
Algorithm 1 Generic off-policy actor-critic

```

1: initialize  $\phi_0$ 
2: initialize  $\theta_0$ 
3: initialize replay buffer  $\mathcal{D} = \emptyset$  as a ring buffer of fixed size
4: initialize  $\mathbf{s} \sim d_0(\mathbf{s})$ 
5: for iteration  $k \in [0, \dots, K]$  do
6:   for step  $s \in [0, \dots, S-1]$  do
7:      $\mathbf{a} \sim \pi_{\theta_k}(\mathbf{a}|\mathbf{s})$  ▷ sample action from current policy
8:      $\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})$  ▷ sample next state from MDP
9:      $\mathcal{D} \leftarrow \mathcal{D} \cup \{(\mathbf{s}, \mathbf{a}, \mathbf{s}', r(\mathbf{s}, \mathbf{a}))\}$  ▷ append to buffer, purging old data if
       buffer too big
10:   end for
11:    $\phi_{k,0} \leftarrow \phi_k$ 
12:   for gradient step  $g \in [0, \dots, G_Q - 1]$  do
13:     sample batch  $B \subset \mathcal{D}$  ▷  $B = \{(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_t)\}$ 
14:     estimate error  $\mathcal{E}(B, \phi_{k,g}) = \sum_i (Q_{\phi_{k,g}}(\mathbf{s}_i, \mathbf{a}_i) - (r_i +$ 
        $\gamma \mathbb{E}_{\mathbf{a}' \sim \pi_k(\mathbf{a}'|\mathbf{s}_i)} Q_{\phi_k}(\mathbf{s}'_i, \mathbf{a}'))^2$ 
15:     update parameters:  $\phi_{k,g+1} \leftarrow \phi_{k,g} - \alpha_Q \nabla_{\phi_{k,g}} \mathcal{E}(B, \phi_{k,g})$ 
16:   end for
17:    $\phi_{k+1} \leftarrow \phi_{k,G_Q}$  ▷ update Q-function parameters
18:    $\theta_{k,0} \leftarrow \theta_k$ 
19:   for gradient step  $g \in [0, \dots, G_\pi - 1]$  do
20:     sample batch of states  $\{\mathbf{s}_i\}$  from  $\mathcal{D}$ 
21:     for each  $\mathbf{s}_i$ , sample  $\mathbf{a}_i \sim \pi_{\theta_{k,g}}(\mathbf{a}|\mathbf{s}_i)$  ▷ do not use actions in the
       buffer!
22:     for each  $(\mathbf{s}_i, \mathbf{a}_i)$ , compute  $\hat{A}(\mathbf{s}_i, \mathbf{a}_i) = Q_{\phi_{k+1}}(\mathbf{s}_i, \mathbf{a}_i) -$ 
        $\mathbb{E}_{\mathbf{a} \sim \pi_{k,g}(\mathbf{a}|\mathbf{s}_i)} [Q_{\phi_{k+1}}(\mathbf{s}_i, \mathbf{a})]$ 
23:      $\nabla_{\theta_{k,g}} J(\pi_{\theta_{k,g}}) \approx \frac{1}{N} \nabla_{\theta_{k,g}} \log \pi_{\theta_{k,g}}(\mathbf{s}_i, \mathbf{a}_i) \hat{A}(\mathbf{s}_i, \mathbf{a}_i)$ 
24:      $\theta_{k,g+1} \leftarrow \theta_{k,g} + \alpha_\pi \nabla_{\theta_{k,g}} J(\pi_{\theta_{k,g}})$ 
25:   end for
26:    $\theta_{k+1} \leftarrow \theta_{k,G_\pi}$  ▷ update policy parameters
27: end for

```
