

Representation Learning with Model-Agnostic Meta-Learning (MAML)

Aryan Mokhtari
ECE Department, UT Austin

ITA Workshop 2022

May 24, 2022

Research Collaborators



Liam Collins



Sewoong Oh

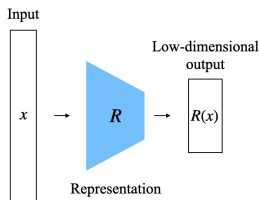


Sanjay Shakkottai

L. Collins, A. Mokhtari, S. Oh, S. Shakkottai. "MAML and ANIL Provably Learn Representations", *ICML 2022*. [<https://arxiv.org/abs/2202.03483>]

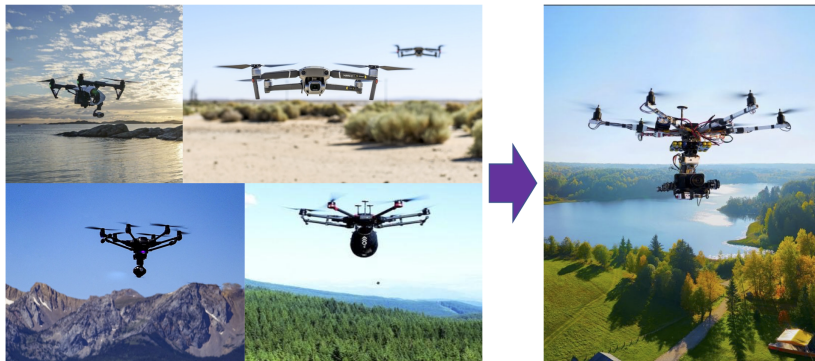
Representation Learning

- A central goal of machine learning:
⇒ **learn useful representations of data.**



- Good representations are useful because they
 - ⇒ reduce data complexity
 - ⇒ enable generalizing models to *new tasks* quickly
- Quickly = with little (labeled) data and computation.

Representation Learning – Generalize to New Tasks



Training Data

Task in new Scenario

Image Credits: bit.ly/3i5m8ay, bit.ly/3w723ZY, bit.ly/3KHMq5E, bit.ly/3i7pREJ, bit.ly/34I1ytT

- Training data in various scenarios; goal is to quickly adapt to the task in a new scenario

Meta-learning and Representation Learning

- *How* can we learn high-quality representations?
- *Meta-learning* methods have recently grown in popularity because they have yielded useful representations in practice.

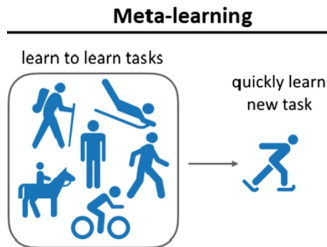


Image credit: <https://meta-world.github.io>, [HRJ21]

- Meta-learning leverages experience from learning a set of meta-training tasks to quickly solve new tasks.
- A popular approach is Model Agnostic Meta-Learning (MAML)

MAML and ERM

Traditional Supervised Learning, Empirical Risk Minimization (ERM)

- Set of tasks: $\mathcal{T} = \{\mathcal{T}_i\}_{i=1}^n$ coming from distribution p
- Select a model θ_{train}^*
- A new task \mathcal{T}_{test} is revealed, drawn according to dist. p
- Performance: $f_{test}(\theta_{train}^*)$

Formally, the training objective is:

- $\min_{\theta} \mathbb{E}_{i \sim p} [f_i(\theta)]$
- $\min_{\theta} \frac{1}{n} \sum_{i=1}^n f_i(\theta)$

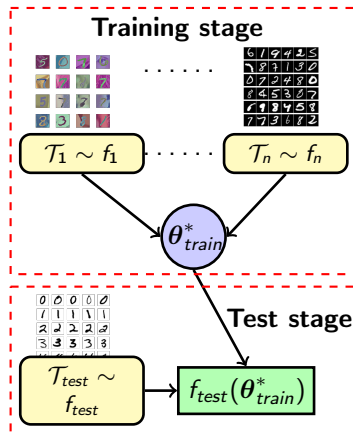


Image Credits: <https://bit.ly/392pda9>,
<https://bit.ly/3EEIElq>

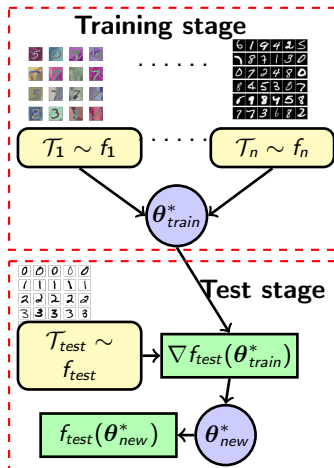
MAML and ERM

Model-Agnostic Meta Learning [FAL17]

- What if **we have budget to slightly update our model** at test time?
- $\mathcal{T} = \{\mathcal{T}_i\}_{i=1}^n$ drawn from distribution p
- Select a model θ_{train}^*
- \mathcal{T}_{test} is revealed, drawn based on p
- A few labeled samples of \mathcal{T}_{test} given
- Performance:
$$f_{test}(\theta_{train}^* - \alpha \nabla f_{test}(\theta_{train}^*))$$

Formally, the training objective is:

- $\min_{\theta} \mathbb{E}_{i \sim p} \mathbb{E}_{(X_i, y_i) \sim D_i} [f_i(\theta - \alpha \nabla f_i(\theta; X_i, y_i))]$
- $\min_{\theta} \frac{1}{n} \sum_{i=1}^n f_i(\theta - \alpha \nabla f_i(\theta; (X_i, y_i)))$



MAML Intuition: Adaptivity

- Empirical Risk Minimization (ERM): $\min_{\theta} \frac{1}{n} \sum_{i=1}^n f_i(\theta)$
- Gradient descent update for ERM.:

$$\theta_{t+1} = \theta_t - \frac{\beta}{n} \sum_{i=1}^n \nabla f_i(\theta_t) \quad (1)$$

Gradient evaluated at same θ_t for all tasks \implies **not adaptive**

- In contrast, for MAML the update is:

$$\theta_{t+1} = \theta_t - \frac{\beta}{n} \sum_{i=1}^n (\mathbf{I} - \alpha \nabla^2 f_i(\theta_t)) \nabla f_i(\theta_{t,i})$$

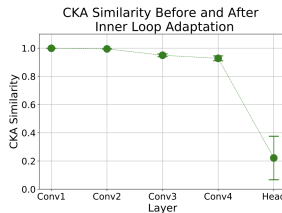
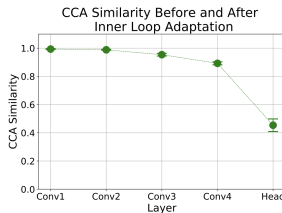
$$\text{where } \theta_{t,i} = \theta_t - \alpha \nabla f_i(\theta_t)$$

$\theta_{t,i}$ adapted to each task \implies MAML finds an **adaptive** solution

- Seems like finding the right initialization for adaptation!

Empirical Observations of MAML

- 1 MAML learns models that can quickly solve new tasks [FAL17, AES19]
 - ▶ in image classification, sinusoid regression, reinforcement learning.
- 2 MAML seems to be learning a representation shared across tasks [RRBV20]
 - ▶ even though it is not designed for representation learning!



The representation learned by MAML does not change significantly when adapted to each task.
Figure credit: [RRBV20]

- 3 Can we formally prove this?

Representation learning in multi-task linear regression (1/2)

- Example: consider multi-task linear regression.
- Task i has ground-truth solution $\theta_{*,i} \in \mathbb{R}^d$:

$$y_i \sim \theta_{*,i}^\top \mathbf{x}_i + z_i$$

- ▶ \mathbf{x}_i is a random feature vector
 - ▶ $z_i \in \mathbb{R}$ is random, mean-zero noise
- Solving each task individually (i.e. finding a $\theta_i \approx \theta_{*,i}$ for each i) would require $\Omega(d)$ samples per task.

Can we do better using shared information across tasks?

Representation learning in multi-task linear regression (2/2)

- Now suppose the $\theta_{*,i}$ lie in a shared k -dimensional subspace, $k \ll d$
- Let the columns of $\mathbf{B}_* \in \mathbb{R}^{d \times k}$ span this subspace, that is, for all tasks there exists $\mathbf{w}_{*,i} \in \mathbb{R}^k$ such that

$$\theta_{*,i} = \mathbf{B}_* \mathbf{w}_{*,i}$$

- ▶ \mathbf{B}_* is the “ground-truth” representation
- If we know $\text{col}(\mathbf{B}_*)$, we can solve new tasks with only $O(k)$ samples

Benefit of Representation Learning

$O(k)$ sample complexity much smaller than $\Omega(d)$

MAML for multi-task linear regression

- Loss function for task i at round t :

$$f_i(\mathbf{B}, \mathbf{w}) := \frac{1}{2} \mathbb{E}_{\mathbf{x}_i, y_i} [(\langle \mathbf{B}\mathbf{w}, \mathbf{x}_i \rangle - y_i)^2]$$

- At each time t , we sample n distinct tasks (may differ across time)
- For each task, we collect $m_{in} + m_{out}$ data samples (approximates the expectation in the population loss in the **inner** and **outer** loops)

Algorithm (MAML)

- (Outer loop) For $t = 1, \dots, T$:
 - Select n tasks satisfying diversity condition ($\text{span}(\{\mathbf{w}_{*,i}\}_{i \in [n]}) = \mathbb{R}^k$)
 - (Inner loop) For $i = 1, \dots, n$:
 - ▶ **Adapt:** $\mathbf{w}_{t,i} = \mathbf{w}_t - \alpha \nabla_{\mathbf{w}} f_i(\mathbf{B}_t, \mathbf{w}_t)$, $\mathbf{B}_{t,i} = \mathbf{B}_t - \alpha \nabla_{\mathbf{B}} f_i(\mathbf{B}_t, \mathbf{w}_t)$
 - $$\begin{bmatrix} \mathbf{w}_{t+1} \\ \bar{\mathbf{B}}_{t+1} \end{bmatrix} = \begin{bmatrix} \mathbf{w}_t \\ \bar{\mathbf{B}}_t \end{bmatrix} - \frac{\beta}{n} \sum_{i=1}^n (\mathbf{I} - \alpha \nabla_{\mathbf{w}, \bar{\mathbf{B}}}^2 f_i(\mathbf{B}_t, \mathbf{w}_t)) \begin{bmatrix} \nabla_{\mathbf{w}} f_i(\mathbf{B}_{t,i}, \mathbf{w}_{t,i}) \\ \nabla_{\mathbf{B}} f_i(\mathbf{B}_{t,i}, \mathbf{w}_{t,i}) \end{bmatrix}$$
- where $\bar{\mathbf{B}}$ is the column-wise vectorization of \mathbf{B} .

Our Main Results

- We consider the multi-task linear regression setting.

Main Results

- *Under standard assumptions, MAML (and variants) recover $\text{col}(\mathbf{B}_*)$ exponentially fast when run on the task population losses.*
- *ANIL and FO-ANIL (simplified MAML variants) require $m = O((\frac{d}{n} + 1)k^3) \ll d$ samples per task to learn the ground-truth subspace.*
- *The key is that MAML and variants' adaptation of the head harnesses **task diversity** to improve the representation in all directions.*
- **First results** showing that MAML and variants provably learn effective representations.

Proof Intuition

- For FO-ANIL (a simplified version of MAML), we have

$$\mathbf{B}_{t+1} = \mathbf{B}_t \underbrace{\left(\mathbf{I}_k - \frac{\beta}{n} \sum_{i=1}^n \mathbf{w}_{t,i} \mathbf{w}_{t,i}^\top \right)}_{\text{prior weight}} + \mathbf{B}_* \underbrace{\frac{\beta}{n} \sum_{i=1}^n \mathbf{w}_{*,t,i} \mathbf{w}_{t,i}^\top}_{\text{signal weight}}$$

- Suppose $\frac{1}{n} \sum_{i=1}^n \mathbf{w}_{t,i} \mathbf{w}_{t,i}^\top$ is full rank (i.e., the $\mathbf{w}_{t,i}$'s are diverse), then:

Key observation

Prior weight reduces energy from \mathbf{B}_t , and **signal weight** boosts energy from \mathbf{B}_* in all directions.

\Rightarrow **Head adaptation** and **task diversity** are critical!

- The **more diverse** the **adapted** $\mathbf{w}_{t,i}$'s (i.e. smaller condition number of $\frac{1}{n} \sum_{i=1}^n \mathbf{w}_{t,i} \mathbf{w}_{t,i}^\top$), the **faster convergence rate**.

Comparison to Empirical Risk Minimization (1/2)

$$\text{ERM: } \min_{\mathbf{B}, \mathbf{w}} \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{B}, \mathbf{w})$$

- In this case we can show:

$$\mathbf{B}_{t+1} = \mathbf{B}_t \underbrace{\left(\mathbf{I}_k - \beta \mathbf{w}_t \mathbf{w}_t^\top \right)}_{\text{prior weight}} + \mathbf{B}_* \underbrace{\frac{\beta}{n} \sum_{i=1}^n \mathbf{w}_{*,t,i} \mathbf{w}_t^\top}_{\text{signal weight}}$$

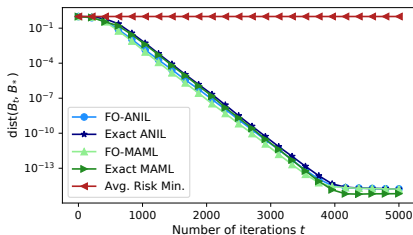
- The **prior weight** is rank $k - 1$, while the **signal weight** is only rank 1.

Key observation

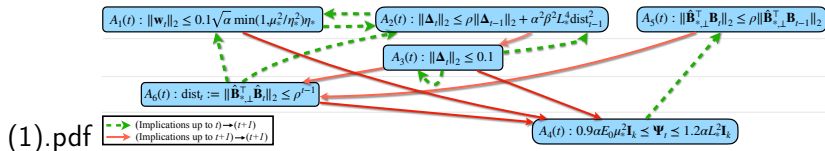
The representation can only move closer to $\text{col}(\mathbf{B}_*)$ in **one** direction on each iteration \rightarrow not clear that it can eventually reach $\text{col}(\mathbf{B}_*)$.

Comparison to Empirical Risk Minimization (2/2)

- Empirically, ERM fails to learn $\text{col}(\mathbf{B}_*)$.



(2).pdf



Inductive logic used in the proof for FO-ANIL.

- We have obtained the **first results** showing that MAML and variants learn effective representations in any setting.
- Inner loop adaptation of the head is **key** to their ability to learn representations.
- Quantifies the benefits of diverse tasks in the training environment.
- Substantial sample complexity improvement can be achieved by learning representations

References

[FAL17] Chelsea Finn, Pieter Abbeel, Sergey Levine. Model-Agnostic Meta-Learning for Fast Adaptation of Neural Networks, *International Conference on Machine Learning*, 2017.

[AES19] Antreas Antoniou, Harrison Edwards, Amos Storkey. How to Train Your MAML, *International Conference on Learning Representations*, 2019.

[RRBV20] Aniruddh Raghu, Maithra Raghu, Samy Bengio, Oriol Vinyals. Rapid Learning or Feature Reuse? Towards Understanding the Effectiveness of MAML, *International Conference on Learning Representations*, 2020.

[HRJ21] Mike Huisman, Jan N. van Rijn, Aske Plaat. A Survey of deep Meta-Learning, *Artificial Intelligence Review* Volume 54, pages 4483–4541, 2021.