

Offline Meta-Reinforcement Learning with
Advantage WeightingEric Mitchell, Rafael Rafailov, Xue Bin Peng,
Sergey Levine, Chelsea Finn

Overview

We introduce the **offline meta-reinforcement learning** (offline meta-RL) problem setting and propose an algorithm that performs well in this setting. Offline meta-RL is analogous to the widely successful supervised **pre-train + fine-tune** transfer learning strategy, in which a model is pre-trained on a large batch of fixed, pre-collected data (possibly from various tasks) and fine-tuned to a new task using relatively little data. That is, in offline meta-RL, we **meta-train on fixed, pre-collected data** from several tasks and **adapt to a new task with a very small amount** of data from the new task. By nature of being offline, algorithms for offline meta-RL can utilize the **largest possible pool of training data** available and eliminate potentially unsafe or costly data collection during meta-training.

This setting inherits the challenges of offline RL, but it **differs significantly** because offline RL does not generally consider a) **transfer to new tasks** or b) **limited data from the test task**, both of which we face in offline meta-RL. Targeting the offline meta-RL setting, we propose an algorithm, **Meta-Actor Critic with Advantage Weighting (MACAW)**. MACAW is an optimization-based meta-learning algorithm that uses simple, supervised regression objectives for both the inner and outer loop of meta-training. Our experiments show that MACAW enables fully offline meta-RL and demonstrates superior performance in a variety of settings including **offline variants of standard meta-RL benchmarks**, **adapting from sub-optimal data**, and **learning from a limited number of training tasks**.

Key Contributions

1. The **fully offline meta-reinforcement learning** problem setting
2. An algorithm for offline meta-reinforcement learning, called **Meta-Actor Critic with Advantage Weighting** or **MACAW**

Offline Meta-RL Problem Setting

We consider the **offline meta-reinforcement learning setting**. An offline meta-RL problem consists of:

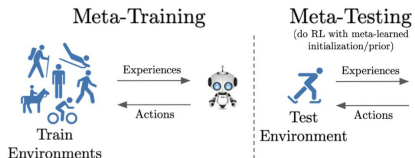
- Training tasks \mathcal{T} sampled from a distribution $p(\mathcal{T})$
- Fixed buffers of offline data $D_i = \{s_{i,j}, a_{i,j}, s'_{i,j}, r_{i,j}\}$ each training task

Each D_i is populated with trajectories sampled from a corresponding behavior policy μ_i

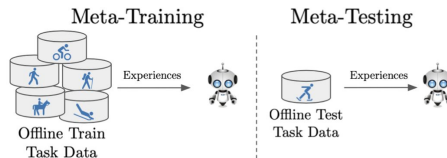
During **offline meta-training**, an agent trains on the data in the fixed training buffers **without interacting with the environment**.

During **fully offline meta-testing**, the agent is given a small amount of interaction data from a test task sampled from $p(\mathcal{T})$ which is used for adaptation **without interacting with the test environment**. The agent is then evaluated by its **average on-policy return** on the test task.

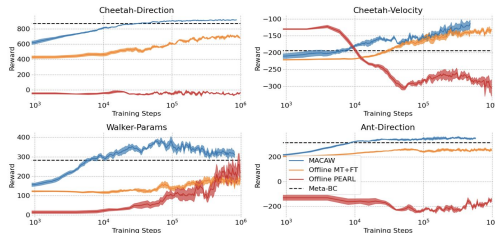
Standard Meta-RL



Offline Meta-RL



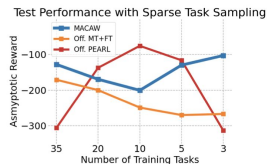
Benchmark Comparison



We evaluate MACAW on offline variants of **standard meta-RL problems** [1], [2]. We compare with an offline variant of PEARL [3], an offline **multi-task + fine-tuning** baseline based on AWR [4], and a **meta-imitation** baseline.

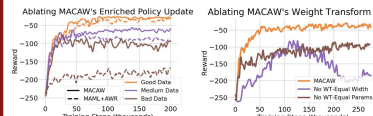
Training Task Sparsity Experiment

Generally, we prefer an offline meta-RL algorithm that can generalize to new tasks when presented with only a **small number of meta-training tasks**. This experiment evaluates the extent to which various algorithms rely on **dense sampling of the space of tasks** during training in order to generalize well.



Surprisingly, Offline PEARL **completely fails** to learn both when training tasks are plentiful and when they are scarce, but learns **relatively effectively** in the middle regime (5-20 tasks). In contrast, MACAW finds a solution of reasonable quality for **any sampling of the task space**, even for very dense or very sparse samplings of the training tasks. In practice, this property is desirable, because it allows the same algorithm to **scale to very large offline datasets while still producing useful adaptation behaviors for small datasets**. Ultimately, MACAW effectively exploits the available data when meta-training tasks are plentiful and shows by far the greatest robustness when tasks are scarce.

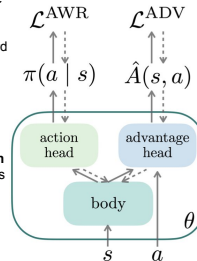
Ablation Studies



The left setting compares MACAW's performance with and without the **enriched policy update** when adaptation data **quality** is varied. When adaptation data is **sub-optimal**, MACAW's enriched policy update enables a **significant improvement** in performance. The right setting compares MACAW's performance with and without **weight transform layers**. The weight transform significantly improves both **speed of learning** as well as **asymptotic performance** of the policy. Both settings utilize the **Cheetah-Vel** problem.

MACAW Policy Architecture

MACAW's policy uses a **multi-headed** architecture in order to accommodate an **auxiliary policy loss** used to enrich the standard AWR [4] policy loss gradient. See **Theorems 1&2 in the paper**. During adaptation in both meta-training and meta-testing, the policy outputs actions and **action advantage estimates**; this eliminates the **policy gradient ambiguity** problem present when adapting using only the standard AWR policy loss.



Method

MACAW meta-learns **initializations** for a **value function** and **policy** during **meta-training** (see Algorithm 1 below). During **meta-testing**, MACAW takes a small number of gradient steps on the inner loop **value function** and **policy** objectives using a small batch of offline data (see Algorithm 2 below).

The **value function** objective is defined as (as in [4]):

$$\mathcal{L}_V(\phi, D) \triangleq \mathbb{E}_{s, a \sim D} [(V_\phi(s) - \mathcal{R}_D(s, a))^2]$$

The **inner loop** policy objective is defined as:

$$\mathcal{L}_\pi \triangleq \mathcal{L}^{\text{AWR}} + \lambda \mathcal{L}^{\text{ADV}}$$

Where \mathcal{L}^{ADV} is an auxiliary **advantage regression** loss designed to **increase inner loop policy update expressiveness**, defined as

$$\mathcal{L}^{\text{ADV}}(\theta, \phi'_i, D) \triangleq \mathbb{E}_{s, a \sim D} [(\hat{A}(s, a) - (\mathcal{R}_D(s, a) - V_{\phi'_i}(s)))^2]$$

and \mathcal{L}^{AWR} is the Advantage-Weighted Regression [4] policy loss:

$$\mathcal{L}^{\text{AWR}}(\theta, \phi, D) \triangleq \mathbb{E}_{s, a \sim D} \left[-\log \pi_\theta(a|s) \exp \left(\frac{1}{T} (\mathcal{R}_D(s, a) - V_\phi(s)) \right) \right]$$

The complete MACAW algorithm for both meta-training and meta-testing is described below.

Algorithm 1 MACAW Meta-Training

1. **Input:** Tasks $\{T_i\}$, offline buffers $\{D_i\}$
2. **Hyperparameters:** learning rates $\alpha_1, \alpha_2, \eta_1, \eta_2$, training iterations n , temperature T
3. Randomly initialize meta-parameters θ, ϕ
4. **for** n steps **do**
5. **for** task $T_i \in \{T_i\}$ **do**
6. Sample disjoint batches $D_i^n, D_i^n \sim D_i$
7. $\phi'_i \leftarrow \phi - \eta_1 \nabla_{\phi} \mathcal{L}_V(\phi, D_i^n)$
8. $\theta' \leftarrow \theta - \alpha_1 \nabla_{\theta} \mathcal{L}_\pi(\theta, \phi'_i, D_i^n)$
9. $\phi \leftarrow \phi - \eta_2 \sum_i [\nabla_{\phi} \mathcal{L}_V(\phi'_i, D_i^n)]$
10. $\theta \leftarrow \theta - \alpha_2 \sum_i [\nabla_{\theta} \mathcal{L}_\pi(\theta', \phi, D_i^n)]$

Algorithm 2 MACAW Meta-Testing

1. **Input:** Test task T_t , experience D_t , meta-policy π_θ , meta-value function V_ϕ
2. **Hyperparameters:** learning rates α_1, η , adaptation iterations n , temperature T
3. Initialize $\theta_t \leftarrow \theta, \phi_t \leftarrow \phi$
4. **for** n steps **do**
5. $\phi_{t+1} \leftarrow \phi_t - \eta_1 \nabla_{\phi} \mathcal{L}_V(\phi_t, D_t)$
6. $\theta_{t+1} \leftarrow \theta_t - \alpha_1 \nabla_{\theta} \mathcal{L}_\pi(\theta_t, \phi_{t+1}, D_t)$

References

- [1] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In International Conference on Machine Learning, 2017a.
- [2] Jonas Rothfuss, Dennis Lee, Ignasi Clavera, Tamim Asfour, and Pieter Abbeel. Prompt: Proximal meta-policy search. arXiv preprint arXiv:1810.06784, 2018.
- [3] Kate Rakelly, Aurick Zhou, Deirdre Quillen, Chelsea Finn, and Sergey Levine. Efficient off-policy meta-reinforcement learning via probabilistic context variables. In International Conference on Machine Learning, 2019.
- [4] Xue Bin Peng, Aviral Kumar, Grace Zhang, and Sergey Levine. Advantage-weighted regression: Simple and scalable off-policy reinforcement learning, 2019.

Correspondence to eric.mitchell@cs.stanford.edu

Paper link:

