

Efficient Off-Policy Meta-Reinforcement Learning via Probabilistic Context Variables

Kate Rakelly^{1*} Aurick Zhou^{1*} Deirdre Quillen¹ Chelsea Finn¹ Sergey Levine¹

ICML 2019

On-policy learning
 Only one policy used throughout the system to both explore and select actions.
 sample inefficiency
 eg: policy gradient

Off-policy learning
 Two policies, one for exploring and the other for action selection.
 eg: DQN

Introduction

□ Meta-Reinforcement Learning

The agent can leverage varied experiences from previous tasks to adapt quickly to the new task at hand.



https://blog.csdn.net/sinat_37422398

Motivation

Most meta-learning RL systems use on-policy learning. The general problem with on-policy learning is sample inefficiency.



Off-policy RL (SAC) is more sample efficient than on-policy approaches (PPO) by 1-2 orders of magnitude (PPO eventually nears SAC performance), figure from Haarnoja et al. 2018.

Tackle the problem of efficient off-policy meta-reinforcement learning.

Method

Meta-training : learn a probabilistic encoder that accumulates the necessary statistics from past experience into the context variables.
 Meta-testing : the context variables can be sampled and held constant for the duration of an episode, enabling temporally-extended exploration.



Method-Context Variable

A task : $\mathcal{T} = \{p(\mathbf{s}_0), p(\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{a}_t), r(\mathbf{s}_t, \mathbf{a}_t)\}$

 $\mathbf{c}_n^{\mathcal{T}} = (\mathbf{s}_n, \mathbf{a}_n, r_n, \mathbf{s}_n')$: one transition in the task \mathcal{T}

z : latent probabilistic context variable

Gaussian factors $\Psi_{\phi}(\mathbf{z}|\mathbf{c}_{n}) = \mathcal{N}(f_{\phi}^{\mu}(\mathbf{c}_{n}), f_{\phi}^{\sigma}(\mathbf{c}_{n}))$ $q_{\phi}(\mathbf{z}|\mathbf{c}_{1:N}) \propto \Pi_{n=1}^{N} \Psi_{\phi}(\mathbf{z}|\mathbf{c}_{n})$

permutation-invariant function of prior experience



Meta-training procedure

Method-Meta-training





Meta-training procedure

$$\begin{aligned} \mathcal{L}_{actor} = \mathbb{E}_{\substack{\mathbf{s} \sim \mathcal{B}, \mathbf{a} \sim \pi_{\theta} \\ \mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{c})}} \left[D_{\mathrm{KL}} \left(\pi_{\theta}(\mathbf{a}|\mathbf{s}, \bar{\mathbf{z}}) \middle\| \frac{\exp(Q_{\theta}(\mathbf{s}, \mathbf{a}, \bar{\mathbf{z}}))}{\mathcal{Z}_{\theta}(\mathbf{s})} \right) \right] \\ \\ \mathcal{L}_{critic} = \mathbb{E}_{\substack{(\mathbf{s}, \mathbf{a}, r, \mathbf{s}') \sim \mathcal{B} \\ \mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{c})}} \left[Q_{\theta}(\mathbf{s}, \mathbf{a}, \mathbf{z}) - (r + \bar{V}(\mathbf{s}', \bar{\mathbf{z}})) \right]^{2} \end{aligned}$$

7/23 Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor

Algorithm 2 PEARL Meta-testing

Require: test task $\mathcal{T} \sim p(\mathcal{T})$ 1: Initialize context $\mathbf{c}^{\mathcal{T}} = \{\}$ 2: **for** k = 1, ..., K **do** 3: Sample $z \sim q_{\phi}(\mathbf{z}|c^{\mathcal{T}})$ 4: Roll out policy $\pi_{\theta}(\mathbf{a}|\mathbf{s}, \mathbf{z})$ to collect data $D_k^{\mathcal{T}} = \{(\mathbf{s}_j, \mathbf{a}_j, \mathbf{s}'_j, r_j)\}_{j:1...N}$ 5: Accumulate context $\mathbf{c}^{\mathcal{T}} = \mathbf{c}^{\mathcal{T}} \cup D_k^{\mathcal{T}}$ 6: **end for**

Experiments-Sample Efficiency and Performance



Figure 3. **Meta-learning continuous control**. Test-task performance vs. samples collected during *meta-training*. Our approach PEARL outperforms previous meta-RL methods both in terms of asymptotic performance and meta-training sample efficiency across six benchmark tasks. Dashed lines correspond to the maximum return achieved by each baseline after 1e8 steps. By leveraging off-policy data during meta-training, PEARL is 20-100x more sample efficient than the baselines, and achieves consistently better or equal final performance compared to the best performing prior method in each environment. See Appendix A for the full timescale version of this plot.

Experiments-Posterior Sampling For Exploration





Figure 4. **Sparse 2D navigation**. The agent must navigate to a previously unseen goal (dark blue, other test goals in light blue) with reward given only when inside the goal radius – radius of 0.2 (illustrated) and 0.8 are tested here. The agent is trained to navigate to a training set of goals, then tested on a distinct set of unseen test goals. By using posterior sampling to explore efficiently, PEARL is able to start adapting to the task after collecting on average only 5 trajectories, outperforming MAESN (Gupta et al., 2018).

Experiments-Ablations



Figure 5. Recurrent encoder ablation. We compare our encoder architecture to a recurrent network. For the RNN, we sample context as trajectories rather than unordered transitions. Sampling the actor-critic batch as de-correlated transitions ("RNN decorrelated") fares much better than sampling trajectories ("RNN correlated"). In addition to better performance, our encoder has a computational advantage.

Figure 6. Context sampling ablation. PEARL samples context batches of recently collected transitions de-correlated with the batches sampled for training the actor-critic. We compare to sampling context from the entire history ("off-policy context"), as well as using the same sampled batch for the context and the actor-critic batch ("off-policy same batch").

Experiments-Ablations



Figure 7. **Deterministic latent context**. We compare PEARL to a variant with deterministic latent context on the sparse reward 2D navigation domain. As expected, without a mechanism for reasoning about uncertainty over tasks, this approach is unable to explore effectively and performs poorly.



Meta-Q-Learning

Rasool Fakoor¹, Pratik Chaudhari², Stefano Soatto¹, Alexander Smola¹

¹ Amazon Web Services

² University of Pennsylvania

Email: {fakoor, soattos, smola}@amazon.com, pratikac@seas.upenn.edu

ICLR 2020

Contribution

Q-learning is competitive with state-ofthe-art meta-RL algorithms if given access to a context variable that is a representation of the past trajectory.

- A multi-task objective to maximize the average reward across the training tasks is an effective method to meta-train RL policies.
- Past data from the meta-training replay buffer can be recycled to adapt the policy on a new task using off-policy updates.



$$\widehat{\theta}_{\text{meta}} = \arg \max_{\theta} \frac{1}{n} \sum_{k=1}^{n} \mathop{\mathbb{E}}_{\tau \sim D^{k}} \left[\ell^{k}(\theta) \right]$$

Method-Context Variable

 \square recurrent context variable z_t

Set z_t to the hidden state at time t of a Gated Recurrent Unit (GRU) model.

 z_t depends on $\{(x_i, u_i, r_i)\}_{i \le t}$.



Algorithm 1: MQL - Meta-training
Input: Set of training tasks \mathcal{D}_{meta}
 Initialize the replay buffer Initialize parameters θ of an off-policy method, e.g., TD3 while not done do
4 // Rollout and update policy
 Sample a task D ~ D_{meta} Gather data from task D using policy π_θ while feeding transitions through context GRU. Add trajectory to the replay buffer. 𝔥 ← Sample mini-batch from buffer
8 Update parameters θ using mini-batch θ and Eqn. (15)
9 $\theta_{\text{meta}} \leftarrow \theta$ 10 return θ_{meta} , replay buffer
$\widehat{\theta}_{\text{meta}} = \arg\min_{\theta} \frac{1}{n} \sum_{k=1}^{n} \mathbb{E}_{\tau \sim D^{k}} \left[\text{TD}^{2}(\theta) \right];$

Method-Meta-training



Analysis

MAML: $\ell^k_{\text{meta}}(\theta) = \ell^k(\theta + \alpha \nabla_{\theta} \ell^k(\theta))$

Taylor series expansion & Gradient:

$$\nabla \ell_{\text{meta}}^{k}(\theta) = \nabla \ell^{k}(\theta) + 2\alpha(m-1) \left(\nabla^{2}\ell^{k}(\theta)\right) \nabla \ell^{k}(\theta) + \mathcal{O}(\alpha^{2}).$$

$$gradient$$

$$\ell^{k}(\theta) + \alpha(m-1) \|\nabla \ell^{k}(\theta)\|_{2}^{2}$$

Method-Meta-testing/Adaptation

 $\Box \text{ Update the policy using the new data}$ $<math display="block">\arg \max_{\theta} \left\{ \sum_{\tau \sim D^{new}} \left[\ell^{new}(\theta) \right] - \frac{\lambda}{2} \|\theta - \widehat{\theta}_{meta}\|_2^2 \right\}.$ (18) $\Box \text{ Exploits the meta-training replay buffer}$ $<math display="block">\arg \max_{\theta} \left\{ \sum_{\tau \sim \mathcal{D}_{meta}} \left[\beta(\tau; D^{new}, \mathcal{D}_{meta}) \, \ell^{new}(\theta) \right] - \frac{\lambda}{2} \|\theta - \widehat{\theta}_{meta}\|_2^2 \right\}.$ (19) $\beta(\tau; D^{new}, \mathcal{D}_{meta}) = \frac{P(\tau \in D^{new})}{P(\tau \in D_{meta})}$ $\Box \text{ Effective Sample Size} \text{ Size } \text{ Size } \text{ Six that knowing but change of the meta-training of the meta-training$

$$\widehat{\text{ESS}} = \frac{1}{m} \frac{\left(\sum_{k=1}^{m} \beta(x_k)\right)^2}{\sum_{k=1}^{m} \beta(x_k)^2} \in [0, 1]$$
$$\lambda = 1 - \widehat{\text{ESS}}$$

$$\arg\max_{\theta} \left\{ \underset{\tau \sim D^{\text{new}}}{\mathbb{E}} \left[\ell^{\text{new}}(\theta) \right] + \underset{\tau \sim \mathcal{D}_{\text{meta}}}{\mathbb{E}} \left[\beta(\tau; D^{\text{new}}, \mathcal{D}_{\text{meta}}) \ \ell^{\text{new}}(\theta) \right] - \left(1 - \widehat{\text{ESS}} \right) \|\theta - \widehat{\theta}_{\text{meta}}\|_2^2 \right\}$$

Method-Meta-testing/Adaptation



$$\arg\max_{\theta} \left\{ \underset{\tau \sim \mathcal{D}_{\text{meta}}}{\mathbb{E}} \left[\beta(\tau; D^{\text{new}}, \mathcal{D}_{\text{meta}}) \ \ell^{\text{new}}(\theta) \right] - \frac{\lambda}{2} \|\theta - \widehat{\theta}_{\text{meta}}\|_2^2 \right\}.$$
(19)

Experiments

Figure 2: Average undiscounted return of TD3 and TD3-context compared with PEARL for validation tasks from four meta-RL environments. The agent fails to learn if the policy is conditioned only on the state. In contrast, everything else remaining same, if TD3 is provided access to context, the rewards are much higher. In spite of not adaptating on the validation tasks, TD3-context is comparable to PEARL.

Experiments

Figure 3: Comparison of the average undiscounted return of MQL (orange) against existing meta-RL algorithms on continuous-control environments. We compare against four existing algorithms, namely MAML (green), RL2 (red), PROMP (purple) and PEARL (blue). In all environments except Walker-2D-Params and Ant-Goal-2D, MQL is better or comparable to existing algorithms in terms of both sample complexity and final returns.

Experiments-Ablations

Figure 4: Ablation studies to examine various components of MQL.

Thanks