

TD3 Algorithm

- 1: Initialize critic networks $Q_{\theta_1}, Q_{\theta_2}$, and actor network π_ϕ with random parameters θ_1, θ_2, ϕ
 - 2: Initialize target networks $\theta'_1 \leftarrow \theta_1, \theta'_2 \leftarrow \theta_2, \phi' \leftarrow \phi$
 - 3: Initialize replay buffer B
 - 4: **for** $t = 1$ to T **do**
 - 5: Select action with exploration noise $a \sim \pi_\phi(s) + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma)$ and observe reward R and new state s'
 - 6: Store transition tuple (s, a, R, s') in B
 - 7: Sample mini-batch of N transitions (s, a, R, s') from B
 - 8: $\tilde{a} \leftarrow \pi_{\phi'}(s') + \epsilon, \epsilon \sim \text{clip}(\mathcal{N}(0, \tilde{\sigma}), -c, c)$
 - 9: $y \leftarrow R + \gamma \min_{i=1,2} Q_{\theta'_i}(s', \tilde{a})$
 - 10: Update critics $\theta_i \leftarrow \arg \min_{\theta_i} N^{-1} \sum (y - Q_{\theta_i}(s, a))^2$
 - 11: **if** $t \bmod d$ **then**
 - 12: Update ϕ by the deterministic policy gradient:
 - 13: $\nabla_\phi J(\phi) = N^{-1} \sum \nabla_a Q_{\theta_1}(s, a)|_{a=\pi_\phi(s)} \nabla_\phi \pi_\phi(s)$
 - 14: Update target networks:
 - 15: $\theta'_i \leftarrow \tau \theta_i + (1 - \tau) \theta'_i$
 - 16: $\phi' \leftarrow \tau \phi + (1 - \tau) \phi'$
 - 17: **end if**
 - 18: **end for**
-