

META-LEARNING WITH IMPLICIT GRADIENTS



Sowmiya Murugiah
sowmiyam@buffalo.edu
50485124



AGENDA

Introduction

Advanced ML Techniques

Meta-Learning

MAML vs iMAML

Formulation

Key Algorithm

Graph

Summary

Introduction

Supervised ML models are data-hungry. If you give large datasets to train large models, you learn a brilliant decision boundary. Traditional ML algorithms start to fall when situations tend to change fast and your model runs out of the training data. The motivation of meta-learning comes from being able to learn from small data. Meta parameters are learned in the outer loop while task-specific models are learned in the inner-loop. By drawing upon implicit differentiation, we develop the implicit MAML algorithm which depends only on the solution to the inner level optimization and not the path taken by the inner loop optimizer



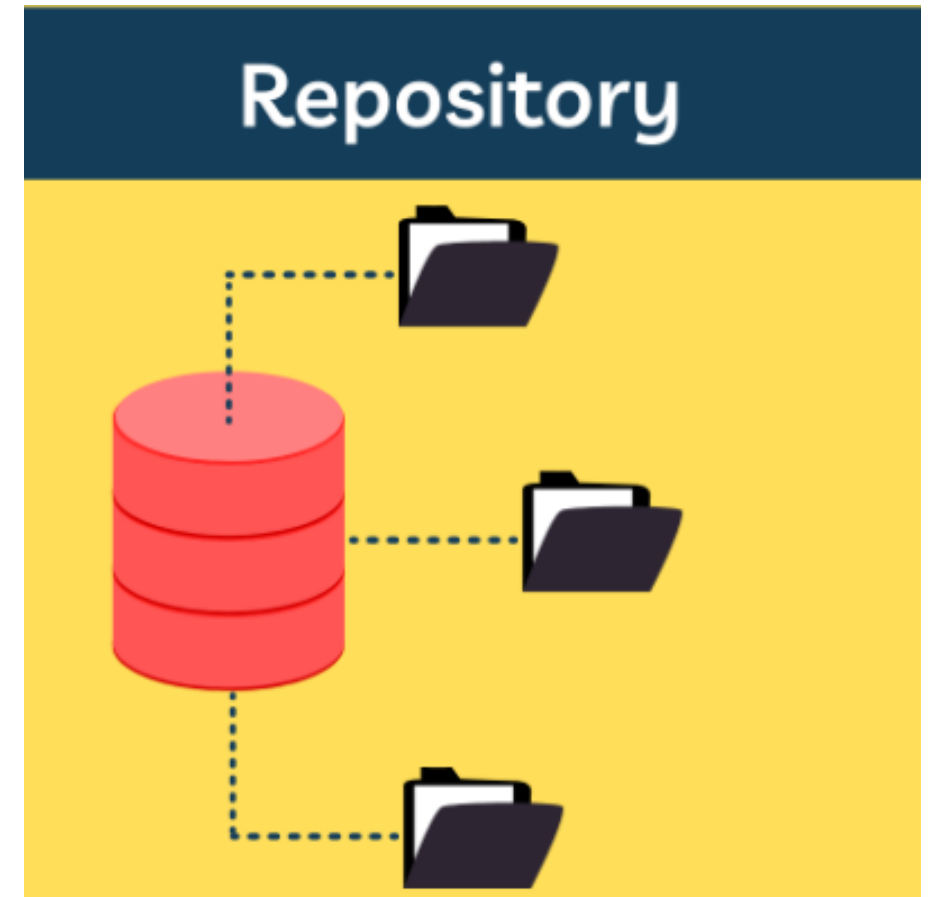
Advanced Machine Learning Techniques

- Transfer Learning – TL
- Multi-Transfer Learning – MTL
- Meta Learning

Transfer Learning

Transfer learning is a technique in machine learning and artificial intelligence where knowledge gained from one task or domain is applied to a different but related task or domain. In transfer learning, a pre-trained model is used as a starting point for a new learning task, rather than starting the learning process from scratch. The pre-trained model has learned features or representations that can be transferred to the new task, reducing the amount of training data required and improving the performance of the new model.

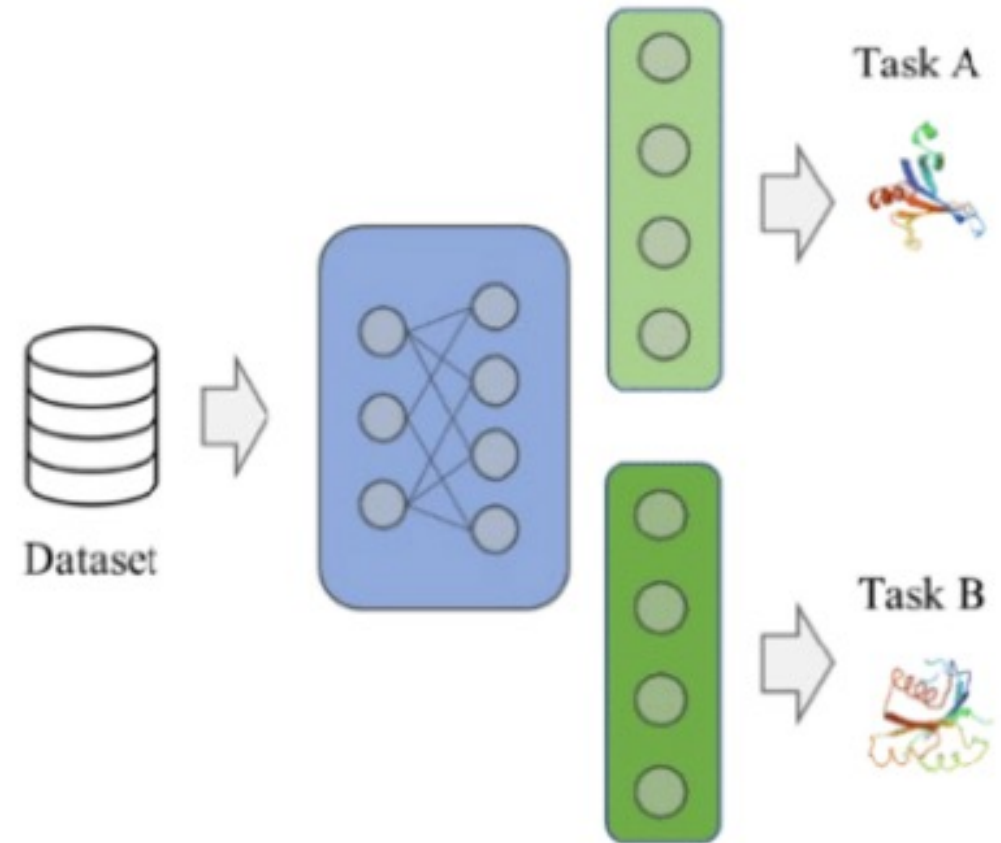
Disadvantage : It can be computationally expensive and time-consuming to train large models



Multi-Transfer Learning

Multi-Transfer Learning is a type of transfer learning that involves transferring knowledge from multiple source domains to a target domain. In traditional transfer learning, a pre-trained model is used to extract features from a single source domain, which are then used to train a new model for a target domain. However, in multi-transfer learning, the pre-trained model is trained on multiple source domains, each with different but related features.

Issue: limit top out accuracy because it is fully not dedicated to one task



WHAT IS META- LEARNING

+

•

○

+

○

•

Meta-Learning

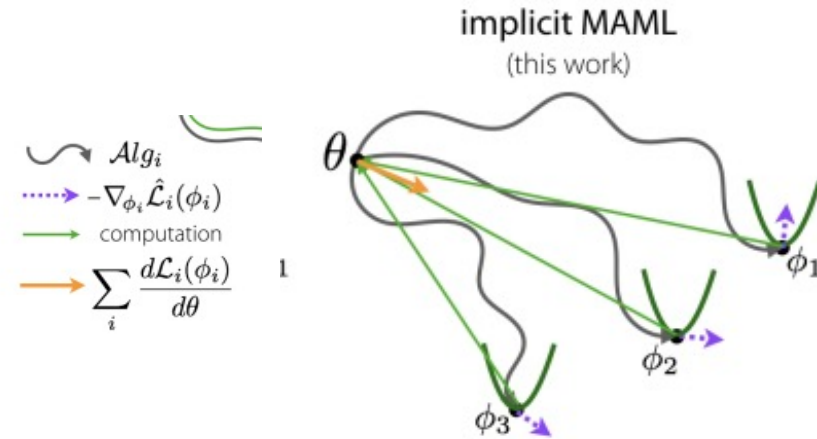
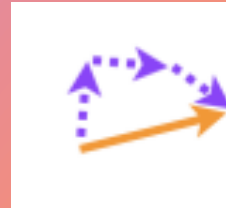
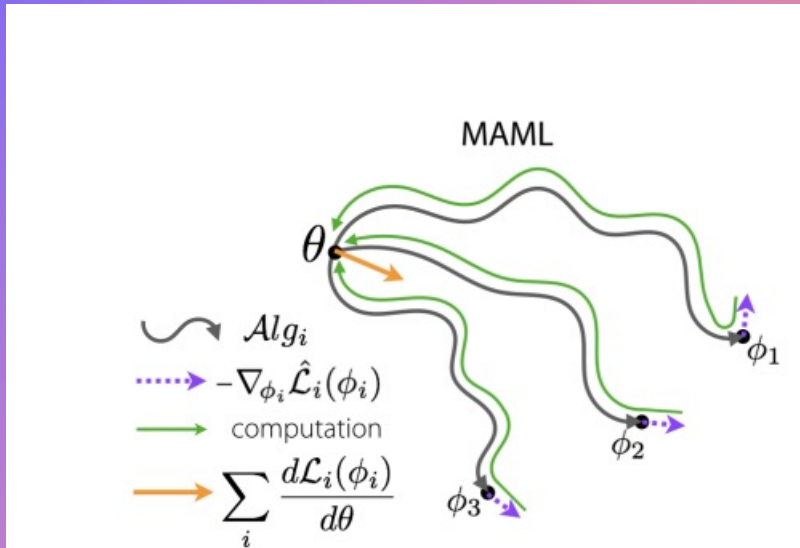
- It is referred 'Learning to Learn'
- It is a subfield of machine learning and cognitive science that focuses on developing algorithms and models that can learn how to learn.
- The goal is to enable machines to quickly adapt to new tasks or environments by acquiring new skills and knowledge from past experiences.
- It uses bi-level optimization criteria. The meta-learning can be divided into 2 parts, inner loop, and outer loop. Inner loop may optimize some cost function however, outer loop optimizes environment.
- One popular method for meta-learning with implicit gradients is MAML

MAML - Model-Agnostic Meta-Learning

- MAML uses a first-order optimization method

iMAML

- iMAML uses a second-order optimization method. It depends only on the solution to the inner level optimization and not the path taken by the inner loop optimizer



Formulation

Outer and Inner Level

$$\overbrace{\theta_{\text{ML}}^* := \operatorname{argmin}_{\theta \in \Theta} F(\theta)}^{\text{outer-level}},$$

$$\text{where } F(\theta) = \frac{1}{M} \sum_{i=1}^M \mathcal{L} \left(\overbrace{\operatorname{Alg}(\theta, \mathcal{D}_i^{\text{tr}})}^{\text{inner-level}}, \mathcal{D}_i^{\text{test}} \right).$$

$$\phi_i \equiv \operatorname{Alg}(\theta, \mathcal{D}_i^{\text{tr}}) = \theta - \alpha \nabla_{\theta} \mathcal{L}(\theta, \mathcal{D}_i^{\text{tr}}). \quad (\text{inner-level of MAML})$$

- θ_{ML} – Best Meta Learning parameter
- F - average validation loss function
- Alg - training algorithm performed on the training set of the particular task while starting from θ
- M – Number of tasks
- T_i – Task
- \mathcal{D}_i - Dataset
- $\mathcal{D}_i^{\text{tr}}$ – Train set
- $\mathcal{D}_i^{\text{test}}$ – Test set
- α – Learning Rate
- Φ_i – model parameters

Section 2.2 Proximal Regularization in the Inner Level

To have sufficient learning in the inner level while also avoiding over-fitting, \mathcal{Alg} needs to incorporate some form of regularization. Since MAML uses a small number of gradient steps, this corresponds to early stopping and can be interpreted as a form of regularization and Bayesian prior [20]. In cases like ill-conditioned optimization landscapes and medium-shot learning, we may want to take many gradient steps, which poses two challenges for MAML. First, we need to store and differentiate through the long optimization path of \mathcal{Alg} , which imposes a considerable computation and memory burden. Second, the dependence of the model-parameters $\{\phi_i\}$ on the meta-parameters (θ) shrinks and vanishes as the number of gradient steps in \mathcal{Alg} grows, making meta-learning difficult.

The old method has this constraint because of back propagation. New method don't have this issue and stays close to the actual parameter chosen by using λ which is a regularization parameter can be scalars, vectors or full matrices. But, in this paper they consider as scalar

$$\mathcal{Alg}^*(\theta, \mathcal{D}_i^{\text{tr}}) = \operatorname{argmin}_{\phi' \in \Phi} \mathcal{L}(\phi', \mathcal{D}_i^{\text{tr}}) + \frac{\lambda}{2} \|\phi' - \theta\|^2.$$

The proximal regularization term in Eq. 3 encourages ϕ_i to remain close to θ , thereby retaining a strong dependence throughout. The regularization strength (λ) plays a role similar to the learning rate (α) in MAML, controlling the strength of the prior (θ) relative to the data ($\mathcal{D}_T^{\text{tr}}$). Like α , the regularization strength λ may also be learned. Furthermore, both α and λ can be scalars, vectors, or full matrices. For simplicity, we treat λ as a scalar hyperparameter. In Eq. 3, we use \star to denote that

Formulation

Bi-level meta learning problem

$$\overbrace{\theta_{\text{ML}}^* := \operatorname{argmin}_{\theta \in \Theta} F(\theta)}^{\text{outer-level}},$$

$$F(\theta) = \frac{1}{M} \sum_{i=1}^M \mathcal{L}_i(\mathcal{A}l g_i^*(\theta)).$$

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

- θ_{ML} – Best Meta Learning parameter
- F - average validation loss function
- $\mathcal{A}l g$ - training algorithm performed on the training set of the particular task while starting from θ
- M – Number of tasks
- T_i – Task
- D_i - Dataset
- D_i^{tr} – Train set
- D_i^{test} – Test set
- α – Learning Rate
- λ - regularization parameter
- Φ_i – model parameters

How does it lead to gradient Descent

- We have to calculate $Df/D\theta$. The conventional approach involved blindly inverting the matrix, which is a computationally expensive process. So, the author suggests not to do that instead multiply the gradient with the invertible matrix.
- η outer step size

$$\frac{dAlg_i^*(\theta)}{d\theta} = \left(\mathbf{I} + \frac{1}{\lambda} \nabla_{\phi}^2 \hat{\mathcal{L}}_i(\phi_i) \right)^{-1}.$$

$$\theta \leftarrow \theta - \eta \frac{1}{M} \sum_{i=1}^M \frac{dAlg_i^*(\theta)}{d\theta} \nabla_{\phi} \mathcal{L}_i(Alg_i^*(\theta)).$$

Algorithm

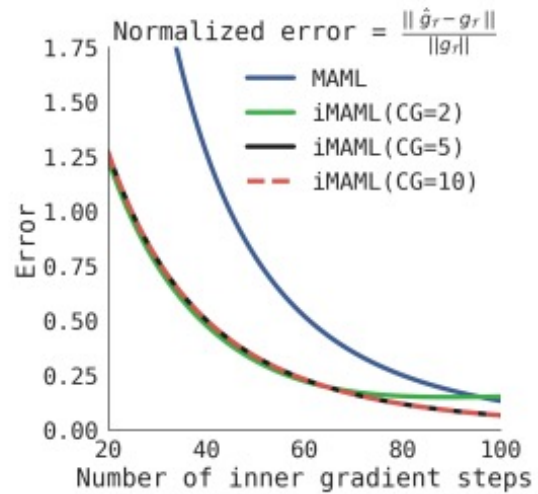
Algorithm 1 Implicit Model-Agnostic Meta-Learning (iMAML)

- 1: **Require:** Distribution over tasks $P(\mathcal{T})$, outer step size η , regularization strength λ ,
 - 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 $\mathbf{g}_i = \text{Implicit-Meta-Gradient}(\mathcal{T}_i, \boldsymbol{\theta}, \lambda)$
 - 6: **end for**
 - 7: Average above gradients to get $\hat{\nabla}F(\boldsymbol{\theta}) = (1/B) \sum_{i=1}^B \mathbf{g}_i$
 - 8: Update meta-parameters with gradient descent: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \hat{\nabla}F(\boldsymbol{\theta})$ // (or Adam)
 - 9: **end while**
-

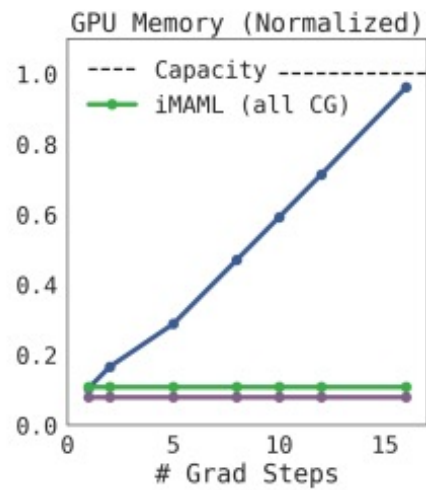
Algorithm 2 Implicit Meta-Gradient Computation

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

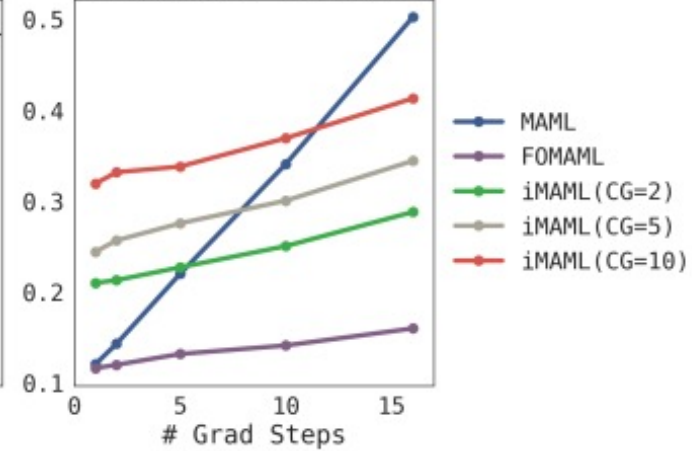
Graph



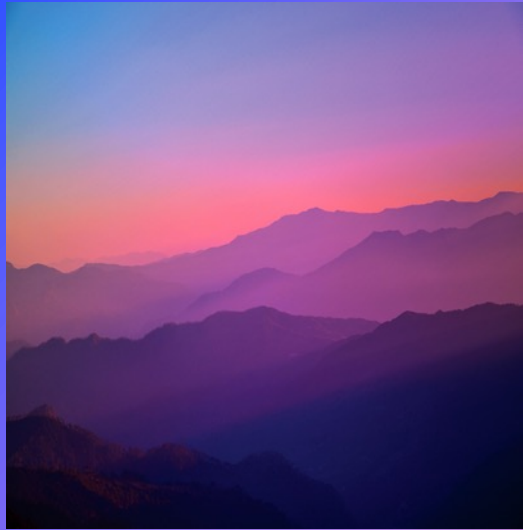
(a)



Compute Time (sec/iter)



(b)



Summary

In summary, meta-learning with implicit gradients is a promising approach for learning how to learn in scenarios where explicit gradients may not be available. These methods can enable models to quickly adapt to new tasks with minimal data, improving their efficiency and generalization. Also, this has the potential to improve the performance and versatility of autonomous systems, such as self-driving cars or robots

+



o



.



THANK YOU