
Algorithm 1: Soft Actor-Critic (SAC)

Input: Initial policy parameters θ , Q-function parameters ϕ_1, ϕ_2 , target Q-function parameters $\bar{\phi}_1 \leftarrow \phi_1, \bar{\phi}_2 \leftarrow \phi_2$, temperature parameter α , polyak averaging coefficient τ , empty replay buffer \mathcal{D}

```
1 for each iteration do
2   Reset environment and observe initial state  $s_0$ ;
3   while not terminal do
4     Sample action  $a_t \sim \pi_\theta(a_t|s_t)$ ;
5     Execute  $a_t$  in environment;
6     Observe reward  $r_t$ , next state  $s_{t+1}$ , and terminal signal  $d_t$ ;
7     Store  $(s_t, a_t, r_t, s_{t+1}, d_t)$  in replay buffer  $\mathcal{D}$ ;
8     if  $d_t$  is True then
9       | Reset environment and observe new initial state  $s_{t+1}$ ;
10    end
11     $s_t \leftarrow s_{t+1}$ ;
12  end
13  for each gradient step do
14    Sample a mini-batch of  $N$  transitions  $(s, a, r, s', d)$  from  $\mathcal{D}$ ;
15    Compute target Q-value:
16
17    
$$y = r + \gamma(1-d) \left( \min_{i=1,2} Q_{\bar{\phi}_i}(s', a') - \alpha \log \pi_\theta(a'|s') \right), \quad a' \sim \pi_\theta(\cdot|s')$$

18
19    Update Q-functions by gradient descent:
20
21    
$$\phi_i \leftarrow \phi_i - \lambda_Q \nabla_{\phi_i} \frac{1}{N} \sum (Q_{\phi_i}(s, a) - y)^2, \quad i = 1, 2$$

22
23    Update policy by gradient ascent:
24
25    
$$\theta \leftarrow \theta - \lambda_\pi \nabla_\theta \frac{1}{N} \sum \left( \alpha \log \pi_\theta(a|s) - \min_{i=1,2} Q_{\phi_i}(s, a) \right), \quad a \sim \pi_\theta(\cdot|s)$$

26
27    Update temperature:
28
29    
$$\alpha \leftarrow \alpha - \lambda_\alpha \nabla_\alpha \frac{1}{N} \sum (-\alpha \log \pi_\theta(a|s) - \alpha \mathcal{H})$$

30
31    Update target Q-functions:
32
33    
$$\bar{\phi}_i \leftarrow \tau \phi_i + (1 - \tau) \bar{\phi}_i, \quad i = 1, 2$$

34  end
35 end
```
