

SAC Algorithm

```
1: Input:  $\theta_1, \theta_2, \phi$ 
2:  $\bar{\theta}_1 \leftarrow \theta_1, \bar{\theta}_2 \leftarrow \theta_2$ 
3:  $D \leftarrow \emptyset$ 
4: for each iteration do
5:   for each environment step do
6:      $a_t \sim \pi_\phi(a_t|s_t)$ 
7:      $s_{t+1} \sim p(s_{t+1}|s_t, a_t)$ 
8:      $D \leftarrow D \cup \{(s_t, a_t, R(s_t, a_t), s_{t+1})\}$ 
9:   end for
10:  for each gradient step do
11:     $\theta_i \leftarrow \theta_i - \lambda_Q \hat{\nabla}_{\theta_i} J_Q(\theta_i)$  for  $i \in \{1, 2\}$ 
12:     $\phi \leftarrow \phi - \lambda_\pi \hat{\nabla}_\phi J_\pi(\phi)$ 
13:     $\alpha \leftarrow \alpha - \lambda_\alpha \hat{\nabla}_\alpha J(\alpha)$ 
14:     $\bar{\theta}_i \leftarrow \tau \theta_i + (1 - \tau) \bar{\theta}_i$  for  $i \in \{1, 2\}$ 
15:  end for
16: end for
17: Output:  $\theta_1, \theta_2, \phi$ 
```

▷ Initial critic and actor networks parameters
▷ Initialize target network weights
▷ Initialize an empty replay buffer

▷ Sample action from the policy
▷ Receive transition from the environment
▷ Store the transition in the replay buffer

▷ Update critic networks Q-functions
▷ Update actor network policy
▷ Adjust temperature
▷ Update target network weights

▷ Optimized critic and actor networks parameters