
Algorithm 1: SAC-DRND (Online)

Input: Number of critic updates G , number of DRND updates per iteration K , target update coefficient ρ , entropy coefficient β , DRND coefficients $\lambda_{\text{actor}}, \lambda_{\text{critic}}$, replay buffer \mathcal{D}

Initialization: Initialize policy parameters ϕ , Q-function parameters φ_1, φ_2 and targets φ'_1, φ'_2 , predictor network θ and targets $\{\theta_n\}_{n=1}^N$

1 Reset environment and observe initial state s_0 ;

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

3 Reset environment and observe initial state s_0 ;

4 **while** not terminal **do**

5 Sample action $a_t \sim \pi_\phi(\cdot|s_t)$;

6 Execute a_t , observe r_t, s_{t+1} ;

7 Store $(s_t, a_t, r_t, b(s_{t+1}, a_t), s_{t+1})$ in \mathcal{D} ;

8 $s_t \leftarrow s_{t+1}$;

9 **end**

10 **for** $k = 1$ to K **do**

11 Sample minibatch $(s, a, r, b, s') \sim \mathcal{D}$;

12 Sample $n \sim \text{Uniform}(1 \dots N)$;

13 Update predictor network θ using:

$$\nabla_{\theta} \frac{2}{|B|} \sum (f_{\theta}(s, a) - f_{\theta_n}(s, a))$$

14 **end**

15 **for** $g = 1$ to G **do**

16 Sample minibatch $(s, a, r, b, s') \sim \mathcal{D}$;

17 Update actor using:

$$\nabla_{\phi} \frac{1}{|B|} \sum_{(s,a)} \left[\min_{i=1,2} Q_{\varphi_i}(s, \tilde{a}_{\phi}(s)) - \beta \log \pi_{\phi}(\tilde{a}_{\phi}(s)|s) - \lambda_{\text{actor}} b(s, \tilde{a}_{\phi}(s)) \right]$$

where $\tilde{a}_{\phi}(s) \sim \pi_{\phi}(\cdot|s)$ via reparameterization;

18 Update each critic Q_{φ_i} using:

$$\nabla_{\varphi_i} \frac{1}{|B|} \sum_{(s,a,r,s')} \left[Q_{\varphi_i}(s, a) - \left(r + \gamma E_{a' \sim \pi_{\phi}(\cdot|s')} \left[Q_{\varphi'_i}(s', a') - \beta \log \pi_{\phi}(a'|s') - \lambda_{\text{critic}} b(s', a') \right] \right)^2 \right]$$

19 Update target critics:

$$\varphi'_i \leftarrow (1 - \rho) \varphi'_i + \rho \varphi_i$$

20 **end**

21 **end**
