

# Multi-Task Offline Reinforcement Learning with Heuristic Conservative Soft Actor-Critic

Project of AI3601 Reinforcement Learning, 2024 Spring, SJTU

---

Group 1: Xiangyuan Xue<sup>\*†</sup>, Kailing Wang<sup>\*</sup>, Jiazi Bu<sup>\*</sup>

June 4, 2024

<sup>\*</sup>Equal contribution <sup>†</sup>Project Leader & Slides Presenter

# Introduction

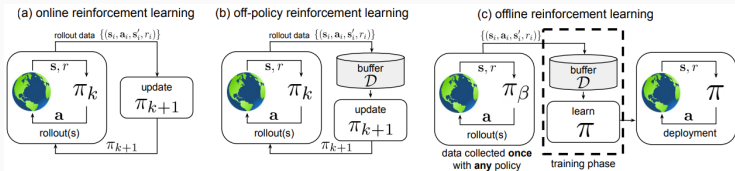
---

## Motivation

- Online RL requires the agent to frequently interact with the environment to gather samples, but such interaction can be extremely expensive.
- Collecting real-world samples can hardly be parallelized, which makes the process remarkably slow.
- The agent may encounter catastrophic failures while training, which is not affordable in some scenarios.
- In the era of large models, scaling up and generalizing the agent to various tasks is promising.

# Introduction

## Problem Definition



- Offline RL requires the agent to learn from a fixed dataset without interacting with the environment.
- Multi-task RL requires the agent to learn a satisfactory policy that can generalize to various tasks.
- We aim to combine offline RL with multi-task RL.

## Challenges

- In the offline setting, bootstrap methods unavoidably introduce distribution shifts, leading to overestimation of Q-values, which makes the traditional algorithms fail to learn a satisfactory policy.
- Most of the proposed multi-task RL methods are designed for the online setting, so they will suffer from the same issues when combined with offline RL and fail to achieve reliable performance.

## Method

---

## Preliminaries

- Markov decision process (MDP):  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, R, \gamma)$
- Bellman equation:

$$V^*(s) = \max_{\pi} \mathbb{E}_{a \sim \pi(\cdot|s)} [Q^{\pi}(s, a)]$$

$$Q^*(s, a) = R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[ \max_{a'} Q^*(s', a') \right]$$

- Offline dataset:  $\mathcal{D} = \{(s_i, a_i, r_i, s'_i)\}_{i=1}^N$
- Multi-task dataset:  $\mathcal{D} = \cup_{i=1}^M \mathcal{D}_i$

## Soft Actor-Critic (SAC)

- SAC maximizes the entropy of the policy to avoid suboptimal policies. The objective function is:

$$J(\pi) = \sum_{t=0}^{\infty} \mathbb{E}_{(s_t, a_t) \sim \rho_{\pi}} [R(s_t, a_t) + \alpha \mathcal{H}(\pi(\cdot | s_t))]$$

- The corresponding modified Bellman operator is:

$$\mathcal{T}^{\pi} Q(s, a) = R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [\mathbb{E}_{a' \sim \pi(\cdot | s')} [Q(s', a') - \alpha \log \pi(a' | s')]]$$

- Update policy by minimizing the KL divergence between the current policy and the softmax of Q-values:

$$\pi'(\cdot | s) = \arg \min_{\mu} D_{\text{KL}} \left( \mu(\cdot | s) \left\| \frac{\exp(Q^{\pi}(s, \cdot))}{Z^{\pi}(s)} \right\| \right)$$



## Soft Actor-Critic (SAC)

- Reparameterization trick is required for Gaussian policy.
- Two Q-networks are used to avoid the overestimation of Q-values, where soft target updates are applied.
- Automating entropy adjustment for maximum entropy:

$$\begin{aligned} \max_{\pi} \sum_{t=0}^{\infty} \mathbb{E}_{(s_t, a_t) \sim \rho_{\pi}} [R(s_t, a_t)] \\ \text{s.t. } \mathbb{E}_{(s_t, a_t) \sim \rho_{\pi}} [-\log \pi(a_t | s_t)] \geq \mathcal{H}_0, \forall t \end{aligned}$$

- Solve the optimization problem by Lagrange duality:

$$\alpha_t^* = \arg \min_{\alpha_t} \mathbb{E}_{a_t \sim \pi(\cdot | s_t)} [-\alpha_t \log \pi(a_t | s_t) - \alpha_t \mathcal{H}_0]$$

## Conservative Q-Learning (CQL)

- The key idea of CQL is to discourage the out-of-distribution Q-values while encouraging the Q-values in the dataset.
- The original form of CQL is formulated as:

$$\min_Q \max_{\mu} \alpha \left( \mathbb{E}_{s \sim \mathcal{D}, a \sim \mu(\cdot|s)} [Q(s, a)] - \mathbb{E}_{s \sim \mathcal{D}, a \sim \pi_{\beta}(\cdot|s)} [Q(s, a)] \right) \\ + \frac{1}{2} \mathbb{E}_{(s, a, s') \sim \mathcal{D}} \left[ \left( Q(s, a) - \mathcal{B}^{\pi_k} Q^k(s, a) \right)^2 \right] + \mathcal{R}(\mu)$$

- Note that  $\alpha$  is the temperature parameter,  $\mathcal{B}$  is the Bellman operator, and  $\mathcal{R}$  is a general regularizer.

## Conservative Q-Learning (CQL)

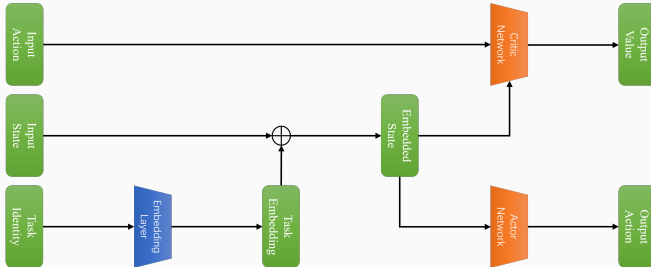
- If  $\mathcal{R}$  is the KL divergence between the current policy and the uniform distribution, it can be rearranged into:

$$\begin{aligned} \min_Q \alpha \mathbb{E}_{s \sim \mathcal{D}} \left[ \log \sum_a \exp(Q(s, a)) - \mathbb{E}_{a \sim \pi_\beta(\cdot|s)} [Q(s, a)] \right] \\ + \frac{1}{2} \mathbb{E}_{(s, a, s') \sim \mathcal{D}} \left[ \left( Q(s, a) - \mathcal{B}^{\pi_k} Q^k(s, a) \right)^2 \right] \end{aligned}$$

- In continuous action spaces, the log-sum-exp operation should be approximated by importance sampling:

$$\begin{aligned} \log \sum_a \exp(Q(s, a)) &= \log \left( \frac{1}{2} \sum_a \exp(Q(s, a)) + \frac{1}{2} \sum_a \exp(Q(s, a)) \right) \\ &= \log \left( \frac{1}{2} \mathbb{E}_{a \sim U(\mathcal{A})} \left[ \frac{\exp(Q(s, a))}{U(a)} \right] + \frac{1}{2} \mathbb{E}_{a \sim \pi(\cdot|s)} \left[ \frac{\exp(Q(s, a))}{\pi(a|s)} \right] \right) \end{aligned}$$

## Task Embedding



- We assign a unique identity  $k$  to each task. The identity will be converted into the task embedding  $e$  and added onto the state  $s$  as the task-aware state  $z$ .

## Heuristic Relabeling

- We share the transitions between tasks by relabeling dataset, where out-of-distribution samples are discarded.
- The threshold is defined by conservative Q-values:

$$\tau(\mathcal{D}_i) = \{Q_i(s, a) | (s, a) \in \mathcal{D}_i\}_{k\%}$$

- We relabel the transition only if its conservative Q-value is no smaller than the threshold, namely:

$$\Delta(s, a) = Q_i(s, a) - \tau(\mathcal{D}_i) \geq 0$$

## Heuristic Relabeling

---

**Algorithm 1** Heuristic Relabeling

---

**Require:** The multi-task offline dataset  $\mathcal{D} = \cup_{i=1}^M \mathcal{D}_i$ .

```
1: Initialize the active dataset  $\mathcal{D}^* \leftarrow \mathcal{D}$ .
2: for  $k \leftarrow 1, 2, \dots, L$  do
3:   Train the agent using samples from  $\mathcal{D}^*$  to obtain the latest critic network  $Q$ .
4:   for  $i \leftarrow 1, 2, \dots, M$  do
5:     Initialize the active dataset  $\mathcal{D}_i^* \leftarrow \mathcal{D}_i$ .
6:     Compute the relabeling threshold  $\tau(\mathcal{D}_i) \leftarrow \{Q(s, a) | (s, a) \in \mathcal{D}_i\}_{k\%}$ .
7:     for  $j \leftarrow 1, 2, \dots, M$  where  $j \neq i$  do
8:       Extend dataset by relabeling  $\mathcal{D}_i^* \leftarrow \mathcal{D}_i^* \cup \{(s, a, r, s') \in \mathcal{D}_j | Q_i(s, a) \geq \tau(\mathcal{D}_i)\}$ .
9:     end for
10:  end for
11:  Merge the active dataset  $\mathcal{D}^* \leftarrow \cup_{i=1}^M \mathcal{D}_i^*$ .
12: end for
```

---

- The algorithm is as the above pseudocode.

# Experiments

---

## Experiment Setup

- Behavioral cloning (BC): weak baseline.
- Conservative soft actor-critic (CSAC): strong baseline.
- Multi-task conservative soft actor-critic (MTCSAC): Apply task embedding based on CSAC.
- Heuristic conservative soft actor-critic (HCSAC): Apply heuristic relabeling based on MTCSAC.
- All the critic and actor networks are implemented as MLPs with 3 hidden layers and 256 hidden units.



# Experiments

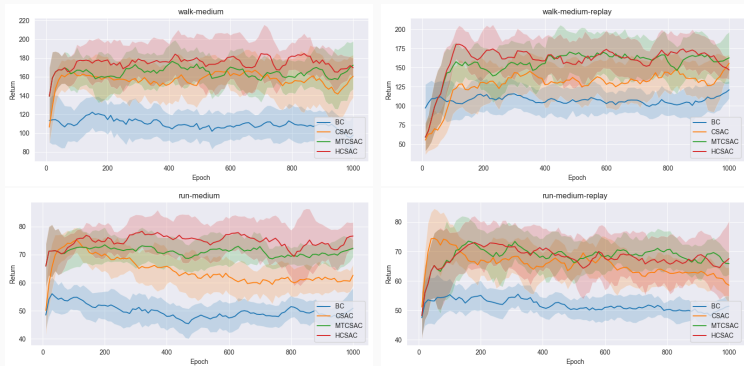
## Experiment Setup

- The experiments are divided into two groups (i.e. **medium** and **medium-replay**) according to the dataset.
- In each group, BC and CSAC are trained on the two tasks (i.e. **walk** and **run**) separately, while MTCSAC and HCSAC are trained on the union of the datasets for the two tasks.
- Train 1000 epochs with a batch size of 256. Online evaluation is conducted every 10 epochs. Learning rate is set to  $3 \times 10^{-4}$  for both critic and actor networks.

## Experiment Setup

- The soft target update rate is set to  $5 \times 10^{-3}$ .
- The discount factor is set to 0.99.
- The weight of the CQL loss term is set to 5.0.
- We sample 10 actions each from uniform distribution and policy distribution to estimate the log-sum-exp term.
- For HCSAC, heuristic relabeling is conducted every 100 epochs, where the quantile is set to 5%.

## Performance Evaluation



- Our methods perform better in terms of training curves.

# Experiments

## Performance Evaluation

- MTCSAC and HCSAC achieve higher average returns with lower standard deviations compared to BC and CSAC.
- Task embedding and heuristic relabeling in general bring superior performance and improve the robustness.

Dataset	Task	BC	CSAC	MTCSAC	HCSAC
medium	walk	113.96 $\pm$ 34.86	163.10 $\pm$ 21.00	170.22 $\pm$ 19.61	<b>173.41 <math>\pm</math> 11.72</b>
	run	54.71 $\pm$ 7.88	<b>77.23 <math>\pm</math> 7.11</b>	72.46 $\pm$ 6.25	76.06 $\pm$ 5.61
	average	84.34 $\pm$ 21.39	120.17 $\pm$ 14.56	121.34 $\pm$ 12.93	<b>124.74 <math>\pm</math> 8.67</b>
replay	walk	136.78 $\pm$ 20.96	131.73 $\pm$ 39.88	<b>184.32 <math>\pm</math> 20.23</b>	158.17 $\pm$ 25.72
	run	54.91 $\pm$ 3.39	<b>76.65 <math>\pm</math> 4.36</b>	66.62 $\pm$ 11.22	69.44 $\pm$ 9.14
	average	95.85 $\pm$ 18.18	104.19 $\pm$ 22.26	<b>125.47 <math>\pm</math> 15.72</b>	113.81 $\pm$ 17.43

## Contribution

- We propose heuristic conservative soft actor-critic to solve the offline multi-task reinforcement learning problem.
- We conduct a series of experiments on the MuJoCo benchmark to verify the effectiveness of our method.
- The results show that task embedding and heuristic relabeling in general bring superior performance and improve the robustness of the agent.

# Conclusion

## Discussion

- The tasks **walk** and **run** both rely on low-level locomotion skills. Sharing the policy may be unnecessary because we expect the agent can learn from a large number of tasks to acquire those essential skills in real-world scenarios.
- In MuJoCo, the vectorized observations already contain all the information needed for robotic control, so the multi-task learning methods cannot be fully exploited.
- Our methods may work even better in more complex environments, where the agent can utilize diverse transitions to learn better representations.

Thank You!