

---

**Algorithm 1** A2C algorithm

---

**Require:** Initialize Actor-Critic network  $G$  with parameter  $\theta$ .

- 1: **for** each episode **do**
  - 2:     Get initial state  $s$
  - 3:     Initialize a storage buffer  $S, A, R, V, S'$
  - 4:     **for**  $i = 1, 2, 3.., N$  **do**
  - 5:         Sample an action  $a \sim G(s)_\pi$  and get the associated value  $v \leftarrow G(s)_v$
  - 6:         Run the action  $a$  through the environment, obtain the reward and  
next state  $r, s' \leftarrow ENV(s, a)$
  - 7:         Collect and store  $S, A, R, V, S' \leftarrow s, a, r, v, s'$
  - 8:          $s \leftarrow s'$
  - 9:     **end for**
  - 10:     Compute the discounted returns  $\hat{V} = \sum_{l=0}^{N-1} \gamma^l r_{t+l}$
  - 11:     Compute an advantage function  $\psi(V, R, S')$
  - 12:     Optimize  $\theta$  to minimize  $-\log(G(A \mid S)_\pi)\psi(V, R, S') + \lambda\|V - \hat{V}\|$
  - 13:     Empty  $S, A, R, V, S'$
-