

Model Agnostic Meta Learning

paper: [Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks ICML 2017](#)

MAML is a simple but effective algorithm aims at learning a good initial parameters of so that the model will have maximal performance on a new task after been updated through one or more gradient steps. This algorithm is model-agnostic and can be applied in different types of tasks including regression, classification and reinforcement learning. This algorithm is discussed in few-shot learning scenario in this paper.

1. General form of MAML algorithm

Algorithm 1 Model-Agnostic Meta-Learning

Require: $p(\mathcal{T})$: distribution over tasks

Require: α, β : step size hyperparameters

1: randomly initialize θ

2: **while** not done **do**

3: Sample batch of tasks $\mathcal{T}_i \sim p(\mathcal{T})$

4: **for all** \mathcal{T}_i **do**

5: Evaluate $\nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$ with respect to K examples

6: Compute adapted parameters with gradient descent: $\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$

7: **end for**

8: Update $\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$

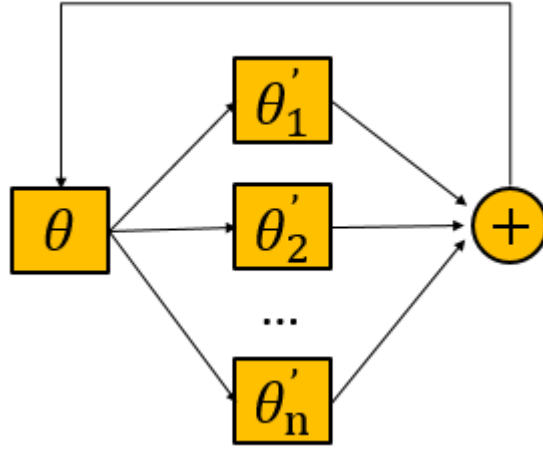
9: **end while**

In step 8, the objective function is:

$$\min_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i}) = \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})})$$

The MAML meta-gradient update involves a gradient through a gradient (In other words, a second-order derivative). The goal of meta-gradient update (Step 8) is to minimize the total loss across all tasks after adapting the initial model to these specific tasks. During every iteration, MAML algorithm forces the model to considering all the tasks it trained with.

In particular, MAML uses the **summation** of all the task-specific gradients to update the initial parameters. The summation loss $\sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$ can be seen as a composite function of parameter θ .



The parameter θ can be seen as a global or general representation for the distribution over tasks, while the parameter θ'_i can be seen as a local or specific representation after adaptation.

This can be viewed from another perspective: MAML use a composite gradient from the individual gradients of specific tasks to update a global state to find a better “base model”.

2. Species of MAML

For supervised classification and regression:

Algorithm 2 MAML for Few-Shot Supervised Learning

Require: $p(\mathcal{T})$: distribution over tasks

Require: α, β : step size hyperparameters

- 1: randomly initialize θ
 - 2: **while** not done **do**
 - 3: Sample batch of tasks $\mathcal{T}_i \sim p(\mathcal{T})$
 - 4: **for all** \mathcal{T}_i **do**
 - 5: Sample K datapoints $\mathcal{D} = \{\mathbf{x}^{(j)}, \mathbf{y}^{(j)}\}$ from \mathcal{T}_i
 - 6: Evaluate $\nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$ using \mathcal{D} and $\mathcal{L}_{\mathcal{T}_i}$ in Equation (2) or (3)
 - 7: Compute adapted parameters with gradient descent:
 $\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$
 - 8: Sample datapoints $\mathcal{D}'_i = \{\mathbf{x}^{(j)}, \mathbf{y}^{(j)}\}$ from \mathcal{T}_i for the meta-update
 - 9: **end for**
 - 10: Update $\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$ using each \mathcal{D}'_i and $\mathcal{L}_{\mathcal{T}_i}$ in Equation 2 or 3
 - 11: **end while**
-

For reinforcement learning:

Algorithm 3 MAML for Reinforcement Learning

Require: $p(\mathcal{T})$: distribution over tasks**Require:** α, β : step size hyperparameters

```
1: randomly initialize  $\theta$ 
2: while not done do
3:   Sample batch of tasks  $\mathcal{T}_i \sim p(\mathcal{T})$ 
4:   for all  $\mathcal{T}_i$  do
5:     Sample  $K$  trajectories  $\mathcal{D} = \{(\mathbf{x}_1, \mathbf{a}_1, \dots, \mathbf{x}_H)\}$  using  $f_\theta$ 
      in  $\mathcal{T}_i$ 
6:     Evaluate  $\nabla_\theta \mathcal{L}_{\mathcal{T}_i}(f_\theta)$  using  $\mathcal{D}$  and  $\mathcal{L}_{\mathcal{T}_i}$  in Equation 4
7:     Compute adapted parameters with gradient descent:
       $\theta'_i = \theta - \alpha \nabla_\theta \mathcal{L}_{\mathcal{T}_i}(f_\theta)$ 
8:     Sample trajectories  $\mathcal{D}'_i = \{(\mathbf{x}_1, \mathbf{a}_1, \dots, \mathbf{x}_H)\}$  using  $f_{\theta'_i}$ 
      in  $\mathcal{T}_i$ 
9:   end for
10:  Update  $\theta \leftarrow \theta - \beta \nabla_\theta \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$  using each  $\mathcal{D}'_i$ 
      and  $\mathcal{L}_{\mathcal{T}_i}$  in Equation 4
11: end while
```

The only difference between two instantiations is Step5 and Step8 where the samples are from data points and trajectories respectively.

3. Experiments

Only parts of the experiments are shown here, more details can be found in the paper.

1) The learning curves in few-shot regression at meta test-time



MAML converges very quickly after one step and to a lower error.

2) Few-shot image classification on MiniImagenet

MiniImagenet (Ravi & Larochelle, 2017)	5-way Accuracy	
	1-shot	5-shot
fine-tuning baseline	28.86 \pm 0.54%	49.79 \pm 0.79%
nearest neighbor baseline	41.08 \pm 0.70%	51.04 \pm 0.65%
matching nets (Vinyals et al., 2016)	43.56 \pm 0.84%	55.31 \pm 0.73%
meta-learner LSTM (Ravi & Larochelle, 2017)	43.44 \pm 0.77%	60.60 \pm 0.71%
MAML, first order approx. (ours)	48.07 \pm 1.75%	63.15 \pm 0.91%
MAML (ours)	48.70 \pm 1.84%	63.11 \pm 0.92%

Note the first order approximation of MAML. Since the second-order derivative is computational expensive, the authors omitted the second-order computation and found this operation only slightly affected the performance.

Apart from higher accuracy, the model with MAML uses fewer parameters compared to matching networks and the meta-learner LSTM, since the algorithm does not introduce any additional parameters beyond the weights of the classifier itself.