**Algorithm 1:** Deep Deterministic Policy Gradient (DDPG)

**Input:** Initial actor parameters $\theta$, critic parameters $\phi$, target actor parameters $\bar{\theta} \leftarrow \theta$, target critic parameters $\bar{\phi} \leftarrow \phi$, polyak averaging coefficient $\tau$, empty replay buffer $\mathcal{D}$, action bounds $a_{\min}$ and $a_{\max}$

**1** **for** *each iteration* **do**

**2**     Reset environment and observe initial state $s_0$;

**3**     Reset Ornstein-Uhlenbeck noise process $\mathcal{N}_t$;

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

**5**         Select action with exploration: $a_t = \pi_\theta(s_t) + \mathcal{N}_t$;

**6**         Clip action: $a_t \leftarrow \text{clip}(a_t, -c, c)$;

**7**         Execute $a_t$ in environment;

**8**         Observe reward $r_t$, next state $s_{t+1}$, and terminal signal $d_t$;

**9**         Store $(s_t, a_t, r_t, s_{t+1}, d_t)$ in replay buffer $\mathcal{D}$;

**10**         **if** $d_t$ *is True* **then**

**11**             Reset environment and observe new initial state $s_{t+1}$;

**12**         **end**

**13**         $s_t \leftarrow s_{t+1}$;

**14**     **end**

**15**     **for** *each gradient step* **do**

**16**         Sample a mini-batch of $N$ transitions $(s, a, r, s', d)$ from $\mathcal{D}$;

**17**         Compute target Q-value:

$$y = r + \gamma(1 - d)Q_{\bar{\phi}}(s', \pi_{\bar{\theta}}(s'))$$

**18**         Update critic by gradient descent:

$$\phi \leftarrow \phi - \lambda_Q \nabla_\phi \frac{1}{N} \sum \left(Q_\phi(s, a) - y\right)^2$$

**19**         Update actor by gradient ascent:

$$\theta \leftarrow \theta + \lambda_\pi \nabla_\theta \frac{1}{N} \sum Q_\phi(s, \pi_\theta(s))$$

**20**         Update target networks:

$$\bar{\phi} \leftarrow \tau\phi + (1 - \tau)\bar{\phi}, \quad \bar{\theta} \leftarrow \tau\theta + (1 - \tau)\bar{\theta}$$

**21**     **end**

**22** **end**