

# DDPG Algorithm

- 1: **Initialize** critic network  $Q(s, a|\theta^Q)$  and actor  $\mu(s|\theta^\mu)$  with weights  $\theta^Q$  and  $\theta^\mu$
- 2: **Initialize** target networks  $Q'$  and  $\mu'$  with weights  $\theta^{Q'} \leftarrow \theta^Q$ ,  $\theta^{\mu'} \leftarrow \theta^\mu$
- 3: **Initialize** replay buffer  $R$
- 4: **for** episode = 1 to  $M$  **do**
- 5:     **Initialize** a random process  $N$  for action exploration
- 6:     **Receive** initial observation state  $s_1$
- 7:     **for**  $t = 1$  to  $T$  **do**
- 8:         **Select** action  $a_t = \mu(s_t|\theta^\mu) + N_t$  according to the current policy and exploration noise
- 9:         **Execute** action  $a_t$  and observe reward  $r_t$  and new state  $s_{t+1}$
- 10:         **Store** transition  $(s_t, a_t, r_t, s_{t+1})$  in  $R$
- 11:         **Sample** a random minibatch of  $N$  transitions  $(s_i, a_i, r_i, s_{i+1})$  from  $R$
- 12:         **Calculate** target Q-value:  $y_i = r_i + \gamma Q'(s_{i+1}, \mu'(s_{i+1}|\theta^{\mu'})|\theta^{Q'})$
- 13:         **Update** critic by minimizing the loss:  $L = \frac{1}{N} \sum_i (y_i - Q(s_i, a_i|\theta^Q))^2$
- 14:         **Update** the actor policy using the sampled policy gradient:

$$\nabla_{\theta^\mu} J \approx \frac{1}{N} \sum_i \nabla_a Q(s, a|\theta^Q) \Big|_{s=s_i, a=\mu(s_i)} \nabla_{\theta^\mu} \mu(s|\theta^\mu) \Big|_{s=s_i}$$

- 15:     **Update** the target networks:

$$\theta^{Q'} \leftarrow \tau \theta^Q + (1 - \tau) \theta^{Q'}$$

$$\theta^{\mu'} \leftarrow \tau \theta^\mu + (1 - \tau) \theta^{\mu'}$$

- 16:     **end for**
- 17: **end for**