#### Meta-Learning with Implicit Gradients

Aravind Rajeswaran, Chelsea Finn, Sham Kakade, Sergey Levine (2019)

**Mayank Mittal** 

 Learn general purpose internal representation that is transferable across different tasks

$$\mathcal{T}_{i} \sim P(\mathcal{T}), \forall i = \{1, ..., M\} \implies \mathcal{D}_{i}^{\mathrm{tr}} = \{(x_{i}^{k}, y_{i}^{k})\}_{k=1}^{K}$$
drawn from 
$$\mathcal{T}_{i} \operatorname{task}$$

distribution over tasks

#### **Meta-Learning objective:**

$$\overbrace{\boldsymbol{\theta}_{\mathrm{ML}}^{*} := \operatorname*{argmin}_{\boldsymbol{\theta} \in \Theta} F(\boldsymbol{\theta})}^{\mathrm{outer-level}}, \text{ where } F(\boldsymbol{\theta}) = \frac{1}{M} \sum_{i=1}^{M} \mathcal{L}\left(\overbrace{\mathcal{A}lg(\boldsymbol{\theta}, \mathcal{D}_{i}^{\mathrm{tr}})}^{\mathrm{inner-level}}, \mathcal{D}_{i}^{\mathrm{test}}\right)$$
$$\underset{\substack{\mathrm{loss}\\\mathrm{function}}}{\overset{\mathrm{task-specific}}{\mathrm{parameters}}}$$

Images from: Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks, Finn et al. (2017)



$$\boldsymbol{\theta}_{\mathrm{ML}}^{\mathrm{outer-level}} = \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{free}} F(\boldsymbol{\theta}), \text{ where } F(\boldsymbol{\theta}) = \frac{1}{M} \sum_{i=1}^{M} \mathcal{L}\left(\overbrace{\mathcal{A}lg(\boldsymbol{\theta}, \mathcal{D}_{i}^{\mathrm{tr}})}^{\mathrm{inner-level}}, \mathcal{D}_{i}^{\mathrm{test}}\right)$$

$$\boldsymbol{\theta}^{\mathrm{k+1}} = \boldsymbol{\theta}^{k} - \eta \frac{1}{M} \sum_{i=1}^{M} \frac{\partial}{\partial \boldsymbol{\theta}} \mathcal{L}(Alg(\boldsymbol{\theta}^{k}, \mathcal{D}_{i}^{\mathrm{tr}}), \mathcal{D}_{i}^{\mathrm{test}})$$

$$\boldsymbol{\varphi}_{i} \equiv \mathcal{A}lg(\boldsymbol{\theta}, \mathcal{D}_{i}^{\mathrm{tr}}) = \boldsymbol{\theta} - \alpha \nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, \mathcal{D}_{i}^{\mathrm{tr}})$$

$$\hat{\mathcal{L}}_{i}(\boldsymbol{\theta}) := \mathcal{L}(\boldsymbol{\phi}, \mathcal{D}_{i}^{\mathrm{test}})$$

$$\boldsymbol{\theta}^{k+1} = \boldsymbol{\theta}^{k} - \eta \frac{1}{M} \sum_{i=1}^{M} d_{\boldsymbol{\theta}} \mathcal{L}_{i}(\mathcal{A}lg_{i}(\boldsymbol{\theta}^{k}))$$

$$d_{\theta}\mathcal{L}_{i}(\mathcal{A}lg_{i}(\theta)) = \frac{d\mathcal{A}lg_{i}(\theta)}{d\theta} \nabla_{\phi}\mathcal{L}_{i}(\phi) \mid_{\phi = \mathcal{A}lg_{i}(\theta)} = \frac{d\mathcal{A}lg_{i}(\theta)}{d\theta} \nabla_{\phi}\mathcal{L}_{i}(\mathcal{A}lg_{i}(\theta))$$

$$(d \times d) \text{-size Jacobian matrix}$$

 $\Rightarrow$ 



Gradients flow over iterative updates of inner-level!

- Need to store optimization path
- Computationally expensive

Hence, heuristics such as first-order MAML and Reptile are used



#### **Motivation**



# Implicit MAML (iMAML)

- In MAML,  $\phi_i \equiv Alg(\theta, D_i^{tr})$  is an iterative gradient descent
  - dependence of model-parameters {φ<sub>i</sub>} on meta-parameters θ shrinks and vanishes as number of gradient steps increases
- Instead, iMAML uses an optimum solution of the regularized loss

$$\mathcal{A}lg^{\star}(\boldsymbol{\theta}, \mathcal{D}_{i}^{\mathrm{tr}}) = \operatorname*{argmin}_{\boldsymbol{\phi}' \in \Phi} \mathcal{L}(\boldsymbol{\phi}', \mathcal{D}_{i}^{\mathrm{tr}}) + \frac{\lambda}{2} ||\boldsymbol{\phi}' - \boldsymbol{\theta}||^{2}$$

encourages dependency between model-parameters and meta-parameters

#### Implicit MAML (iMAML)

$$\boldsymbol{\theta}_{\mathrm{ML}}^{*} := \operatorname*{argmin}_{\boldsymbol{\theta} \in \Theta} F(\boldsymbol{\theta}) \,, \text{ where } F(\boldsymbol{\theta}) = \frac{1}{M} \sum_{i=1}^{M} \mathcal{L}_{i} \big( \mathcal{A} l g_{i}^{*}(\boldsymbol{\theta}) \big), \text{ and} \\ \mathcal{A} l g_{i}^{*}(\boldsymbol{\theta}) := \operatorname*{argmin}_{\boldsymbol{\phi}' \in \Phi} G_{i}(\boldsymbol{\phi}', \boldsymbol{\theta}), \text{ where } G_{i}(\boldsymbol{\phi}', \boldsymbol{\theta}) = \hat{\mathcal{L}}_{i}(\boldsymbol{\phi}') + \frac{\lambda}{2} ||\boldsymbol{\phi}' - \boldsymbol{\theta}||^{2}$$

**Issue:** Still the meta-gradient requires to computation of the Jacobian, i.e.

$$\boldsymbol{d}_{\boldsymbol{\theta}}\mathcal{L}_{i}(\mathcal{A}lg_{i}(\boldsymbol{\theta})) = \frac{d\mathcal{A}lg_{i}(\boldsymbol{\theta})}{d\boldsymbol{\theta}} \nabla_{\boldsymbol{\phi}}\mathcal{L}_{i}(\boldsymbol{\phi}) \mid_{\boldsymbol{\phi}=\mathcal{A}lg_{i}(\boldsymbol{\theta})} = \boxed{\frac{d\mathcal{A}lg_{i}(\boldsymbol{\theta})}{d\boldsymbol{\theta}}} \nabla_{\boldsymbol{\phi}}\mathcal{L}_{i}(\mathcal{A}lg_{i}(\boldsymbol{\theta}))$$
$$(d \times d) \text{-size Jacobian matrix}$$

Soln: Use implicit Jacobian! 🙂

$$\frac{d\mathcal{A}lg_i^{\star}(\boldsymbol{\theta})}{d\boldsymbol{\theta}} = \left(\boldsymbol{I} + \frac{1}{\lambda} \nabla_{\boldsymbol{\phi}}^2 \hat{\mathcal{L}}_i(\boldsymbol{\phi}_i)\right)^{-1}$$

#### **Implicit Jacobian Lemma**

**Lemma 1.** (Implicit Jacobian) Consider  $Alg_i^{\star}(\theta)$  as defined in Eq. 4 for task  $\mathcal{T}_i$ . Let  $\phi_i = Alg_i^{\star}(\theta)$  be the result of  $Alg_i^{\star}(\theta)$ . If  $\left(\mathbf{I} + \frac{1}{\lambda}\nabla_{\phi}^2 \hat{\mathcal{L}}_i(\phi_i)\right)$  is invertible, then the derivative Jacobian is

$$\frac{d\mathcal{A}lg_i^{\star}(\boldsymbol{\theta})}{d\boldsymbol{\theta}} = \left(\boldsymbol{I} + \frac{1}{\lambda} \nabla_{\boldsymbol{\phi}}^2 \hat{\mathcal{L}}_i(\boldsymbol{\phi}_i)\right)^{-1}.$$
(6)

where  $\mathcal{A}lg_i^{\star}(\boldsymbol{\theta}) := \operatorname*{argmin}_{\boldsymbol{\phi}' \in \Phi} G_i(\boldsymbol{\phi}', \boldsymbol{\theta}), \text{ where } G_i(\boldsymbol{\phi}', \boldsymbol{\theta}) = \hat{\mathcal{L}}_i(\boldsymbol{\phi}') + \frac{\lambda}{2} ||\boldsymbol{\phi}' - \boldsymbol{\theta}||^2$ 

**Proof.** Apply stationary point condition:

$$\nabla_{\phi'} G(\phi', \theta) \mid_{\phi'=\phi} = 0 \implies \nabla \hat{\mathcal{L}}(\phi) + \lambda(\phi - \theta) = 0 \implies \phi = \theta - \frac{1}{\lambda} \nabla \hat{\mathcal{L}}(\phi),$$

which is an implicit equation that often arises in proximal point methods. When the derivative exists, we can differentiate the above equation to obtain:

$$\frac{d\phi}{d\theta} = \boldsymbol{I} - \frac{1}{\lambda} \nabla^2 \hat{\mathcal{L}}(\phi) \frac{d\phi}{d\theta} \implies \left( \boldsymbol{I} + \frac{1}{\lambda} \nabla^2 \hat{\mathcal{L}}(\phi) \right) \frac{d\phi}{d\theta} = \boldsymbol{I}.$$

### **Practical Algorithm**

**Meta gradient:**  $d_{\theta}\mathcal{L}_{i}(\mathcal{A}lg_{i}(\theta)) = \frac{d\mathcal{A}lg_{i}(\theta)}{d\theta} \nabla_{\phi}\mathcal{L}_{i}(\phi) \mid_{\phi = \mathcal{A}lg_{i}(\theta)} = \frac{\frac{d\mathcal{A}lg_{i}(\theta)}{d\theta}}{\int_{\phi}} \nabla_{\phi}\mathcal{L}_{i}(\mathcal{A}lg_{i}(\theta))$  $\frac{d\mathcal{A}lg_{i}^{\star}(\theta)}{d\theta} = \left(\mathbf{I} + \frac{1}{\lambda}\nabla_{\phi}^{2}\hat{\mathcal{L}}_{i}(\phi_{i})\right)^{-1}$ 

#### **Issues:**

- obtaining exact solution to inner level optimization problem may be not feasible, in practice
- computing Jacobian by explicitly forming and inverting the matrix may be intractable for large networks

**Soln:** Use approximations!

#### **Practical Algorithm**

**Meta gradient:** 
$$d_{\theta}\mathcal{L}_{i}(\mathcal{A}lg_{i}(\theta)) = \frac{d\mathcal{A}lg_{i}(\theta)}{d\theta} \nabla_{\phi}\mathcal{L}_{i}(\phi) \mid_{\phi = \mathcal{A}lg_{i}(\theta)} = \underbrace{\frac{d\mathcal{A}lg_{i}(\theta)}{d\theta}}_{g_{i}} \nabla_{\phi}\mathcal{L}_{i}(\mathcal{A}lg_{i}(\theta))$$
  
 $g_{i}$   
 $\frac{d\mathcal{A}lg_{i}^{\star}(\theta)}{d\theta} = \left(\mathbf{I} + \frac{1}{\lambda}\nabla_{\phi}^{2}\hat{\mathcal{L}}_{i}(\phi_{i})\right)^{-1}$ 

#### Soln: Use approximations! 🙂

1. Approximate algorithm solution:

Let  $Alg_i(\boldsymbol{\theta})$  be a  $\delta$ -accurate approximation of  $Alg_i^{\star}(\boldsymbol{\theta})$ , i.e.  $\|Alg_i(\boldsymbol{\theta}) - Alg_i^{\star}(\boldsymbol{\theta})\| \leq \delta$ 

2. Approximate matrix inversion product:

Let  $\boldsymbol{g}_i$  be a vector such that  $\|\boldsymbol{g}_i - \left(\boldsymbol{I} + \frac{1}{\lambda} \nabla_{\phi}^2 \hat{\mathcal{L}}_i(\phi_i)\right)^{-1} \nabla_{\phi} \mathcal{L}_i(\phi_i)\| \leq \delta'$ where  $\phi_i = \mathcal{A}lg_i(\boldsymbol{\theta})$  and  $\mathcal{A}lg_i$  is based on definition 1. Solved using Conjugate Gradient (CG) method  $\boldsymbol{\Box}_w \boldsymbol{w}^{\top} \left(\boldsymbol{I} + \frac{1}{\lambda} \nabla_{\phi}^2 \hat{\mathcal{L}}_i(\phi_i)\right) \boldsymbol{w} - \boldsymbol{w}^{\top} \nabla_{\phi} \mathcal{L}_i(\phi_i)$ 

### **Practical Algorithm**

Implicit gradient accuracy: (valid under other regularity assumptions)

Assuming that  $\delta < \mu/(2\rho)$ , we have that:

$$\|\boldsymbol{g}_i - \boldsymbol{d}_{\boldsymbol{\theta}} \mathcal{L}_i(\mathcal{A} l g_i^{\star}(\boldsymbol{\theta}))\| \leq \left(2\frac{\lambda\rho}{\mu^2}B + \frac{\lambda L}{\mu}\right)\delta + \delta'$$

#### Soln: Use approximations! 🙂

1. Approximate algorithm solution:

Let  $Alg_i(\boldsymbol{\theta})$  be a  $\delta$ -accurate approximation of  $Alg_i^{\star}(\boldsymbol{\theta})$ , i.e.  $\|Alg_i(\boldsymbol{\theta}) - Alg_i^{\star}(\boldsymbol{\theta})\| \leq \delta$ 

2. Approximate matrix inversion product:

Let  $g_i$  be a vector such that  $||g_i - (I + \frac{1}{\lambda} \nabla_{\phi}^2 \hat{\mathcal{L}}_i(\phi_i))^{-1} \nabla_{\phi} \mathcal{L}_i(\phi_i)|| \le \delta'$ where  $\phi_i = \mathcal{A}lg_i(\theta)$  and  $\mathcal{A}lg_i$  is based on definition 1. Solved using Conjugate Gradient (CG) method  $\mathbf{I} = \mathbf{W} \mathbf{W}^{\mathsf{T}} \Big( \mathbf{I} + \frac{1}{\lambda} \nabla_{\phi}^2 \hat{\mathcal{L}}_i(\phi_i) \Big) \mathbf{W} - \mathbf{W}^{\mathsf{T}} \nabla_{\phi} \mathcal{L}_i(\phi_i)$ 

# **Implicit MAML Algorithm**

Algorithm 1 Implicit Model-Agnostic Meta-Learning (iMAML)

- 1: **Require:** Distribution over tasks  $P(\mathcal{T})$ , outer step size  $\eta$ , regularization strength  $\lambda$ ,
- 2: while not converged do
- 3: Sample mini-batch of tasks  $\{\mathcal{T}_i\}_{i=1}^B \sim P(\mathcal{T})$
- 4: **for** Each task  $\mathcal{T}_i$  **do**
- 5: Compute task meta-gradient  $g_i = \text{Implicit-Meta-Gradient}(\mathcal{T}_i, \theta, \lambda)$
- 6: end for 7: Average above gradients to get  $\hat{\nabla}F(\theta) = (1/B)\sum_{i=1}^{B} g_i$ 8: Update meta-parameters with gradient descent:  $\theta \leftarrow \theta - \eta \hat{\nabla}F(\theta)$  // (or Adam) 9: end while 0 uter Level/ Meta-Optimiza tion Step

Algorithm 2 Implicit Meta-Gradient Computation

- 1: Input: Task  $\mathcal{T}_i$ , meta-parameters  $\boldsymbol{\theta}$ , regularization strength  $\lambda$
- 2: Hyperparameters: Optimization accuracy thresholds  $\delta$  and  $\delta'$
- 3: Obtain task parameters  $\phi_i$  using iterative optimization solver such that:  $\|\phi_i Alg_i^{\star}(\theta)\| \leq \delta$
- 4: Compute partial outer-level gradient  $v_i = \nabla_{\phi} \mathcal{L}_{\mathcal{T}}(\phi_i)$
- 5: Use an iterative solver (e.g. CG) along with reverse mode differentiation (to compute Hessian vector products) to compute  $\boldsymbol{g}_i$  such that:  $\|\boldsymbol{g}_i (\boldsymbol{I} + \frac{1}{\lambda}\nabla^2 \hat{\mathcal{L}}_i(\boldsymbol{\phi}_i))^{-1} \boldsymbol{v}_i\| \leq \delta'$
- Inner Level/ Adaptation Step

6: Return:  $g_i$ 

#### **Computation Bounds**

**Theory:** iMAML computes meta-gradient by small error in memory efficient way

Algorithm	Compute	Memory	Error
MAML (GD + full back-prop)	$\kappa \log \left( \frac{D}{\delta} \right)$	$\operatorname{Mem}(\nabla \hat{\mathcal{L}}_i) \cdot \kappa  \log\left(\frac{D}{\delta}\right)$	0
MAML (Nesterov's AGD + full back-prop)	$\sqrt{\kappa}\log\left(\frac{D}{\delta}\right)$	$Mem(\nabla \hat{\mathcal{L}}_i) \cdot \sqrt{\kappa}  \log\left(\frac{D}{\delta}\right)$	0
Truncated back-prop [53] (GD)	$\kappa \log \left( \frac{D}{\delta} \right)$	$\operatorname{Mem}(\nabla \hat{\mathcal{L}}_i) \cdot \kappa  \log\left(\frac{1}{\epsilon}\right)$	$\epsilon$
Implicit MAML (this work)	$\sqrt{\kappa}\log\left(\frac{D}{\delta}\right)$	$\operatorname{Mem}( abla \hat{\mathcal{L}}_i)$	δ

#### **Experiment:**



Figure 2: Accuracy, Computation, and Memory tradeoffs of iMAML, MAML, and FOMAML. (a) Metagradient accuracy level in synthetic example. Computed gradients are compared against the exact meta-gradient per Def 3. (b) Computation and memory trade-offs with 4 layer CNN on 20-way-5-shot Omniglot task. We implemented iMAML in PyTorch, and for an apples-to-apples comparison, we use a PyTorch implementation of MAML from: https://github.com/dragen1860/MAML-Pytorch

#### **Results on few-shots learning**

Table 2: Omniglot results. MAML results are taken from the original work of Finn et al. [15], and first-order MAML and Reptile results are from Nichol et al. [43]. iMAML with gradient descent (GD) uses 16 and 25 steps for 5-way and 20-way tasks respectively. iMAML with Hessian-free uses 5 CG steps to compute the search direction and performs line-search to pick step size. Both versions of iMAML use  $\lambda = 2.0$  for regularization, and 5 CG steps to compute the task meta-gradient.

Algorithm	5-way 1-shot	5-way 5-shot	20-way 1-shot	20-way 5-shot
MAML [15]	$98.7\pm0.4\%$	99.9 $\pm$ 0.1 %	$95.8\pm0.3\%$	$98.9\pm0.2\%$
first-order MAML [15]	$98.3\pm0.5\%$	$99.2\pm0.2\%$	$89.4\pm0.5\%$	$97.9\pm0.1\%$
Reptile [43]	$97.68 \pm 0.04\%$	$99.48 \pm 0.06\%$	$89.43 \pm 0.14\%$	$97.12 \pm 0.32\%$
iMAML, GD (ours)	$99.16 \pm 0.35\%$	$99.67 \pm 0.12\%$	$94.46 \pm 0.42\%$	$98.69\pm0.1\%$
iMAML, Hessian-Free (ours)	$99.50\pm0.26\%$	$99.74 \pm 0.11\%$	$\textbf{96.18} \pm \textbf{0.36\%}$	$\textbf{99.14} \pm \textbf{0.1\%}$

Table 3: Mini-ImageNet 5-way-1-shot accuracy

Algorithm	5-way 1-shot
MAML	$48.70 \pm 1.84 \%$
Reptile	$\begin{array}{c} 48.07 \pm 1.75 \ \% \\ 49.97 \pm 0.32 \ \% \end{array}$
iMAML GD (ours) iMAML HF (ours)	$\begin{array}{c} 48.96 \pm 1.84 \ \% \\ 49.30 \pm 1.88 \ \% \end{array}$

#### **Future Work**

- Broader classes for inner loop optimization
  - dynamic programming
  - energy-based models
  - actor-critic RL
- More flexible regularizers
  - Learn a vector- or matrix-valued λ (regularization coefficient)

#### References

- <u>Model-Agnostic Meta-Learning for Fast</u> <u>Adaptation of Deep Networks</u>, Finn et al. (2017)
- <u>Meta-Learning with Implicit Gradients</u>, Rajeswaran et al. (2019)
- <u>Notes on iMAML</u>, Ferenc Huszar (2019)

# Thank you!

Н





 $\phi_2$