

---

**Algorithm 1** Actor-Critic Algorithm(On-Policy)

---

**Require:**

- Learning rate for Actor  $\alpha_\theta$ ,
  - Learning rate for Critic  $\alpha_\omega$ ,
  - Initial parameters  $\theta, \omega$ ,
  - Number of iterations  $N$
- 1:  $\mathcal{D} \leftarrow \emptyset$
  - 2: **for** iteration = 1 to  $N$  **do**
  - 3:   Run policy  $\pi_\theta$  for  $T$  steps to collect data  $\mathcal{D}$ .
  - 4:    $\mathcal{L}_{\text{Critic}} \leftarrow \text{Comput Critic Loss}(\mathcal{D}, V_\omega)$
  - 5:    $\omega \leftarrow \omega + \alpha_\omega \nabla_\omega \mathcal{L}_{\text{Critic}}$
  - 6:    $\mathcal{L}_{\text{Actor}} \leftarrow \text{Compute Actor Loss}(\mathcal{D}, \pi_\theta)$
  - 7:    $\theta \leftarrow \theta + \alpha_\theta \nabla_\theta \mathcal{L}_{\text{Actor}}$
  - 8: **end for**
-