

$$Q^\pi(s_t, a_t) = x \quad (1)$$

$$\nabla \quad (2)$$

$$\begin{aligned}\nabla_{\theta^\mu} J &\approx \mathbb{E}_{s_t \sim S} [\nabla_{\theta^\mu} Q(s, a \mid \theta^Q) \mid_{s=s_t, a=\pi(s_t \mid \theta_t^\mu)}] \\ &= \mathbb{E}_{s_t \sim S} [\nabla_a Q(s, a \mid \theta^Q) \mid_{s=s_t, a=\pi(s_t \mid \theta_t^\mu)}] \\ \nabla_{\theta^\mu} \pi(s \mid \theta^\mu) \mid_{s=s_t}\end{aligned}\tag{3}$$

$$L(\theta^E) = \mathbb{E} \left[(E(o_t | \theta_t^E) - a_t)^2 \right] \quad (4)$$

1 introduction

[illegible]

Algorithm 1 PACEE

(The encoder layers are included in actor network)
 Randomly initialize critic network and actor network with weights $\theta_Q, \theta_1^\mu, \theta_2^\mu, \dots, \theta_n^\mu$
 Randomly initialize experience network with weights θ^E
 Initialize target network $\theta^{Q'} \leftarrow \theta^Q, \theta_1^{\mu'} \leftarrow \theta_1^\mu, \theta_2^{\mu'} \leftarrow \theta_2^\mu, \dots, \theta_n^{\mu'} \leftarrow \theta_n^\mu$
 Randomly initialize experience network with weights θ^E
 Initialize replay buffer D_1, D_2, \dots, D_n, D'
 For $episode = 1, 2, \dots, M$ do:
 Initialize episode return $G = 0$
 Initialize an empty episode track Tr
 Receive initial state s_1
 For $t = 1, 2, \dots, T$ do:
 Judge the stage of the time step n
 Select action and get observation: $a_t, o_t = \pi(s_t, \theta_{t,n}^\mu)$
 Add positive guidance: $a_t = a_t + (E(o_t | \theta_{t,n}^E) - a_t) \xi$
 Execute action a_t and observe reward r_{t+1} and new state s_{t+1}
 Accumulate return $G = G + r_{t+1}$
 Store $(s_t, a_t, r_{t+1}, s_{t+1})$ and $(o_t, a_t, r_{t+1}, s_{t+1})$ in D_n and Tr respectively
 Sample a random minibatch of N transitions $(s_i, a_i, r_{i+1}, s_{i+1})$ from D_n
 Calculate: $q_{i+1} = Q(s_{i+1}, \pi(s_{i+1} | \theta_{t,n}^{\mu'}) | \theta_t^{Q'})$
 Add positive guidance: $q_{i+1} = q_{i+1} + (Q(s_{i+1}, E(o_{i+1} | \theta_{t,n}^E) \theta_t^{Q'}) - q_{i+1}) \xi$
 Set: $y_i = r_{i+1} + \gamma q_{i+1} \phi$
 Calculate gradients wrt θ^Q and update critic network:

$$d\theta^Q \leftarrow \frac{1}{N} \nabla_{\theta^Q} \sum_i \left[(Q(s_i, a_i | \theta_t^Q) - y_i)^2 \right]$$

 Calculate gradients wrt θ_n^μ and update actor network:

$$d\theta_n^\mu \leftarrow \frac{1}{N} \sum_i \left[\nabla_a Q(s, a | \theta_t^Q) \Big|_{s=s_i, a=\pi(s_i | \theta_{t,n}^\mu)} \nabla_{\theta_n^\mu} \pi(s | \theta_{t,n}^\mu) \Big|_{s=s_i} \right]$$

 Sample a random minibatch of N transitions $(o_k, a_k, r_{k+1}, s_{k+1})$ from D'
 Calculate gradients wrt θ^E and update experience network:

$$d\theta^E \leftarrow \frac{1}{N} \nabla_{\theta^E} \sum_k \left[(E(o_k | \theta_t^E) - a_k)^2 \right]$$

 Update the target network: $\theta_{t+1}^{Q'} \leftarrow \tau \theta_t^Q + (1 - \tau) \theta_t^{Q'}, \theta_{t+1,n}^{\mu'} \leftarrow \tau \theta_{t,n}^\mu + (1 - \tau) \theta_{t,n}^{\mu'}$
 End for
 If $G \geq \bar{R}_K$, then store Tr into D'
 End for
