Schmidhuber, 1987; Thrun & Pratt, 1998).

Among meta-learning algorithms, model-agnostic meta-learning (MAML), a gradient-based meta-learning method proposed by Finn et al. (2017), has recently been extensively studied. For example, MAML is used for continual learning (Finn et al., 2019; Jerfel et al., 2019; Spigler, 2019; Al-Shedivat et al., 2018), reinforcement learning (Finn et al., 2017; Al-Shedivat et al., 2018; Gupta et al., 2018; Deleu & Bengio, 2018; Liu & Theodorou, 2019) and probabilistic inference (Finn et al., 2018; Yoon et al., 2018; Grant et al., 2018). The reason why MAML is widely used is because MAML is simple but efficient and applicable to a wide range of tasks independent of model architecture and the loss function. However, MAML is notorious for being hard to train (Antoniou et al., 2019). One of the reasons why training MAML is hard is the existence of two learning rates in MAML: the inner learning rate α and meta-learning rate β . A learning rate is known to be one of the most important parameters, and tuning this parameter may be challenging even if the simple gradient descent (GD) method is used. Nevertheless, we do not yet know the relationship between these two learning rates and have little guidance on how to tune them. Hence, guidelines for choosing these parameters are urgently needed.

In this paper, we investigate the MAML algorithm and propose a guideline for selecting the learning rates. First, in Section 2 we briefly explain by using the first-order approximation how MAML can be regarded as optimization with the negative gradient penalty. Because the gradient norm is related to the shape of the loss surface, a bias towards a larger gradient norm can make training unstable. Next, based on the approximation explained in Section 2, in Section 3, we derive the necessary conditions that α and β must satisfy for MAML to converge to minima with some simplifications. Furthermore, by removing a constraint, we derive a sufficient condition for the necessary condition for MAML to converge to minima with fewer simplifications. We find that the upper bound β_c of meta-learning rate depends on inner learning rate α . In particular, β_c of $\alpha \approx \alpha_c$ is larger than that

of $\alpha = 0$, where α_c is the upper bound of α . Note that β_c of $\alpha = 0$ corresponds to that of vanilla GD. This is verified by experiments in Section 5. These results imply a guideline for selecting the learning rates: first, search for the largest possible α ; next, tune β .

2 MAML AS OPTIMIZATION WITH NEGATIVE GRADIENT PENALTY

2.1 MAML

The goal of MAML is to find a representation that can rapidly adapt to new tasks with a small quantity of data. In other words, MAML performs optimization for parameters $\theta \in \mathbb{R}^d$ that the optimizer can use to quickly reach the optimal parameter θ_τ^* for task τ with few data points. To this end, MAML takes the following steps to update θ . First, it samples a batch of tasks from task distribution $P(\tau)$ and updates θ for each task τ with stochastic gradient descent (SGD):

$$\theta'_\tau = \theta - \alpha \nabla_{\theta} L_\tau(\theta), \quad (1)$$

where α is a step size referred to as the inner learning rate, and $L_\tau(\theta)$ is the loss of τ . The data used for this update is called training data. Next, MAML resamples data, referred to as test data, from each τ and computes the loss at the updated parameters θ'_τ , obtaining $L_\tau(\theta'_\tau)$ for each task. Finally, to determine θ that can be adapted to θ'_τ for all tasks, θ is updated with the gradient of a sum of loss values $L_\tau(\theta'_\tau)$ over all tasks. In other words,

$$\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_{\tau \sim P(\tau)} L_\tau(\theta'_\tau), \quad (2)$$

where β is the learning rate called the meta-learning rate. Though learning rates α and β can be tuned during training or different for each task in practice, we will think them as fixed scalar hyperparameters.

2.2 NEGATIVE GRADIENT PENALTY

Unless otherwise noted, we will consider the case of only one step being made per update, and the data are not resampled to compute the loss for updating θ . The case of multiple steps and that of training data and test data being separated are considered in Appendix A. The gradient of the loss at θ'_τ is $\mathbf{g}_\tau(\theta'_\tau) = \nabla_{\theta} L_\tau(\theta'_\tau) = \frac{\partial L_\tau}{\partial \theta'_\tau} \nabla_{\theta} \theta'_\tau$, where $\mathbf{g}(\cdot)$ is the gradient of $L(\cdot)$ with respect to θ . If α is small, we can assume that $\frac{\partial L_\tau}{\partial \theta'_\tau} \mathbf{I} = \mathbf{g}_\tau(\theta)$; this seems to hold since α is usually small (Finn et al., 2017). Then,

$$\nabla_{\theta} L_\tau(\theta'_\tau) = \frac{\partial L_\tau}{\partial \theta'_\tau} \nabla_{\theta} \theta'_\tau = \frac{\partial L_\tau}{\partial \theta'_\tau} (\mathbf{I} - \alpha \nabla_{\theta}^2 L_\tau) \quad (3)$$

$$\approx \mathbf{g}_\tau(\theta) - \alpha \mathbf{H}_\tau(\theta) \mathbf{g}_\tau(\theta). \quad (4)$$

The above is known as the first-order approximation, which has been mentioned by Finn et al. (2017) and extensively studied by Nichol et al. (2018) and Fallah et al. (2019). It is known that the error induced by the first-order approximation does not degrade the performance so much in practice (Finn et al., 2017). Also, Fallah et al. (2019) theoretically proved that the first-order approximation of MAML does not affect convergence result when each task is similar to each other or α is small enough.

For simplicity, we will assume that only one task is considered during training, omitting task index τ . Therefore, instead of $\sum_{\tau \sim P(\tau)} L_\tau(\theta'_\tau)$, we will consider $L(\theta')$ as the MAML loss. Since the MAML loss is just a sum of task-specific loss, extension to the case of multiple tasks being considered is straightforward, as provided in Appendix A.1.1. Because $\nabla_{\theta}(\mathbf{g}(\theta)^\top \mathbf{g}(\theta)) = 2\mathbf{H}(\theta)\mathbf{g}(\theta)$, if we define $\tilde{L}(\theta) = L(\theta')$,

$$\tilde{L}(\theta) = L(\theta) - \frac{\alpha}{2} \mathbf{g}(\theta)^\top \mathbf{g}(\theta). \quad (5)$$

The above means that MAML can be regarded as optimization with the negative gradient penalty with some simplifications. We will analyze this MAML loss in Section 3. It can also be interpreted

as a Taylor series expansion of the MAML loss for the first-order term, up to scale:

$$\tilde{L}(\boldsymbol{\theta}) = L(\boldsymbol{\theta} - \alpha \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta})) \quad (6)$$

$$\begin{aligned} &\approx L(\boldsymbol{\theta}) - \alpha \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) \quad (\text{Taylor series expansion}) \\ &= L(\boldsymbol{\theta}) - \alpha \mathbf{g}(\boldsymbol{\theta})^\top \mathbf{g}(\boldsymbol{\theta}). \end{aligned} \quad (7)$$

The fact that MAML is optimization with the negative gradient penalty is worth keeping in mind. Because the goal of gradient-based optimization is to find a point where the gradient is zero, a bias that favors a larger gradient is highly likely to make training unstable; this can be a cause of instability of MAML (Antoniou et al., 2019). In fact, as shown in Fig. 1, the gradient norm becomes larger during training, as do the gradient inner products, as Guiroy et al. (2019) observed.

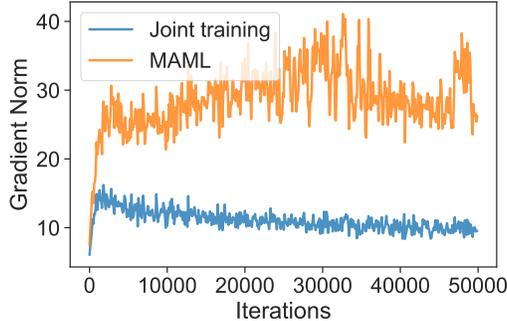


Figure 1. Gradient norm during training. We compute the norm per task and subsequently compute their average. *Joint training* shows when $\alpha = 0$, and *MAML* is when $\alpha = 1e-2$. These results are computed using training data, but those determined using test data behave similarly. The total number of iterations is 50000, $\beta = 1e-3$ and the Adam optimizer is used. Other settings are the same as those in Section 5.2.

3 LEARNING RATE FOR CONVERGENCE

In this section, we will derive the conditions that learning rate α and β should satisfy to converge. To this end, we will assume that only one step is taken for update and training data and test data are not distinguished as we did in Section 2. Also, we do not consider SGD but steepest GD to derive the conditions. In other words, same training data are assumed to be used continuously for updating parameters during training. First, we will consider the case of only single task being considered. Next, we will consider the case of multiple tasks being considered.

3.1 SINGLE TASK

3.1.1 NECESSARY CONDITION FOR INNER LEARNING RATE α

First, we derive the condition that learning rate α should satisfy. To this end, we will consider the condition that a fixed point is a minimum. Taking the Taylor series for the second-order term at a fixed point $\boldsymbol{\theta}^*$, the MAML loss is

$$\tilde{L}(\boldsymbol{\theta}) \approx \tilde{L}(\boldsymbol{\theta}^*) + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \tilde{\mathbf{H}}(\boldsymbol{\theta} - \boldsymbol{\theta}^*). \quad (8)$$

where $\tilde{\mathbf{H}} = \mathbf{H} - \alpha(\mathbf{T}\mathbf{g} + \mathbf{H}^2)$ is the Hessian matrix of $\tilde{L}(\boldsymbol{\theta})$ at $\boldsymbol{\theta}^*$ and $\mathbf{g} = \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}^*) \in \mathbb{R}^d$, $\mathbf{H} = \nabla_{\boldsymbol{\theta}}^2 L(\boldsymbol{\theta}^*) \in \mathbb{R}^{d \times d}$, and $\mathbf{T} = \nabla_{\boldsymbol{\theta}}^3 L(\boldsymbol{\theta}^*) \in \mathbb{R}^{d \times d \times d}$. The calculation of $\tilde{\mathbf{H}}$ is presented in Appendix B. We calculated the magnitudes of $\mathbf{T}\mathbf{g}$ and \mathbf{H}^2 numerically and observed that $\mathbf{T}\mathbf{g}$ was much smaller than \mathbf{H}^2 in practice. Hence, we will ignore $\mathbf{T}\mathbf{g}$ while deriving the conditions and will thus assume that $\tilde{\mathbf{H}} = \mathbf{H} - \alpha\mathbf{H}^2$. Further details are provided in Appendix C, and the case of $\mathbf{T}\mathbf{g}$ being considered is provided in Appendix D. Since $\mathbf{P}\boldsymbol{\Lambda}_{\tilde{\mathbf{H}}}\mathbf{P}^\top = \mathbf{P}[\boldsymbol{\Lambda}_{\mathbf{H}} - \alpha\boldsymbol{\Lambda}_{\mathbf{H}}^2]\mathbf{P}^\top$ where $\boldsymbol{\Lambda}_{\tilde{\mathbf{H}}}$

is a diagonal matrix with entries that are eigenvalues of $\tilde{\mathbf{H}}$ and \mathbf{P} is a matrix with rows that are eigenvectors of $\tilde{\mathbf{H}}$, the necessary condition for $\boldsymbol{\theta}^*$ to reach a minimum is

$$\forall i, \lambda(\tilde{\mathbf{H}})_i = \lambda(\mathbf{H})_i - \alpha\lambda(\mathbf{H})_i^2 \geq 0 \quad (9)$$

$$\Rightarrow \forall i, \alpha \leq \frac{1}{\lambda(\mathbf{H})_i}, (\lambda(\mathbf{H})_i \neq 0) \text{ or } \lambda(\mathbf{H})_i = 0. \quad (10)$$

Note that $\lambda(\mathbf{A})_i$ represents the i th eigenvalue of matrix \mathbf{A} . Hence, if we define $1/0$ to be ∞ , the necessary condition that α must satisfy for $\boldsymbol{\theta}^*$ to reach a minimum is

$$\forall i, \alpha \leq \frac{1}{\lambda(\mathbf{H})_i}. \quad (11)$$

Therefore, α_c is the inverse of the largest eigenvalue of \mathbf{H} .

3.1.2 NECESSARY CONDITION FOR META-LEARNING RATE β

Next, we derive the necessary condition that meta-learning rate β must satisfy for MAML to converge to the minima discussed above. This is an extension of research of LeCun et al. (1998). If we denote $\mathbf{P}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)$ by \mathbf{v} , the MAML loss is $\tilde{L}(\mathbf{v}) \approx \tilde{L}(0) + \frac{1}{2}\mathbf{v}^\top \boldsymbol{\Lambda}_{\tilde{\mathbf{H}}}\mathbf{v}$. Because the gradient of $\tilde{L}(\mathbf{v})$ for \mathbf{v} is $\nabla_{\mathbf{v}}\tilde{L}(\mathbf{v}) = \boldsymbol{\Lambda}_{\tilde{\mathbf{H}}}\mathbf{v}$, the update equation of \mathbf{v} is

$$\mathbf{v}(t+1) = \mathbf{v}(t) - \beta\boldsymbol{\Lambda}_{\tilde{\mathbf{H}}}\mathbf{v}(t) = (\mathbf{I} - \beta\boldsymbol{\Lambda}_{\tilde{\mathbf{H}}})\mathbf{v}(t), \quad (12)$$

where $\mathbf{v}(t)$ is the value of \mathbf{v} during iteration t ($t = 0, \dots, M$), and M is the total number of iterations. Assuming that Eq. 11 holds, the necessary condition that β must satisfy is as follows: for all i ,

$$|1 - \beta\lambda(\mathbf{H} - \alpha\mathbf{H}^2)_i| = |1 - \beta(\lambda(\mathbf{H})_i - \alpha\lambda(\mathbf{H})_i^2)| < 1 \quad (13)$$

$$\Rightarrow -1 + \beta(\lambda(\mathbf{H})_i - \alpha\lambda(\mathbf{H})_i^2) < 1 \quad (\because \lambda(\mathbf{H})_i - \alpha\lambda(\mathbf{H})_i^2 \geq 0 \text{ holds because of Eq. 11}) \quad (14)$$

$$\Rightarrow \beta < \frac{2}{\lambda(\mathbf{H})_i - \alpha\lambda(\mathbf{H})_i^2}. \quad (15)$$

Consequently, the necessary condition for MAML to converge to minima is as follows:

$$\forall i, \alpha \leq \frac{1}{\lambda(\mathbf{H})_i} \wedge \beta \leq \frac{2}{\lambda(\mathbf{H})_i - \alpha\lambda(\mathbf{H})_i^2}. \quad (16)$$

Vanilla GD with learning rate β corresponds to MAML if $\alpha = 0$. In this case, $\beta < \frac{2}{\lambda_{max}}$ is necessary for the optimizer to converge, where λ_{max} is the largest eigenvalue of \mathbf{H} , because $2/\lambda_{max}$ is smaller than any other $2/\lambda_i$ (LeCun et al., 1998). Though this holds for MAML as well, this is not the case if α is close to α_c . The reason is that β_c diverges as α approaches $\frac{1}{\lambda(\mathbf{H})_i}$, or α_c as Eq. 15 indicates. Hence, for MAML we must consider not only the largest but also other eigenvalues and in particular, the second-largest eigenvalue. In short, unlike when vanilla GD is employed, β_c depends on α in the case of MAML, and β_c is expected to be larger if α is close to α_c , as shown in Fig. 2. This finding is validated by experiments presented in Section 5.

3.2 MULTIPLE TASKS

We discussed the case of only one task being available during training in Section 3.1. In this section, we derive upper bounds of α and β that apply if multiple tasks $\tau \sim P(\tau)$ are considered. However, instead of the necessary condition we derive a sufficient condition for the necessary condition of convergence. Assumptions, except that multiple tasks are considered, are the same as those in Section 3.1.

3.2.1 CONDITION FOR α

Now, we define the meta-objective as a sum of task-specific objectives: $\hat{L}(\boldsymbol{\theta}) = \sum_{\tau \sim P(\tau)} \tilde{L}_\tau(\boldsymbol{\theta})$. Then, at a fixed point $\boldsymbol{\theta}^*$,

$$\hat{L}(\boldsymbol{\theta}) \approx \hat{L}(\boldsymbol{\theta}^*) + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \hat{\mathbf{H}}(\boldsymbol{\theta} - \boldsymbol{\theta}^*) \quad (17)$$

$$= \sum_{\tau \sim P(\tau)} \left(L_\tau(\boldsymbol{\theta}^*) + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top (\mathbf{H}_\tau - \alpha \mathbf{H}_\tau^2)(\boldsymbol{\theta} - \boldsymbol{\theta}^*) \right) \quad (18)$$

$$= \sum_{\tau \sim P(\tau)} \left(L_\tau(\boldsymbol{\theta}^*) + \frac{1}{2}(\mathbf{P}_\tau(\boldsymbol{\theta} - \boldsymbol{\theta}^*))^\top \mathbf{P}_\tau (\boldsymbol{\Lambda}_{\mathbf{H}_\tau} - \alpha \boldsymbol{\Lambda}_{\mathbf{H}_\tau^2}) \mathbf{P}_\tau^\top (\mathbf{P}_\tau(\boldsymbol{\theta} - \boldsymbol{\theta}^*)) \right), \quad (19)$$

where $\hat{\mathbf{H}}$ is the Hessian matrix of $\hat{L}(\boldsymbol{\theta})$ at $\boldsymbol{\theta}$. Note that we ignore $\mathbf{T}_\tau \mathbf{g}_\tau$ as we did in Section 3.1.1. Since $\mathbf{H}_\tau - \alpha \mathbf{H}_\tau^2$ is not necessarily simultaneously diagonalizable for each task τ , we cannot exactly express eigenvalues of $\hat{\mathbf{H}}$ as function of $\lambda(\mathbf{H}_\tau)$ and α . Therefore, instead of the exact value of α_c , we will derive an upper bound of α_c .

If $\boldsymbol{\theta}^*$ is a minimum for all $\tilde{L}_\tau(\boldsymbol{\theta}^*)$, $\boldsymbol{\theta}^*$ is a minimum for $\hat{L}(\boldsymbol{\theta})$ as well. Hence, if all eigenvalues of $\tilde{\mathbf{H}}_\tau = \mathbf{H}_\tau - 2\alpha \mathbf{H}_\tau^2$ are nonnegative for all tasks, those of $\hat{\mathbf{H}}$ are nonnegative as well. Therefore, if the inequality

$$\forall \tau, i, \quad \alpha \leq \frac{1}{\lambda(\mathbf{H}_\tau)_i} \quad (20)$$

holds, it guarantees that the necessary condition that $\boldsymbol{\theta}^*$ is a minimum of $\hat{L}(\boldsymbol{\theta})$ is satisfied. Note that this is a condition if $\boldsymbol{\theta}^*$ is a minimum for all $\tilde{L}_\tau(\boldsymbol{\theta}^*)$. Since $\boldsymbol{\theta}^*$ can be a minimum of $\hat{L}(\boldsymbol{\theta})$ even when it is not a minimum for all $\tilde{L}_\tau(\boldsymbol{\theta}^*)$, the upper bound of α seems to be larger than that in Eq. 20 in practice.

3.2.2 CONDITION FOR β

The analysis for β is also basically the same as that performed in Section 3.1.2. Denoting $\mathbf{P}_\tau(\boldsymbol{\theta} - \boldsymbol{\theta}^*)$ by \mathbf{v}_τ , we obtain

$$\sum_{\tau} \mathbf{v}_\tau(t+1) = \sum_{\tau} (\mathbf{I} - \beta [\boldsymbol{\Lambda}_{\mathbf{H}_\tau} - \alpha \boldsymbol{\Lambda}_{\mathbf{H}_\tau^2}]) \mathbf{v}_\tau(t). \quad (21)$$

Since $\tilde{\mathbf{H}}_\tau$ is not always simultaneously diagonalizable for each task, as mentioned in Section 3.2.1, \mathbf{P}_τ differs from task to task in general. Hence, we will consider an upper bound of β_c as we did for α_c in Section 3.2.1. Accordingly, if both Eq. 20 and

$$|1 - \beta(\lambda(\mathbf{H}_\tau)_i - \alpha \lambda(\mathbf{H}_\tau)_i^2)| < 1 \quad (22)$$

hold for any eigenvalue λ_i of any task τ , it guarantees that β satisfies the necessary condition for convergence. Therefore, a sufficient condition for the necessary condition for MAML to converge to minima in the case of multiple tasks is as follows:

$$\forall \tau, i, \quad \alpha \leq \frac{1}{\lambda(\mathbf{H}_\tau)_i} \quad \wedge \quad \beta \leq \frac{2}{\lambda(\mathbf{H}_\tau)_i - \alpha \lambda(\mathbf{H}_\tau)_i^2} \quad (23)$$

4 RELATED WORKS

Several papers have investigated MAML and proposed various algorithms (Nichol et al., 2018; Guioy et al., 2019; Eshratifar et al., 2018; Antoniou et al., 2019; Fallah et al., 2019; Khodak et al., 2019; Vuorio et al., 2018; Finn et al., 2019; Deleu & Bengio, 2018; Liu & Theodorou, 2019; Deleu & Bengio, 2018; Grant et al., 2018). Nichol et al. (2018) studied the first-order MAML family in detail and showed that the MAML gradient could be decomposed into two terms: a term related to joint training and a term responsible for increasing the inner product between gradients for different tasks.

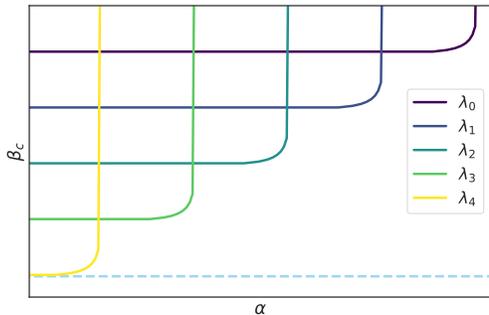


Figure 2. Curves of β_c as a function of α for eigenvalues of the Hessian, $\lambda_0 < \dots < \lambda_4$. Parameter β is supposed to be smaller than β_c for both λ_4 and λ_3 . Hence, β should be chosen from the colored area. Since α must satisfy $\alpha < \frac{1}{\lambda_i}$, α should also be in the colored region. The dashed line shows β_c if $\alpha = 0$. If $\alpha \approx \alpha_c$, β_c is larger than that at $\alpha = 0$.

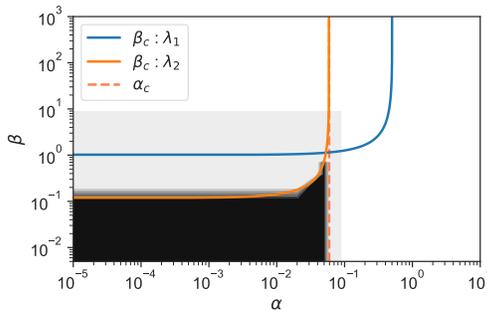


Figure 3. Training loss of linear regression. The area colored in black is when the loss is below $1e-2$, and that in gray is when the loss is over $1e-2$. Uncolored region is not considered. $\beta_c : \lambda_i$ shows β_c of λ_i , where $\lambda_1 < \lambda_2$. The dashed line is α_c . Theoretical β_c and α_c correspond to empirical ones.

Guiroy et al. (2019) investigated the generalization ability of MAML. The researchers observed that generalization was correlated with the average gradient inner product and that flatness of the loss surface, often thought to be an indicator of strong generalizability in normal neural network training, was not necessarily related to generalizability in the case of MAML. Eshratifar et al. (2018) also noted that the average gradient inner product was important. Hence, the authors proposed an algorithm that considered the relative importance of each parameter based on the magnitude of the inner product between the task-specific gradient and the average gradient. Although the above studies were cognizant of the importance of the inner product of the gradients, they did not explicitly insert the negative gradient inner product, which is the negative squared gradient norm with simplifications, as a regularization term. To consider MAML as optimization with a regularization term is a contribution of our study.

Antoniou et al. (2019) enumerated five factors that could cause training MAML to be difficult. Then, they authors proposed an algorithm to address all of these problems and make training MAML easier and more stable. Behl et al. (2019), like us, pointed out that tuning the inner learning rate α and meta-learning rate β was troublesome. The authors approached this problem by proposing an algorithm that tuned learning rates automatically during training.

Fallah et al. (2019) studied the convergence theory of MAML. They proposed a method for selecting meta-learning rate by approximating smoothness of the loss. Based on this result, they proved that MAML can find an ε -first-order stationary point after sufficient number of iterations. On the other hand, we studied the relationship between conditions that inner learning rate α and meta-learning rate β must satisfy and showed that how the largest possible β is affected by the value of α .

5 EXPERIMENTS

In this section, we will present the results of experiments to confirm our expectation that MAML allows larger β if α is close to its upper bound. First, we will show the result of linear regression with simplifications used in Section 3.1. Because linear regression is convex optimization, the result is expected to exactly match the theory presented in Section 3.1. Second, to check if our expectation is confirmed in practice as well, we will present results of the practical case without any simplification. In particular, we conducted sinusoid regression and classification of Omniglot and MiniImagenet datasets with a multilayer perceptron and a convolutional neural network (CNN). Note that the meta-objective used for experiments is not a sum of task-specific objectives but a mean of them.

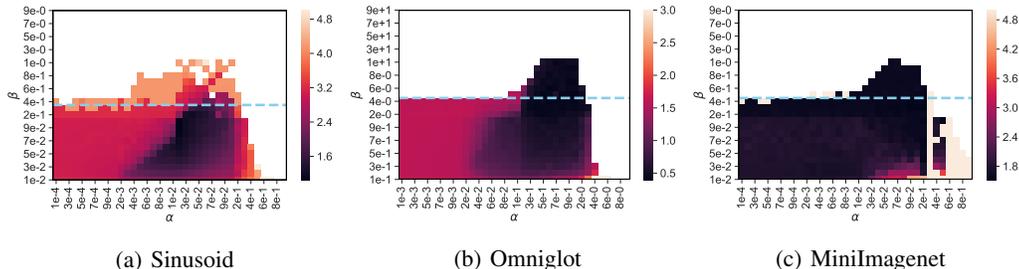


Figure 4. Training losses for (a) sinusoid regression, (b) Omniglot classification, and (c) MiniImagenet classification at various values of α and β after a fixed number of iterations. The area with no color represents the diverged losses, and the dashed line indicates the values of β above which the loss diverges for $\alpha = 0$. The maximum possible β is larger if α is close to the value above which the losses diverge than that at $\alpha = 0$.

5.1 LINEAR REGRESSION

We performed a linear regression, where the task is to regress a linear function with scale parameter in the range of $[0, 5.0]$ and bias parameter in the range of $[0, 5.0]$ based on data points in the range of $[-5.0, 5.0]$. The true function has the same architecture as that of the model. We employed the steepest gradient descent method to minimize the mean squared loss, where 1 step was taken during update. Only one task was considered during training and the same data was used to update the task-specific parameter and the meta parameter as we did in Section 3.1. Using these settings, we computed the training loss after 500 iterations with α in the range of $[1e-5, 9e-2]$ and β in the range of $[5e-3, 9e+0]$. The eigenvalues are those of the Hessian matrix of the training loss at the end of the training, where $\alpha = 5e-2$ and $\beta = 7e-1$. We chose this training loss because it was thought to be the closest to minima. Fig. 3 shows the training losses at various values of α and β . Horizontal axis indicates α and vertical axis indicates β . The curves are β_c of two eigenvalues and the dashed line shows α_c . In the case of linear regression with simplifications used in Section 3.1, the result of numerical experiment shows good agreement with upper bounds of β_c and α_c that we derived in Section 3.1, as shown in Fig. 3.

5.2 SINUSOID REGRESSION

We conducted a sinusoid regression, where each task is to regress a sine wave with amplitude in the range of $[0.1, 5.0]$ and phase in the range of $[0, \pi]$ based on data points in the range of $[-5.0, 5.0]$. A multilayer perceptron with two hidden units of size 40 and ReLU was trained with SGD. The batch size of data was 10, the number of tasks was 100, and 1 step was taken for update. Using these settings, we computed the training loss after 500 iterations with α in the range of $[1e-4, 9e-1]$ and β in the range of $[1e-2, 9e+0]$. Fig. 4 (a) shows the training losses with various values of α and β . The dashed line indicates β of $\alpha = 0$ over which training loss diverges. According to Fig. 4 (a), if α is close to the value above which the losses diverge, a larger β can be used. As explained above, we did not put any simplification as we did in Section 3 and 5.1 in the experiment, meaning that we used different data for updating task-specific parameter and meta parameter, considered multiple tasks, and employed not steepest GD but SGD as optimizer. Despite simplifications, surprisingly, this result confirms the expectation that MAML allows larger β if α is close to α_c .

5.3 CLASSIFICATION

We performed classification of the Omniglot and MiniImagenet datasets (Lake et al., 2011; ravi & Larochelle, 2017), which are benchmark datasets for few-shot learning. The model used was essentially the same as that Finn et al. (2017), and hence, Vinyals et al. (2016) used. The task is a five-way one-shot classification, where the query size is 15, the number of update steps is two, and the task batch size is 32 for Omniglot and four for MiniImagenet. In this setup, we computed the training losses after 100 iterations for the Omniglot dataset and one epoch for the MiniImagenet dataset with various values of α and β ; for Omniglot, α was in the range of $[1e-3, 9e-0]$ and β was in the range of $[1e-1, 9e+1]$, and for MiniImagenet, α was in the range of $[1e-4, 9e-1]$ and β was in

the range of $[1e-2, 9e-0]$. Fig. 4 (b) and (c) show the training losses of classification task at various values of α and β . The dashed line indicates β of $\alpha = 0$ above which training loss diverges. As shown in Fig. 4 (b) and (c), the maximum β is larger at large α . Even though the model architecture is composed of convolutional layer, max-pooling, and batch normalization (Ioffe & Szegedy, 2015) and practical dataset is used for training, our expectation is confirmed in this case as well. This result confirms that our theory is applicable in practice.

Our experimental result confirms that larger α is good for stabilizing MAML training. According to Fig. 4 (a), (b) and (c), moreover, while large β does not necessarily make training loss smaller, employing large α leads to smaller training loss comparatively. These result has a practical implication for tuning the learning rates: first, the largest possible α should be identified, and β may be subsequently tuned based on the value of α . Once you identify α_c , MAML is likely to work well even if meta-learning rate is roughly chosen. Taking large α is desirable for the goal of MAML as well. The aim of MAML is to find a good initial parameter which quickly adapt to new tasks and the quickness is determined by inner learning rate α because α determines the step size from initial parameter to task-specific parameter when model is fine-tuned. Therefore, identifying α_c is not only good for robustifying the model against divergence but also good for finding good initial parameter.

6 CONCLUSIONS

We regard MAML as training with the negative gradient penalty. Based on this formulation, we derived conditions that the inner learning rate α and the meta-learning rate β must satisfy for MAML to converge to minima with some simplifications. We showed that the upper bound of β required for MAML to converge to minima depends on α . Moreover, we found that if α is close to its upper bound α_c , the maximum possible meta-learning rate β_c is larger than that used while training with ordinary SGD. This finding is validated by experiments, confirming that our theory is applicable in practice. According to this result, we propose a guideline for determining α and β ; first, search for α close to α_c ; next, tune β based on the selected value of α .

REFERENCES

- Maruan Al-Shedivat, Trapit Bansal, Yura Burda, Ilya Sutskever, and Igor Mordatch. Continuous Adaptation Via Meta-Learning In Nonstationary And Competitive Environments. *arXiv:1710.03641*, 2018.
- Antreas Antoniou, Harrison Edwards, and Amos Storkey. How to train your maml. *International Conference on Learning Representations (ICLR)*, 2019.
- Harkirat Singh Behl, Atılım Güneş Baydin, and Philip H.S. Torr. Alpha maml: Adaptive model-agnostic meta-learning. *36th International Conference on Machine Learning*, 2019.
- Rajendra Bhatia. Linear algebra to quantum cohomology: The story of alfred horn’s inequalities. *The American Mathematical Monthly*, 108:289–318, 2001.
- Tristan Deleu and Yoshua Bengio. The effects of negative adaptation in Model-Agnostic Meta-Learning. *arXiv preprint arXiv:1812.02159*, 2018.
- Amir Erfan Eshratifar, David Eigen, and Massoud Pedram. Gradient agreement as an optimization objective for meta-learning. *arXiv preprint arXiv:1810.08178*, 2018.
- Alireza Fallah, Aryan Mokhtari, and Asuman Ozdaglar. On the Convergence Theory of Gradient-Based Model-Agnostic Meta-Learning Algorithms. *arXiv preprint arXiv:1908.10400*, 2019.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic metalearning for fast adaptation of deep networks. *International Conference on Machine Learning (ICML)*, 2017.
- Chelsea Finn, Kelvin Xu, and Sergey Levine. Probabilistic Model-Agnostic Meta-Learning. *Neural Information Processing Systems 2018*, 2018.
- Chelsea Finn, Aravind Rajeswaran, Sham Kakade, and Sergey Levine. Online Meta-Learning. *arXiv preprint arXiv:1902.08438*, 2019.

- Erin Grant, Chelsea Finn, Sergey Levine, Trevor Darrell, and Thomas Griffiths. Recasting Gradient-Based Meta-Learning as Hierarchical Bayes. *arXiv preprint arXiv:1801.08930*, 2018.
- Simon Guiroy, Vikas Verma, and Christopher Pal. Towards understanding generalization in gradient-based meta-learning. *arXiv preprint arXiv:1907.07287*, 2019.
- Abhishek Gupta, Benjamin Eysenbach, Chelsea Finn, and Sergey Levine. Unsupervised Meta-Learning for Reinforcement Learning. *arXiv:1806.04640*, 2018.
- Sergey Ioffe and Christian Szegedy. Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. *arXiv:1502.03167*, 2015.
- Ghassen Jerfel, Erin Grant, Thomas L. Griffiths, and Katherine Heller. Reconciling meta-learning and continual learning with online mixtures of tasks. *arXiv:1812.06080*, 2019.
- Mikhail Khodak, Maria-Florina Balcan, and Amee Talwalkar. Provable Guarantees for Gradient-Based Meta-Learning. *arXiv preprint arXiv:1902.10644*, 2019.
- Diederik P. Kingma and Jimmy Lei Ba. Adam: A Method for Stochastic Optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Brenden M. Lake, Ruslan Salakhutdinov, Jason Gross, and Joshua B. Tenenbaum. One shot learning of simple visual concepts. In *Proceedings of the 33th Annual Meeting of the Cognitive Science Society, CogSci 2011, Boston, Massachusetts, USA, July 20-23, 2011*, 2011.
- Yann LeCun, Leon Bottou, B. Genevieve Orr, and Klaus-Robert Müller. Efficient backprop. *Neural networks: Tricks of the trade*, pp. 9 – 50, 1998.
- Guan-Hong Liu and Evangelos A. Theodorou. Deep Learning Theory Review: An Optimal Control And Dynamical Systems Perspective. *arXiv preprint arXiv:1908.10920*, 2019.
- Alex Nichol, Joshua Achiam, and John Schulman. On first-order meta-learning algorithms. *arXiv preprint arXiv:1803.02999*, 2018.
- Sachin ravi and Hugo Larochelle. Optimization as a Model for Few-shot Learning. *The International Conference on Learning Representations 2017*, 2017.
- Jurgen Schmidhuber. Evolutionary principles in self-referential learning, or on learning how to learn: the meta-meta-... hook. *PhD thesis, Technische Universitat Munchen*, 1987.
- Giacomo Spigler. Meta-learned priors slow down catastrophic forgetting in neural networks. *arXiv:1909.04170*, 2019.
- Sebastian Thrun and Lorien Pratt. *Learning to learn*. Springer, 1998.
- Oriol Vinyals, Charles Blundell, Timothy Lillicrap, koray kavukcuoglu, and Daan Wierstra. Matching networks for one shot learning. In *Advances in Neural Information Processing Systems 29*, pp. 3630–3638, 2016.
- Risto Vuorio, Shao-Hua Sun, Hexiang Hu, and Joseph J. Lim. Toward Multimodal Model-Agnostic Meta-Learning. *arXiv preprint arXiv:1812.07172*, 2018.
- Jaesik Yoon, Taesup Kim, Ousmane Dia, Sungwoong Kim, Yoshua Bengio, and Sungjin Ahn. Bayesian Model-Agnostic Meta-Learning. *Neural Information Processing Systems 2018*, 2018.

A NEGATIVE GRADIENT PENALTY: GENERAL CASE

In Section 2, we explained the simplest case of training data and test data being the same, only one task being considered and only one step being taken for update. In this section, we analyze the cases of waiving one of these simplifications. Here, we will interpret MAML loss as a Taylor series expansion of the MAML loss for the first-order term.

A.1 WHEN THE TRAINING DATA AND TEST DATA ARE DIFFERENT

When using different data to compute the losses for updating the task-specific parameters and the meta parameters as done in practical applications, the simplified loss is as follows:

$$\tilde{L}(\boldsymbol{\theta}) = L^{test}(\boldsymbol{\theta} - \alpha \nabla_{\boldsymbol{\theta}} L^{train}(\boldsymbol{\theta})) \quad (24)$$

$$\approx L^{test}(\boldsymbol{\theta}) - \alpha \mathbf{g}(\boldsymbol{\theta})^{test\top} \mathbf{g}(\boldsymbol{\theta})^{train}. \quad (25)$$

where L^{train} and L^{test} represent the loss values for training and test data, respectively, and g^{train} and g^{test} are gradients of these losses.

A.1.1 WHEN MULTIPLE TASKS ARE CONSIDERED

When multiple tasks are considered, as MAML does, the meta-loss $\hat{L}(\boldsymbol{\theta})$ is just a sum of losses $\tilde{L}_{\tau}(\boldsymbol{\theta})$ for all tasks τ :

$$\hat{L}(\boldsymbol{\theta}) = \sum_{\tau \sim P(\tau)} \tilde{L}_{\tau}(\boldsymbol{\theta}) = \sum_{\tau \sim P(\tau)} (L_{\tau}^{test}(\boldsymbol{\theta}) - \alpha \mathbf{g}_{\tau}(\boldsymbol{\theta})^{test\top} \mathbf{g}_{\tau}(\boldsymbol{\theta})^{train}). \quad (26)$$

A.2 WHEN k STEPS ARE TAKEN DURING THE UPDATE

If task-specific parameters are updated with k -step SGD, the loss can be written as follows;

$$\tilde{L}_k(\boldsymbol{\theta}) = L_k(\boldsymbol{\theta}) - \alpha \left(\sum_i^k \mathbf{g}_k^{\top}(\boldsymbol{\theta}) \mathbf{g}_i(\boldsymbol{\theta}) \right). \quad (27)$$

Note that $L_i(\boldsymbol{\theta})$ is the loss computed with the data at the i th step, and $\mathbf{g}_i(\boldsymbol{\theta})$ is the gradient of the loss.

B CALCULATION OF \tilde{H}

Because \tilde{H} is the Hessian matrix of \tilde{L} at $\boldsymbol{\theta}^*$, we derive the Hessian of Eq. 8. Then,

$$\tilde{H} = \nabla_{\boldsymbol{\theta}}^2 \tilde{L}(\boldsymbol{\theta}) \quad (28)$$

$$= \nabla_{\boldsymbol{\theta}}^2 \left(L(\boldsymbol{\theta}) - \frac{\alpha}{2} \mathbf{g}(\boldsymbol{\theta})^{\top} \mathbf{g}(\boldsymbol{\theta}) \right) \quad (29)$$

$$= \mathbf{H}(\boldsymbol{\theta}) - \alpha \nabla_{\boldsymbol{\theta}} (\mathbf{H}(\boldsymbol{\theta}) \mathbf{g}(\boldsymbol{\theta})) \quad (30)$$

$$= \mathbf{H}(\boldsymbol{\theta}) - \alpha (\nabla_{\boldsymbol{\theta}} \mathbf{H}(\boldsymbol{\theta}) \mathbf{g}(\boldsymbol{\theta}) + \mathbf{H}(\boldsymbol{\theta}) \mathbf{H}(\boldsymbol{\theta})) \quad (31)$$

$$= \mathbf{H} - \alpha (\mathbf{Tg} + \mathbf{H}^2). \quad (32)$$

C MAGNITUDE OF Tg AND H^2

As we showed in Section 3, especially large eigenvalues of \tilde{H} are important for the upper bounds of learning rates. Therefore, if $\lambda(\mathbf{Tg} + \mathbf{H}^2)_{max} \approx \lambda(\mathbf{H}^2)_{max}$, we can ignore \mathbf{Tg} when deriving the necessary condition. We calculate the maximum and the second-largest eigenvalues of \mathbf{Tg} , \mathbf{H}^2 and $\mathbf{Tg} + \mathbf{H}^2$ of the trained model. As shown in Fig. 5 (a), $\lambda(\mathbf{Tg} + \mathbf{H}^2)_{max}$ is almost equal to $\lambda(\mathbf{H}^2)_{max}$, and $\lambda(\mathbf{Tg})_{max}$ is by far smaller than them. Therefore, ignoring $\lambda(\mathbf{Tg})_{max}$ is reasonable when the conditions are derived. Furthermore, we calculate the Frobenius norm of \mathbf{Tg} and \mathbf{H}^2 . As Fig. 5 (b) indicates, the Frobenius norm of \mathbf{Tg} is much smaller than that of \mathbf{H}^2 , meaning that

\mathbf{Tg} is negligible in the sense of the magnitude of the norm as well. These results confirm that we can neglect \mathbf{Tg} when considering $\tilde{\mathbf{H}}$. The trained model we use was trained with essentially the same condition that we explain in Section 5.2. However, unlike in Section 5.2, the total number of iterations is 50000. Parameters α and β are 1e-2 and 1e-3 respectively.

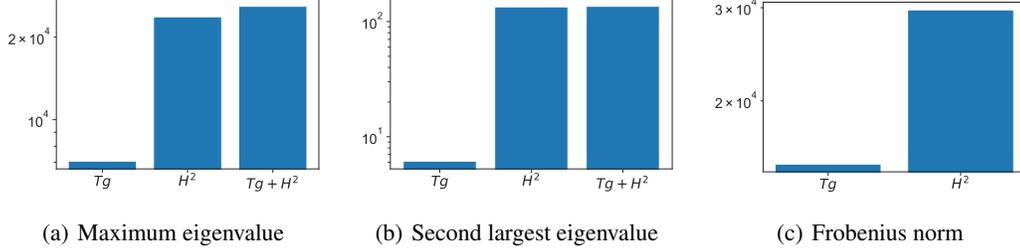


Figure 5. (a): The maximum eigenvalues of \mathbf{Tg} , \mathbf{H}^2 and $\mathbf{Tg} + \mathbf{H}^2$. It is clear that the maximum eigenvalue of $\mathbf{Tg} + \mathbf{H}^2$ is almost the same as that of \mathbf{H}^2 , while that of \mathbf{Tg} is much smaller than them. (b): The second-largest eigenvalues of \mathbf{Tg} , \mathbf{H}^2 and $\mathbf{Tg} + \mathbf{H}^2$. Like (a), the second-largest eigenvalue of $\mathbf{Tg} + \mathbf{H}^2$ is almost equal to that of \mathbf{H}^2 . (c): The Frobenius norm of \mathbf{Tg} and \mathbf{H}^2 . The Frobenius norm of \mathbf{Tg} is much smaller than that of \mathbf{H}^2 .

D CONVERGENCE CONDITION WHEN \mathbf{Tg} IS CONSIDERED

We assumed that \mathbf{Tg} was negligible in Section 3. In this section, we derive a condition that α and β should satisfy under some assumptions when \mathbf{Tg} is considered. However, here instead of the necessary condition we derive a sufficient condition for the necessary condition for MAML to converge to minima.

D.1 CONDITION FOR INNER LEARNING RATE α TO SATISFY

When \mathbf{Tg} is considered, the Hessian matrix of the MAML loss \tilde{L} is $\tilde{\mathbf{H}} = \mathbf{H} - \alpha(\mathbf{Tg} + \mathbf{H}^2)$. Because $\mathbf{H} - \alpha(\mathbf{Tg} + \mathbf{H}^2)$ is a real symmetric matrix, it can be diagonalized. Then, the necessary condition that a fixed point θ^* is a minimum is

$$\forall i, \lambda(\tilde{\mathbf{H}})_i = \lambda(\mathbf{H} - \alpha(\mathbf{Tg} + \mathbf{H}^2))_i \geq 0. \quad (33)$$

Note that \mathbf{Tg} and \mathbf{H} are not simultaneously diagonalizable, so we cannot decompose $\lambda(\tilde{\mathbf{H}})_i$ into eigenvalues of each matrix as we did in Section 3. Therefore, we have to consider the relationship among $\lambda(\tilde{\mathbf{H}})_i$, $\lambda(\mathbf{H})_i$ and $\lambda(\mathbf{Tg})_i$. In general, it is known that $n \times n$ Hermitian matrices \mathbf{A} and \mathbf{B} satisfy the following equation (Bhatia, 2001):

$$\lambda^\downarrow(\mathbf{A}) + \lambda^\uparrow(\mathbf{B}) \prec \lambda(\mathbf{A} + \mathbf{B}), \quad (34)$$

where $\lambda(\mathbf{A})$ represents a vector with elements that are eigenvalues of \mathbf{A} , \uparrow indicates the operation of sorting a vector in the ascending order, and \downarrow indicates that in the descending order. If two real vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ are related in the following way, \mathbf{x} is said to be majorized by \mathbf{y} and the relationship is written as $\mathbf{x} \prec \mathbf{y}$:

$$\sum_{i=1}^k x_i^\downarrow \leq \sum_{i=1}^k y_i^\downarrow \quad \text{for } 1 \leq k \leq d, \quad (35)$$

$$\sum_{i=1}^d x_i^\downarrow = \sum_{i=1}^d y_i^\downarrow. \quad (36)$$

We define $\mathbf{A} = \mathbf{H} - \alpha\mathbf{H}^2$ and $\mathbf{B} = -\alpha\mathbf{Tg}$; then,

$$\lambda^\downarrow(\mathbf{A}) + \lambda^\uparrow(\mathbf{B}) \prec \lambda(\mathbf{A} + \mathbf{B}) \quad (37)$$

$$\Rightarrow \lambda^\downarrow(\mathbf{A}) - \alpha\lambda^\downarrow(\mathbf{Tg}) \prec \lambda(\mathbf{A} + \mathbf{B}) \quad (38)$$

$$\Rightarrow \sum_{i=1}^k (\lambda^\downarrow(\mathbf{A})_i - \alpha\lambda^\downarrow(\mathbf{Tg})_i) \leq \sum_{i=1}^k \lambda^\downarrow(\mathbf{A} + \mathbf{B})_i \quad \text{for } 1 \leq k \leq d. \quad (39)$$

Let us suppose that for top k eigenvalues of \mathbf{A} , \mathbf{Tg} , and $\mathbf{A} + \mathbf{B}$, conditions $\lambda(\mathbf{A} + \mathbf{B})_i \approx c_{\mathbf{A}+\mathbf{B}}$, $\lambda(\mathbf{Tg})_i \approx c_{\mathbf{Tg}}$ and $\lambda(\mathbf{A})_i \approx c_{\mathbf{A}}$ hold, where $c_{\mathbf{A}+\mathbf{B}}$, $c_{\mathbf{A}}$ and $c_{\mathbf{Tg}}$ are some constant values. Then,

$$\lambda(\mathbf{A})_i - \alpha\lambda(\mathbf{Tg})_i \leq \lambda(\mathbf{A} + \mathbf{B})_i \quad \text{for } 1 \leq i \leq k \quad (40)$$

$$\Rightarrow \lambda(\mathbf{H} - \alpha\mathbf{H}^2)_i - \alpha\lambda(\mathbf{Tg})_i \leq \lambda(\mathbf{H} - \alpha(\mathbf{Tg} + \mathbf{H}^2))_i \quad \text{for } 1 \leq i \leq k. \quad (41)$$

The necessary condition for a fixed point θ^* to be a minimum is

$$\forall i, \lambda(\mathbf{H} - \alpha\mathbf{H}^2)_i - \alpha\lambda(\mathbf{Tg})_i \geq 0 \quad (42)$$

Now, we assume that all eigenvalues of \mathbf{H} below a threshold are 0; $\lambda(\mathbf{H})_i \approx 0$ for $k < i \leq n$. In fact, many eigenvalues of the loss computed by a trained model such as a deep neural network are known to be very small. Thus, if the following condition is satisfied, it is enough to say that the necessary condition is satisfied:

$$\begin{cases} \lambda(\mathbf{H} - \alpha\mathbf{H}^2)_i - \alpha\lambda(\mathbf{Tg})_i \geq 0 & \text{for } 1 \leq i \leq k \\ \alpha\lambda(\mathbf{Tg})_i \leq 0 & \text{for } k \leq i \leq d \end{cases} \quad (43)$$

$$\Rightarrow \begin{cases} \lambda(\mathbf{H})_i - \alpha\lambda(\mathbf{H})_{min}^2 - \alpha\lambda(\mathbf{Tg})_i \geq 0 & \text{for } 1 \leq i \leq k \\ \lambda(\mathbf{Tg})_i \leq 0 & \text{for } k \leq i \leq d \end{cases} \quad (44)$$

$$\Leftrightarrow \begin{cases} \lambda(\mathbf{H})_i - \alpha(\lambda(\mathbf{H})_{min}^2 + \lambda(\mathbf{Tg})_i) \geq 0 & \text{for } 1 \leq i \leq k \\ \lambda(\mathbf{Tg})_i \leq 0 & \text{for } k \leq i \leq d \end{cases} \quad (45)$$

$$\Rightarrow \begin{cases} \lambda(\mathbf{H})_i - \alpha\lambda(\mathbf{Tg})_i \geq 0 & \text{for } 1 \leq i \leq k \\ \lambda(\mathbf{Tg})_i \leq 0 & \text{for } k \leq i \leq d \end{cases} \quad (\because \lambda(\mathbf{H})_{min} = 0) \quad (46)$$

In other words, if the smaller $d - k$ eigenvalues of \mathbf{Tg} are not negative values, the fixed point is not guaranteed to be a minimum, and if they are negative, satisfying Eq. 46 is sufficient for the necessary condition.

D.2 CONDITION FOR META LEARNING RATE β TO SATISFY

Next, we derive a condition for meta-learning rate β to satisfy. Our analysis will be based on the assumptions identical to those mentioned above. If the inequalities in Eq. 46 hold, it is sufficient for β to satisfy

$$\forall i, -1 + \beta (\lambda(\mathbf{H} - \alpha\mathbf{H}^2)_i - \alpha\lambda(\mathbf{Tg})_i) \leq -1 + \beta\lambda(\mathbf{H} - \alpha(\mathbf{Tg} + \mathbf{H}^2))_i < 1 \quad (47)$$

$$\Rightarrow \forall i, \beta < \frac{2}{\lambda(\mathbf{H})_i - \alpha\lambda(\mathbf{Tg})_i} \quad (\because \lambda(\mathbf{H})_{min} = 0) \quad (48)$$

for the necessary condition of convergence.

E EXPERIMENT WITH ADAM

We conducted the same experiment as that explained in Section 5.2 and 5.3, but this time, we used Adam optimizer instead of SGD for the meta-optimizer (Kingma & Ba, 2014). The results are different from those we showed in Section 5. For training with Adam, MAML no longer allows us to use larger β_c . Because Adam has more parameters that we have to consider than SGD does, this result is not at all surprising. Nonetheless, it is important to keep these facts in mind when Adam is used for MAML optimization.

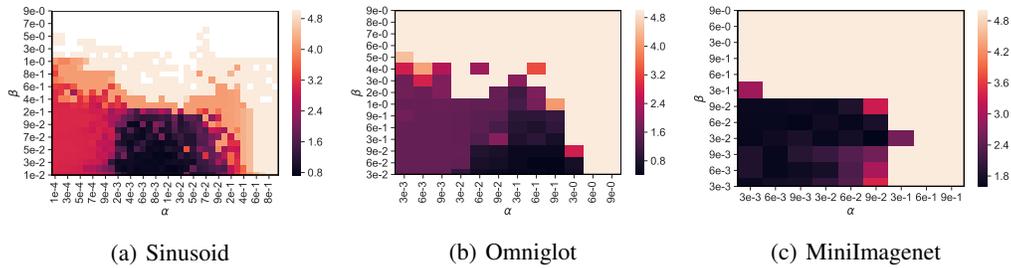


Figure 6. Training losses for (a) the sinusoid regression, (b) Omniglot classification and (c) MiniImagenet classification for various values of the inner learning rate α and meta-learning rate β after a fixed number of iterations. The area with no color represents the diverged losses.