



# Paper Presentation

Group 2

309505018 郭俊廷  
309551026 黃柏愷



---

# **Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor**

---

ICML 2018

# Outline



1. Problem we want to solve
2. What's different between Soft Actor-Critic and others
3. Advantages of Maximum Entropy
4. Meticulous hyperparameter tuning
5. Experiment Result

## Problem we want to solve



Some famous model-free algorithm in Reinforcement Learning, like DDPG and PPO, will learn a policy to do a continuous control or a series of action.

PPO (Proximal Policy Optimization) is an on-policy algorithm, it performs very well on discrete and continuous control problems. But the weakness of PPO is sample inefficiency problem.

DDPG (Deep Deterministic Policy Gradient) is an off-policy algorithm, more sample efficient than PPO. But the weakness of DDPG is a deterministic policy, that means it will select the best action when policy gradient, and lacks of stochastic.

## Problem we want to solve



Soft Actor-Critic (SAC) is an off-policy algorithm, and it used “Maximum Entropy” method to update policy.

Compared to PPO, it solved sample inefficiency problem.

Compared to DDPG, it used stochastic policy which performs better than deterministic policy in many environment.

The most important, Soft Actor-Critic (SAC) is an Open Source project, so that we can replicate it easy. It can use in robotics to solve real world problem.

## What's different between Soft Actor-Critic and others



Original Policy Gradient

$$\pi^* = \arg \max_{\pi} \mathbb{E}_{(s_t, a_t) \sim \rho_{\pi}} \left[ \sum_t R(s_t, a_t) \right]$$

Add Maximum Entropy method

$$\pi^* = \arg \max_{\pi} \mathbb{E}_{(s_t, a_t) \sim \rho_{\pi}} \left[ \underbrace{\sum_t R(s_t, a_t)}_{\text{reward}} + \underbrace{\alpha H(\pi(\cdot | s_t))}_{\text{entropy}} \right]$$



## Soft Policy Iteration

$$\mathcal{T}^{\pi} Q(\mathbf{s}_t, \mathbf{a}_t) \triangleq r(\mathbf{s}_t, \mathbf{a}_t) + \gamma \mathbb{E}_{\mathbf{s}_{t+1} \sim p} [V(\mathbf{s}_{t+1})],$$

where

$$V(\mathbf{s}_t) = \mathbb{E}_{\mathbf{a}_t \sim \pi} [Q(\mathbf{s}_t, \mathbf{a}_t) - \log \pi(\mathbf{a}_t | \mathbf{s}_t)]$$

$$r_{soft}(s_t, a_t) = r(s_t, a_t) + \gamma \alpha \mathbb{E}_{s_{t+1} \sim \rho} H(\pi(\cdot | s_{t+1}))$$

$$Q(s_t, a_t) = r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1}, a_{t+1}} [Q(s_{t+1}, a_{t+1})]$$

$$\begin{aligned} Q_{soft}(s_t, a_t) &= r(s_t, a_t) + \gamma \alpha \mathbb{E}_{s_{t+1} \sim \rho} H(\pi(\cdot | s_{t+1})) + \gamma \mathbb{E}_{s_{t+1}, a_{t+1}} [Q_{soft}(s_{t+1}, a_{t+1})] \\ &= r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim \rho, a_{t+1} \sim \pi} [Q_{soft}(s_{t+1}, a_{t+1})] + \gamma \alpha \mathbb{E}_{s_{t+1} \sim \rho} H(\pi(\cdot | s_{t+1})) \\ &= r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim \rho, a_{t+1} \sim \pi} [Q_{soft}(s_{t+1}, a_{t+1})] + \gamma \mathbb{E}_{s_{t+1} \sim \rho} \mathbb{E}_{a_{t+1} \sim \pi} [-\alpha \log \pi(a_{t+1} | s_{t+1})] \\ &= r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim \rho} [\mathbb{E}_{a_{t+1} \sim \pi} [Q_{soft}(s_{t+1}, a_{t+1}) - \alpha \log(\pi(a_{t+1} | s_{t+1}))]] \\ &= r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1}, a_{t+1}} [Q_{soft}(s_{t+1}, a_{t+1}) - \alpha \log(\pi(a_{t+1} | s_{t+1}))] \end{aligned}$$





## Soft Policy Iteration

$$\begin{aligned}\pi_{\text{new}}(\cdot | \mathbf{s}_t) &= \arg \min_{\pi' \in \Pi} D_{\text{KL}}(\pi'(\cdot | \mathbf{s}_t) \parallel \exp(Q^{\pi_{\text{old}}}(\mathbf{s}_t, \cdot) - \log Z^{\pi_{\text{old}}}(\mathbf{s}_t))) \\ &= \arg \min_{\pi' \in \Pi} J_{\pi_{\text{old}}}(\pi'(\cdot | \mathbf{s}_t)).\end{aligned}$$




$$J_{\pi_{\text{old}}}(\pi_{\text{new}}(\cdot | \mathbf{s}_t)) \leq J_{\pi_{\text{old}}}(\pi_{\text{old}}(\cdot | \mathbf{s}_t)).$$

$$\mathbb{E}_{\mathbf{a}_t \sim \pi_{\text{new}}} [\log \pi_{\text{new}}(\mathbf{a}_t | \mathbf{s}_t) - Q^{\pi_{\text{old}}}(\mathbf{s}_t, \mathbf{a}_t) + \log Z^{\pi_{\text{old}}}(\mathbf{s}_t)] \leq$$

$$\mathbb{E}_{\mathbf{a}_t \sim \pi_{\text{old}}} [\log \pi_{\text{old}}(\mathbf{a}_t | \mathbf{s}_t) - Q^{\pi_{\text{old}}}(\mathbf{s}_t, \mathbf{a}_t) + \log Z^{\pi_{\text{old}}}(\mathbf{s}_t)]$$

$$\mathbb{E}_{\mathbf{a}_t \sim \pi_{\text{new}}} [Q^{\pi_{\text{old}}}(\mathbf{s}_t, \mathbf{a}_t) - \log \pi_{\text{new}}(\mathbf{a}_t | \mathbf{s}_t)] \geq V^{\pi_{\text{old}}}(\mathbf{s}_t)$$


$$\begin{aligned} Q^{\pi_{\text{old}}}(\mathbf{s}_t, \mathbf{a}_t) &= r(\mathbf{s}_t, \mathbf{a}_t) + \gamma \mathbb{E}_{\mathbf{s}_{t+1} \sim p} [V^{\pi_{\text{old}}}(\mathbf{s}_{t+1})] \\ &\leq r(\mathbf{s}_t, \mathbf{a}_t) + \gamma \mathbb{E}_{\mathbf{s}_{t+1} \sim p} [\mathbb{E}_{\mathbf{a}_{t+1} \sim \pi_{\text{new}}} [Q^{\pi_{\text{old}}}(\mathbf{s}_{t+1}, \mathbf{a}_{t+1}) - \log \pi_{\text{new}}(\mathbf{a}_{t+1} | \mathbf{s}_{t+1})]] \\ &\quad \vdots \\ &\leq Q^{\pi_{\text{new}}}(\mathbf{s}_t, \mathbf{a}_t), \end{aligned}$$

# Advantages of Maximum Entropy Reinforcement Learning



1. Through maximum entropy, policy not only learns a way to solve tasks, but all. Therefore, such a policy is more conducive to learning new tasks.
2. Stronger exploration ability, it is easier to find better modes under multimodal reward.
3. More robust and generalized. Because it is necessary to explore various optimal possibilities in different ways, it is easier to make adjustments in the face of interference.

# Implementation of SAC

---

## Algorithm 1 Soft Actor-Critic

---

**Input:**  $\theta_1, \theta_2, \phi$

$\theta_1 \leftarrow \theta_1, \theta_2 \leftarrow \theta_2$

$\mathcal{D} \leftarrow \emptyset$

**for** each iteration **do**

**for** each environment step **do**

$\mathbf{a}_t \sim \pi_\phi(\mathbf{a}_t | \mathbf{s}_t)$

$\mathbf{s}_{t+1} \sim p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$

$\mathcal{D} \leftarrow \mathcal{D} \cup \{(\mathbf{s}_t, \mathbf{a}_t, r(\mathbf{s}_t, \mathbf{a}_t), \mathbf{s}_{t+1})\}$

**end for**

**for** each gradient step **do**

$\theta_i \leftarrow \theta_i - \lambda_Q \hat{\nabla}_{\theta_i} J_Q(\theta_i)$  for  $i \in \{1, 2\}$

$\phi \leftarrow \phi - \lambda_\pi \hat{\nabla}_\phi J_\pi(\phi)$

$\alpha \leftarrow \alpha - \lambda \hat{\nabla}_\alpha J(\alpha)$

$\bar{\theta}_i \leftarrow \tau \theta_i + (1 - \tau) \bar{\theta}_i$  for  $i \in \{1, 2\}$

**end for**

**end for**

**Output:**  $\theta_1, \theta_2, \phi$

▷ Initial parameters

▷ Initialize target network weights

▷ Initialize an empty replay pool

▷ Sample action from the policy

▷ Sample transition from the environment

▷ Store the transition in the replay pool

▷ Update the Q-function parameters

▷ Update policy weights

▷ Adjust temperature

▷ Update target network weights

▷ Optimized parameters

---

$$J_Q(\theta) = \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t) \sim \mathcal{D}} \left[ \frac{1}{2} \left( Q_\theta(\mathbf{s}_t, \mathbf{a}_t) - \hat{Q}(\mathbf{s}_t, \mathbf{a}_t) \right)^2 \right] \quad \hat{Q}(\mathbf{s}_t, \mathbf{a}_t) = r(\mathbf{s}_t, \mathbf{a}_t) + \gamma \mathbb{E}_{\mathbf{s}_{t+1} \sim p} [V_{\bar{\psi}}(\mathbf{s}_{t+1})]$$

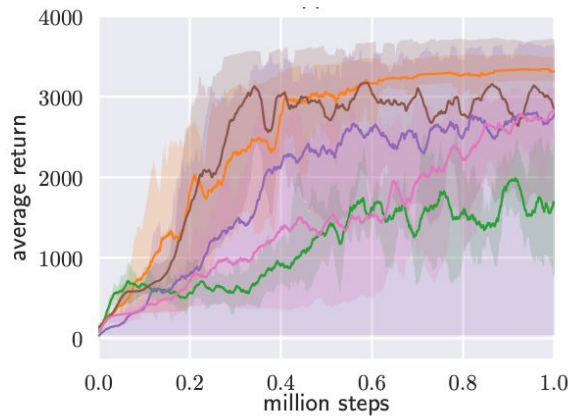
$$J_\pi(\phi) = \mathbb{E}_{\mathbf{s}_t \sim \mathcal{D}} \left[ \text{D}_{\text{KL}} \left( \pi_\phi(\cdot | \mathbf{s}_t) \parallel \frac{\exp(Q_\theta(\mathbf{s}_t, \cdot))}{Z_\theta(\mathbf{s}_t)} \right) \right]$$

$$J_\pi(\phi) = \mathbb{E}_{\mathbf{s}_t \sim \mathcal{D}, \epsilon_t \sim \mathcal{N}} [\log \pi_\phi(f_\phi(\epsilon_t; \mathbf{s}_t) | \mathbf{s}_t) - Q_\theta(\mathbf{s}_t, f_\phi(\epsilon_t; \mathbf{s}_t))]$$

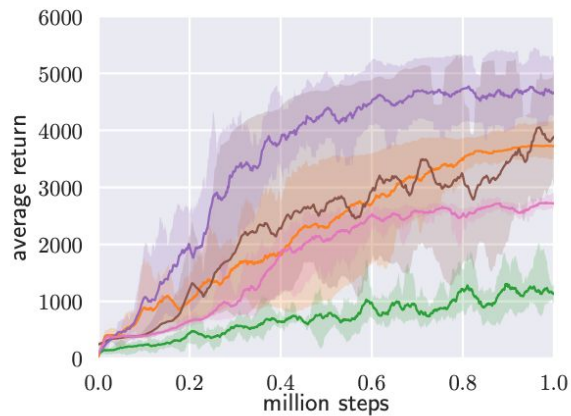
$$\tilde{a}_\theta(s, \xi) = \tanh \left( \mu_\theta(s) + \sigma_\theta(s) \odot \xi \right), \quad \xi \sim \mathcal{N}(0, I).$$



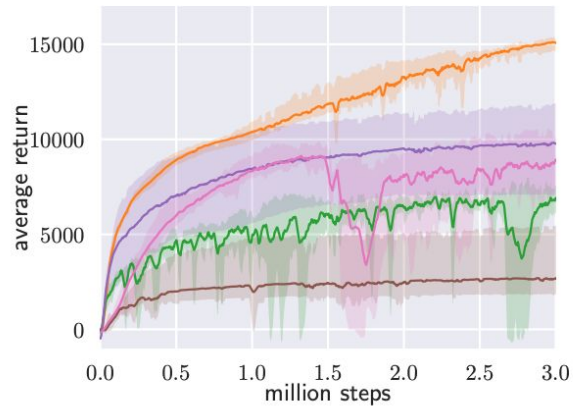
# Experiment Result



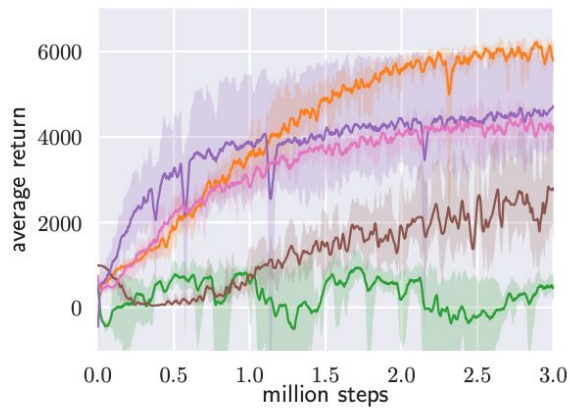
(a) Hopper-v1



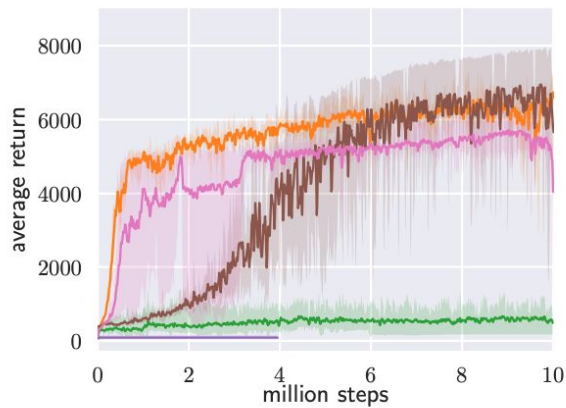
(b) Walker2d-v1



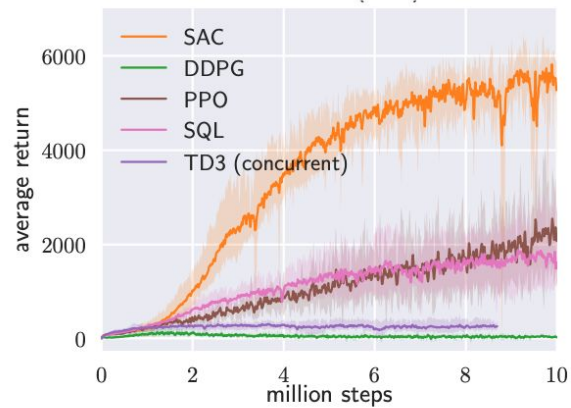
(c) HalfCheetah-v1



(d) Ant-v1



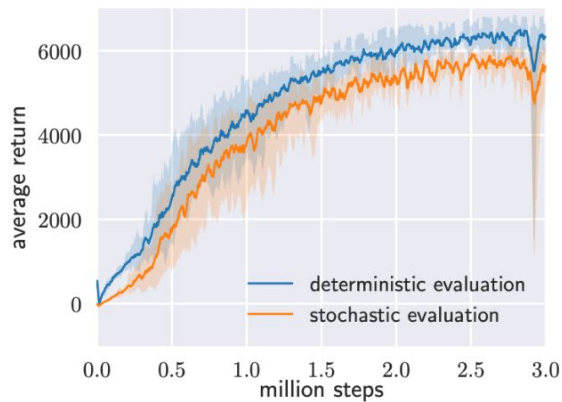
(e) Humanoid-v1



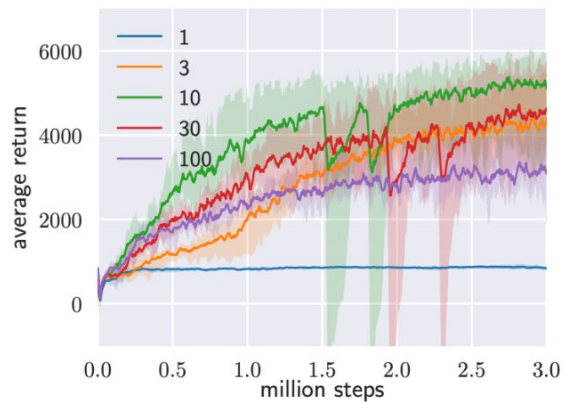
(f) Humanoid (rllab)



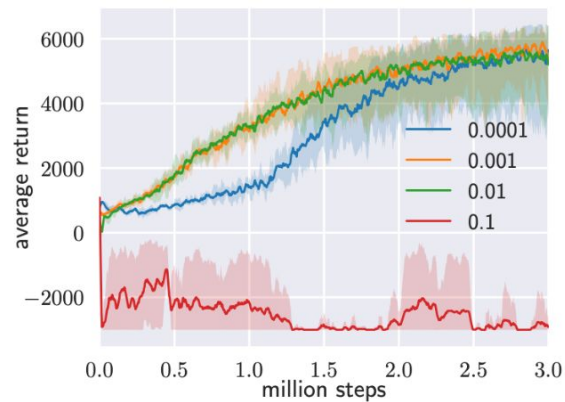
# Ablation Study



(a) Evaluation



(b) Reward Scale



(c) Target Smoothing Coefficient ( $\tau$ )