

ON THE CONDITIONS OF MAML CONVERGENCE

SHIRO TAKAGI¹ YOSHIHIRO NAGANO¹ YUKI YOSHIDA¹ MASATO OKADA¹

¹THE UNIVERSITY OF TOKYO

We derived the necessary conditions of inner learning rate α and meta-learning rate β for a simplified MAML to locally converge to local min from any point in the vicinity of the local min

We found that maximum possible β is larger when α is close to its maximum possible value

MAML AS NEGATIVE GRADIENT PENALTY

MAML

Update parameters to find a representation that can rapidly adapt to new tasks with a small quantity of data.

1. Inner loop

2. Meta-loop

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

Inner lr Meta-lr [Finn et al. 2017]

Assumptions

1. Only one step is taken per update
2. Only one task is considered
3. Data are not resampled

Approximation

Ignore higher order derivative terms for simplicity.

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

$$\approx g_{\tau}(\theta) - \alpha H_{\tau}(\theta) g_{\tau}(\theta)$$

Negative gradient penalty

An approximated MAML loss can be regarded as the loss with the negative gradient penalty.

$$\tilde{L}(\theta) = L(\theta') \approx L(\theta) - \frac{\alpha}{2} g(\theta)^{\top} g(\theta)$$

Then, update equation of \mathbf{v} is

$$\mathbf{v}(t+1) = \mathbf{v}(t) - \beta \Lambda_{\tilde{H}} \mathbf{v}(t) \quad [\text{LeCun et al. 1998}]$$

To reach a minimum, β should satisfy the following condition:

$$\forall i, |1 - \beta \lambda(H - \alpha H^2)_i| < 1$$

The necessary condition of β is:

$$\forall i, \quad \beta < \frac{2}{\lambda(H)_i - \alpha \lambda(H)_i^2}$$

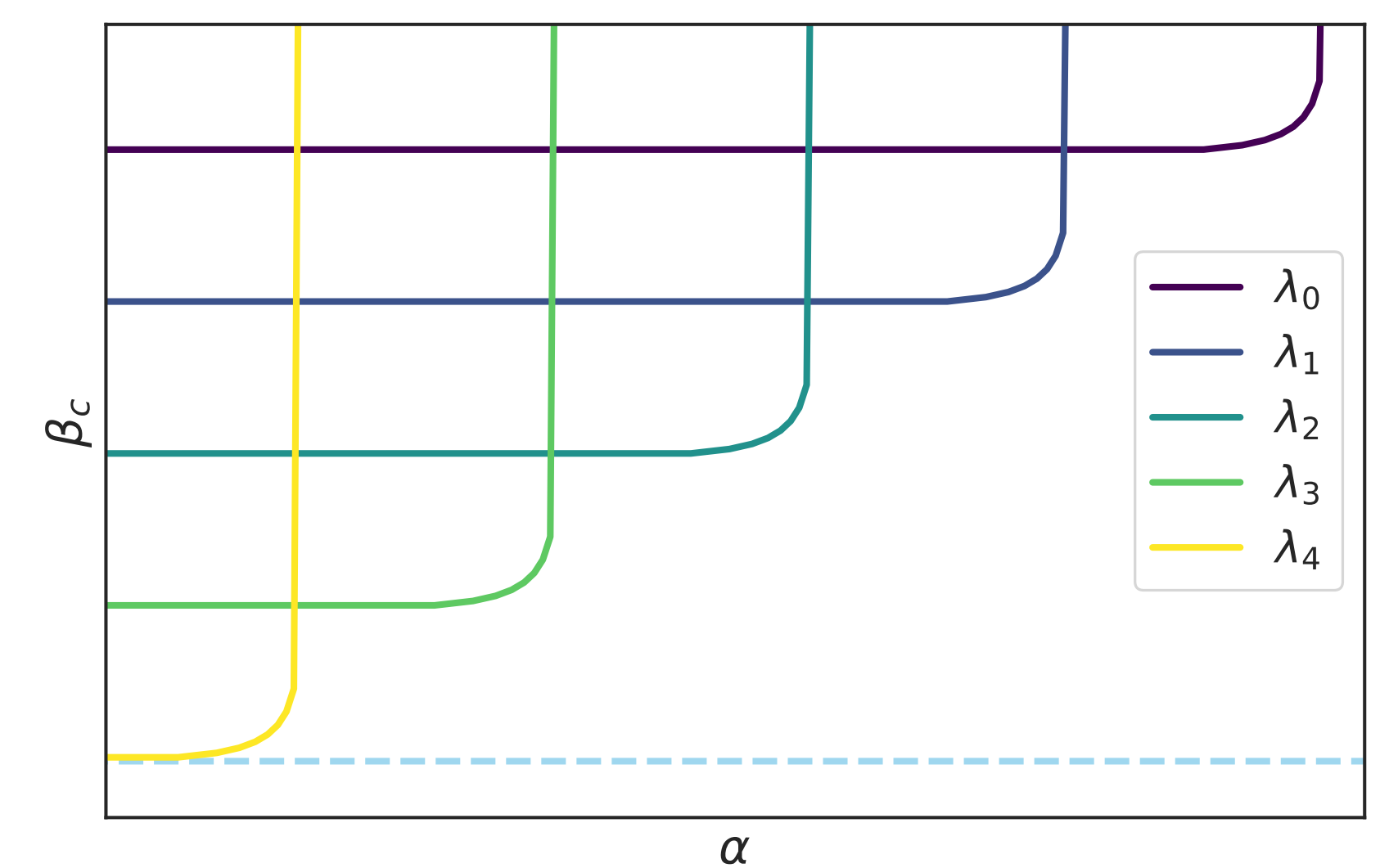
Condition for simplified MAML to locally converge

The condition for the simplified MAML to locally converge to local min from any point in the vicinity of the local min:

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

As α approaches α_c , upper bound β_c diverges.

$\Rightarrow \beta_c$ is larger when α is close to α_c



CONVERGENCE CONDITION

Simplified MAML Loss

Taking the Taylor series for the second-order term at a fixed point θ^* , the simplified MAML loss is

$$\tilde{L}(\theta) \approx \tilde{L}(\theta^*) + \frac{1}{2} (\theta - \theta^*)^{\top} \tilde{H} (\theta - \theta^*)$$

Condition for inner learning rate α

Since the necessary condition θ^* to be a local minimum is that all eigenvalues λ of the Hessian \tilde{H} at θ^* are non-negative, α should satisfy the following condition:

$$\forall i, \lambda(\tilde{H})_i \approx \lambda(H)_i - \alpha \lambda(H)_i^2 \geq 0$$

$$\tilde{H} = H - \alpha (Tg + H^2) \approx H - \alpha (H^2)$$

The necessary condition of α is:

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

Condition for meta-learning rate β

\tilde{H} is diagonalizable: $\tilde{H} = P \Lambda_{\tilde{H}} P^{\top}$

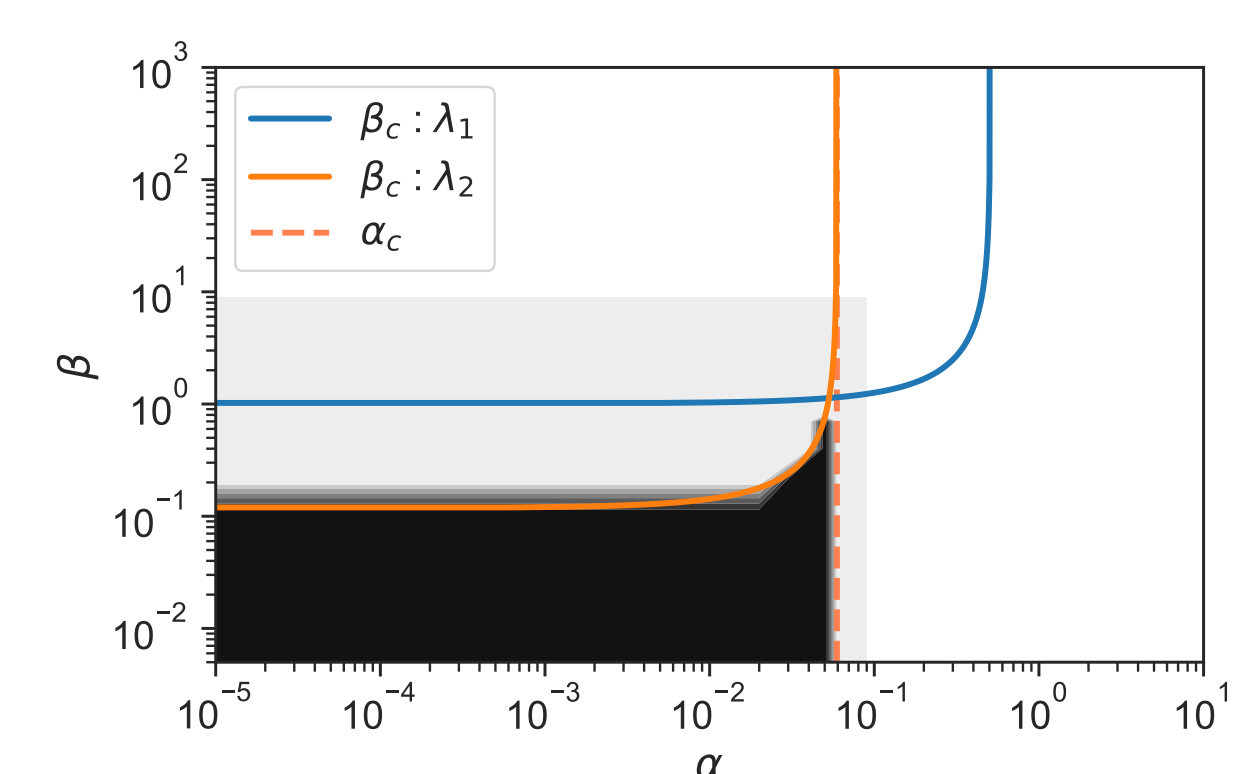
Denoting $P^{\top}(\theta - \theta^*)$ by \mathbf{v} , the simplified MAML loss is

$$\tilde{L}(\mathbf{v}) \approx \tilde{L}(0) + \frac{1}{2} \mathbf{v}^{\top} \Lambda_{\tilde{H}} \mathbf{v}$$

EXPERIMENTS

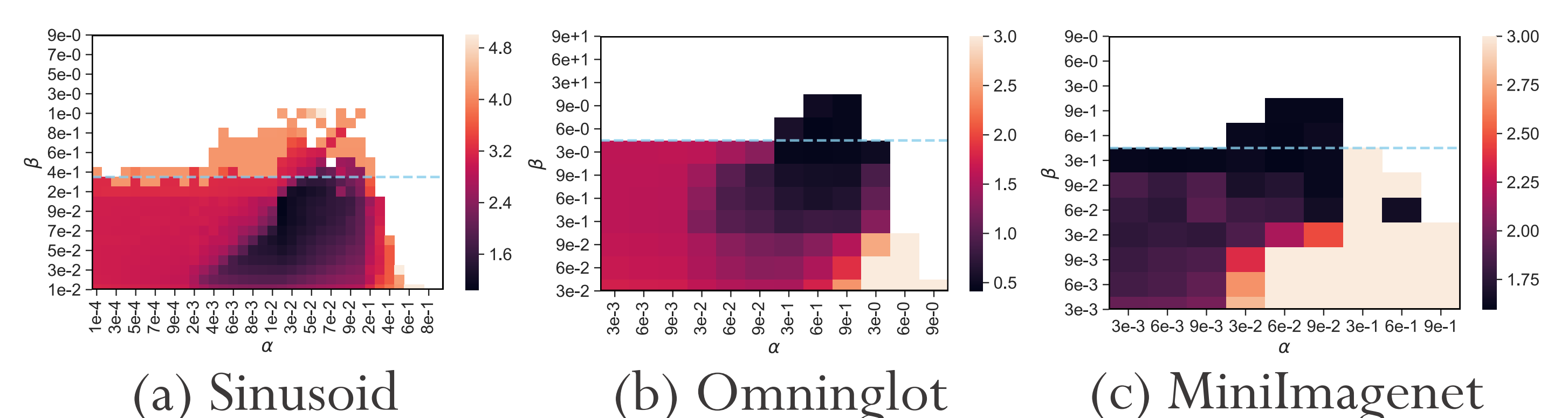
Linear regression

Linear regression with all assumptions being hold. Theoretical β_c and α_c match empirical ones.



Few-shot learning

Training error of (a) Sinusoid regression and (b) Omniglot and (c) MiniImagenet classification with various α and β .



The largest possible β is larger when α is close to its maximum possible value.

\Rightarrow Our theory explains the experimental result.