

# Mode-Agnostic Meta-Learning for Fast Adaptation of Deep Networks

Scaramuzzino Giovanna

September 6, 2021

# Overview

- 1 Introduction
  - What is Meta Learning?
  - Formalizing Meta-Learning
  - Few Shot Learning
- 2 Model-Agnostic Meta-Learning
  - Algorithm
  - Gradient of Gradient
  - First order approximation
- 3 Experiments and results
  - Regression problem
  - Classification problem
- 4 Conclusion

# What is Meta Learning?

Deep Learning technics have seen great success in a variety of fields but there are limitations:

- quantities of data (large training set to train our model)
- compute resources

**Meta Learning** is an alternative paradigm where machine learning models gains experience over multiple learning process and uses this experience to improve its future learning performance.

Meta learning produces a versatile AI model that can learn to perform various tasks without having to train them from scratch.

**Goal:** learner quickly learn a new task from a small amount of new data

Two stages to update the model:

- **base learning**: inner learning algorithm solve a task
- **meta learning**: an outer algorithm updates the inner learning algorithm such that the model it learns improves an outer objective.

# Conventional Machine Learning

In supervised learning, we are given a training dataset  $\mathcal{D} = \{(x_1, y_1), \dots, (x_n, y_n)\}$ .

We can train a predictive model  $\hat{y} = f_{\theta}(x)$  by solving:

$$\theta^* = \arg \min_{\theta} \mathcal{L}(\mathcal{D}; \theta, \omega) \quad (1)$$

Where:

- $\mathcal{L}$ : loss function
- $\theta$ : parameter of model
- $\omega$ : the optimizer chosen

# Meta Learning: Task-Distribution View

Define:

- $\mathcal{T} = \{\mathcal{D}, \mathcal{L}\}$  task, with  $\mathcal{D}$ :dataset,  $\mathcal{L}$ :loss function
- $p(\mathcal{T})$ : distribution of tasks
- $\omega$ : meta-knowledge
- $M$  source tasks used in meta-training stage as
$$\mathcal{D}_{source} = \{(\mathcal{D}_{source}^{train}, \mathcal{D}_{source}^{val})^{(i)}\}_{i=1}^M$$
- $Q$  target tasks used in meta-testing stage as
$$\mathcal{D}_{target} = \{(\mathcal{D}_{target}^{train}, \mathcal{D}_{target}^{test})^{(i)}\}_{i=1}^Q$$

We have to solve:

$$\min_{\omega} \mathbb{E}_{\mathcal{T} \sim p(\mathcal{T})} \mathcal{L}(D; \omega) \quad (2)$$

Assume access to a set of source tasks sampled from a distribution of tasks  $p(\mathcal{T})$ , **Meta-training** step can be written as:

$$\omega^* = \arg \max_{\omega} \log p(\omega | \mathcal{D}_{source}) \quad (3)$$

# Meta Learning: Task-Distribution and Bilevel Optimization Views

How to solve the meta-training step?

Using **bilevel optimization**, can be formalised as follows:

$$\omega^* = \arg \min_{\omega} \sum_{i=1}^M \mathcal{L}^{meta}(\theta^{*(i)}(\omega), \omega, \mathcal{D}_{source}^{val(i)}) \quad (4)$$

$$s.t. \quad \theta^{*(i)}(\omega) = \arg \min_{\theta} \mathcal{L}^{task}(\theta, \omega, \mathcal{D}_{source}^{train(i)}) \quad (5)$$

where  $\mathcal{L}^{meta}$  and  $\mathcal{L}^{task}$  refer to the outer and inner objectives

**Meta-testing** use the meta-knowledge to train the base model on each previously unseen target task i:

$$\theta^{*(i)} = \arg \max_{\theta} \log p(\theta | \omega^*, \mathcal{D}_{target}^{train(i)}) \quad (6)$$

$\mathcal{D}_{target}^{test(i)}$  can be used to evaluate the accuracy of the meta-learner.

# Few Shot Learning

Few-Shot learning is learning from fewer data points

N-way K-shot Learning:

- N is the number of the classes
- K is the number of data points in each of the classes in the dataset

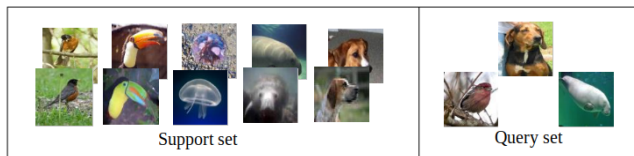
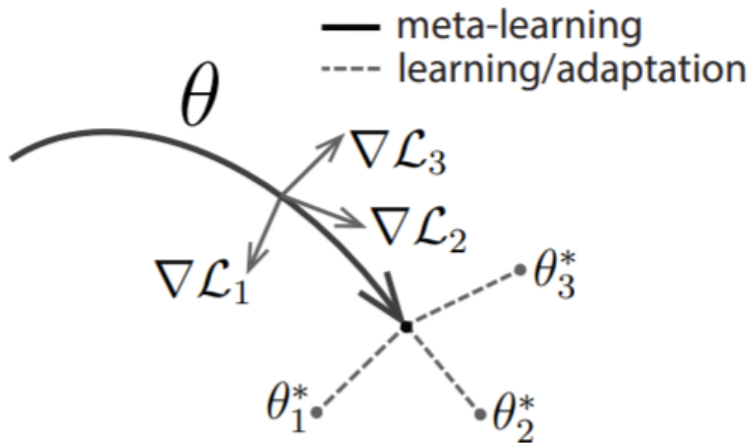


Figure: Example of 5way-2shot - Minilmagenet Dataset



# Model-Agnostic Meta-Learning

- The authors proposed an algorithm **compatible with any model** trained with gradient descent and applicable to a variety of different learning problems: classification, regression, reinforcement learning.
- **Key idea:** train the model's initial parameters such that a **small number of gradient updates** will lead to fast learning on a new task using a **small amount of data**.



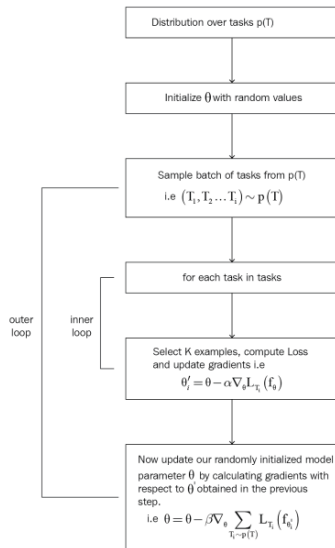
# Algorithm

## Algorithm 2 MAML for Few-Shot Supervised Learning

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

**Require:**  $\alpha, \beta$ : 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  $K$  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**



Mean-Squared-Error for regression:

$$\mathcal{L}_{\mathcal{T}_i}(f_{\phi}) = \sum_{\mathbf{x}^{(j)}, \mathbf{y}^{(j)} \sim \mathcal{T}_i} \|f_{\phi}(\mathbf{x}^{(j)}) - \mathbf{y}^{(j)}\|_2^2 \quad (2)$$

Cross-Entropy for classification:

$$\mathcal{L}_{\mathcal{T}_i}(f_{\phi}) = \sum_{\mathbf{x}^{(j)}, \mathbf{y}^{(j)} \sim \mathcal{T}_i} \mathbf{y}^{(j)} \log f_{\phi}(\mathbf{x}^{(j)}) + (1 - \mathbf{y}^{(j)}) \log(1 - f_{\phi}(\mathbf{x}^{(j)})) \quad (3)$$

# Gradient of Gradient

From line 10 in Algorithm 2:

$$\theta = \theta - \beta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i}) \quad \text{Recall: } \theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$$

$$= \theta - \beta \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i}) \quad (\mathcal{L} \text{ is differentiable})$$

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

$$= \theta - \beta \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \boxed{(1 - \alpha \nabla_{\theta}^2 \mathcal{L}_{\mathcal{T}_i}(f_{\theta}))} \nabla_{\theta'_i} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$$



Calculation of Hessian matrix is required.

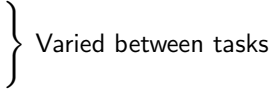
It Makes MAML be slow when backpropagating.

# First order approximation

The Authors to reduce the cost try to compute the 1 order approximation  
Ignoring the second-order term that appears in the update of MAML

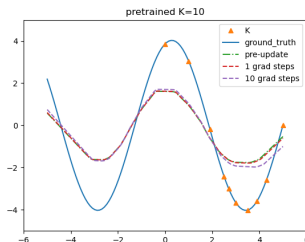
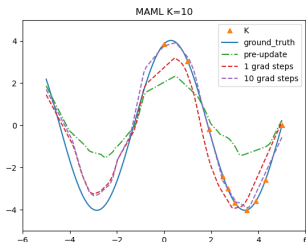
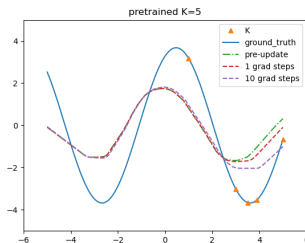
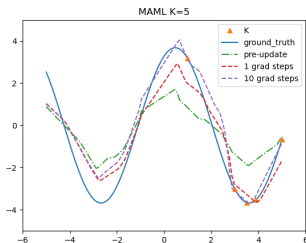
$$\begin{aligned}\theta &= \theta - \beta \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \boxed{(1 - \alpha \nabla_{\theta}^2 \mathcal{L}_{\mathcal{T}_i}(f_{\theta}))} \nabla_{\theta'_i} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})) \\ &= \theta - \beta \nabla_{\theta'_i} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})\end{aligned}$$

# Regression problem

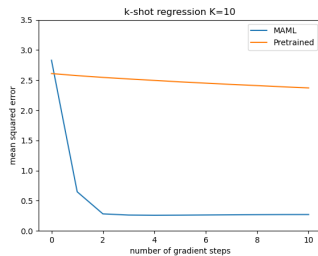
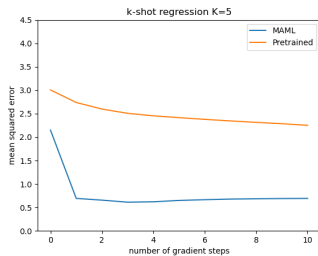
- Sine wave function with:
  - Amplitude in  $[0.1, 5.0]$
  - Phase in  $[0, \pi]$
  - Datapoints sampled uniformly from  $[-5.0, 5.0]$
- Loss function: Mean Squared Error
- Neural Network: 2 hidden layers with 40 units and ReLU
- Training:
  - 1 gradient step in the inner loop
  - Examples K: 5 or 10
  - Optimizer: Adam with  $\alpha = 0.01$
- Evaluation (finetuning the model learned by MAML):
  - 10 gradient steps in the inner loop
  - K: 5 and 10 datapoints

# Regression problem: results

- To evaluate performance MAML model was compared with model pretraining on all of the tasks



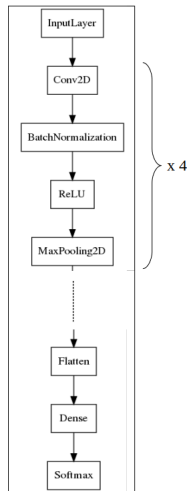
# Quantitative sinusoid regression results





# Classification problem

- Datasets:
  - Omniglot:
    - 50 different alphabets of 1623 characters
    - each character have 20 instances drawn by 20 different people
    - 1200 characters for training, 423 for test
    - images downsampled to 28x28 and augmented with rotations by multiples of 90 degrees
  - Minilmagenet:
    - 64 training classes, 12 validation classes, 24 test classes
- Loss: cross-entropy



# Classification experiments

- N-way K-shot classification
- Outer Learning rate (lr) ( $\beta$ ): 0.001
- models trained for 60000 iterations

<b>Minilmagenet</b>		<b>Train</b>		<b>Test</b>	
	meta-batch-size	gradient steps	inner lr	gradient steps	inner lr
<b>5-way 1-shot</b>	4 tasks	5	0.01	10	0.01
<b>5-way 5-shot</b>	2 tasks	5	0.01	10	0.01

<b>Omniglot</b>		<b>Train</b>		<b>Test</b>	
	meta-batch-size	gradient steps	inner lr	gradient steps	inner lr
<b>5-way 1-shot</b> <b>5-way 5-shot</b>	32 tasks	1	0.4	3	0.4

# Classification results

Minilmagenet	5-way Accuracy	
	1-shot	5-shot
author's model	48.70±1.84%	63.11±0.92%
my result	45.55±0.4%	63.04±0.48%
author's model first order approx.	48.07 ±1.75%	63.15±0.91%
my result first order approx.	43.5±0.4%	59.4±0.49%

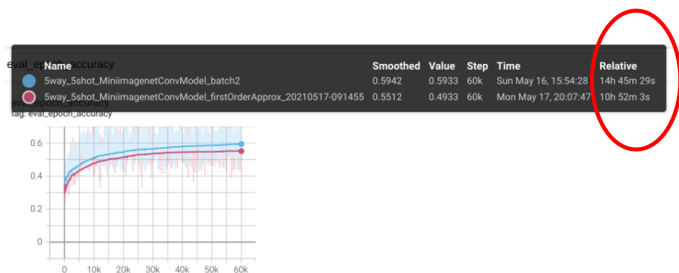
Omniglot	5-way Accuracy	
	1-shot	5-shot
author's model	98.7±0.4%	99.9 ±0.1%
my result	97.41±0.1%	99.09±0.3%

## Metrics

- Accuracy =  $\frac{\text{number of correct prediction}}{\text{total number of prediction}}$
- $\pm$  shows 95% confidence intervals over tasks:

$$\text{average} \pm 1.96 * \frac{\text{std}}{\sqrt{\text{number of observations}}}$$

# Classification results - approximation



**Figure:** Accuracy on validation set: pink line represent model with 1 order approximation

Approximation speed up network computation by 33% while maintaining approximately the same performance

# Conclusion

## Goal

Implement Model-Agnostic Meta learning algorithm

## Final considerations

- My models reach comparables results with reference to the proposed one
- differences may be due to:
  - different choice of some parameters not indicated in the paper
  - different split of the dataset
  - use of early stopping
  - the authors could train the model several times and then choose the best one