# 10 - Policy Iteration & Policy Gradients

## 知识点 & 题目

### Policy Evaluation

**Input:** $\pi$ the policy for evaluation, $V^\pi$ value function, and MDP $M = \langle S, s_0, A, P_a(s' \mid s), r(s, a, s') \rangle$
**Output:** Value function $V^\pi$

Repeat
    $\Delta \leftarrow 0$
    For each $s \in S$
        $\underbrace{V'^\pi(s) \leftarrow \sum_{s' \in S} P_{\pi(s)}(s' \mid s) \left[ r(s, a, s') + \gamma V^\pi(s') \right]}_{\text{Policy evaluation equation}}$
        $\Delta \leftarrow \max(\Delta, |V'^\pi(s) - V^\pi(s)|)$
    $V^\pi \leftarrow V'^\pi$
Until $\Delta \leq \theta$

### Policy Improvement

$$Q^\pi(s, a) = \sum_{s' \in S} P_a(s' \mid s) \left[ r + \gamma V^\pi(s') \right]$$

$$\text{If } Q^\pi(s, a) > Q^\pi(s, \pi(s))$$

$$\pi(s) \leftarrow a$$

### Policy Iteration

**Input:** MDP $M = \langle S, s_0, A, P_a(s' \mid s), r(s, a, s') \rangle$
**Output:** Policy $\pi$

Set $V^\pi$ to arbitrary value function; e.g., $V^\pi(s) = 0$ for all $s$.

Set $\pi$ to arbitrary policy; e.g. $\pi(s) = a$ for all $s$, where $a \in A$ is an arbitrary action.

Repeat $\quad O\left(|A|^{|S|}\right) \qquad\qquad O\left(|S|^3\right)$

$\qquad$ Compute $V^\pi(s)$ for all $s$ using policy evaluation
$\qquad$ For each $s \in S \quad O\left(|S|^2 \cdot |A|\right)$
$\qquad\qquad \pi(s) \leftarrow \text{argmax}_{a \in A(s)} Q^\pi(s, a)$
Until $\pi$ does not change

## POLICY ITERATION: EXAMPLE

Assume $\gamma = 0.9$

$$V^\pi(s) = \sum_{s' \in S} P_{\pi(s)}(s' \mid s) \left[ r(s, a, s') + \gamma V^\pi(s') \right]$$

$\pi(2,2) = \text{up}$

$V^\pi(3,2) = 1$

$V^\pi(2,2) = \begin{array}{l} 0.8 \\ 0.1 \\ 0.1 \end{array} \begin{array}{l} [0 \quad + \quad 0.9 \cdot 0] \, \nwarrow \\ [0 \quad + \quad 0.9 \cdot 0] \quad = 0.09 \\ [0 \quad + \quad 0.9 \cdot 1] \end{array}$

$Q^\pi(2,2, \text{Right}) = \begin{array}{l} 0.8 \\ 0.1 \\ 0.1 \end{array} \begin{array}{l} [0 + 0.9 \cdot 1] \, \nwarrow \\ [c + 0.9 \cdot 0] \, + \\ [c + 0.9 \cdot c] \end{array} = 0.72$

$\pi(2,2) \leftarrow \text{Right}$



## Policy Gradients

**Algorithm – REINFORCE**

**Input:** A differentiable policy $\pi_\theta(s, a)$, an MDP $M = \langle S, s_0, A, P_a(s' \mid s), r(s, a, s') \rangle$
**Output:** Policy $\pi_\theta(s, a)$

Repeat
$\qquad$ Generate episode $(s_0, a_0, r_1, \ldots s_{T-1}, a_{T-1}, r_T)$ by following $\pi_\theta$
$\qquad$ For each $(s_t, a_t)$ in the episode
$\qquad\qquad G \leftarrow \sum_{k=t+1}^{T} \gamma^{k-t-1} r_k$
$\qquad\qquad \theta \leftarrow \theta + \alpha \gamma^t G \nabla \ln \pi_\theta(s, a)$
Until some time limit or until $\pi_\theta$ converges

- Continuous action space
- G: estimate of Q(s,a)
  - Instability
- Generally converge to local optima

## Q Actor Critic

**Input:** An MDP $M = \langle S, s_0, A, P_a(s' \mid s), r(s, a, s') \rangle$

**Input:** A differentiable actor policy $\pi_\theta(s, a)$

**Input:** A differentiable critic Q-function $Q(s, a)$

**Output:** Policy $\pi_\theta(s, a)$

Initialise actor $\pi$ parameters $\theta$ and critic parameters $w$ arbitrarily

Repeat (for each episode)

    $s \leftarrow$ the first state in episode $e$

    Select action $a \sim \pi_\theta(s)$

    Repeat (for each step in episode $e$)

        Execute action $a$ in state $s$

        Observe reward $r$ and new state $s'$

        Select action $a' \sim \pi_\theta(s')$

        $\delta \leftarrow r + \gamma \cdot Q_w(s', a') - Q_w(s', a')$

        $w \leftarrow w + \alpha_w \cdot \delta \cdot \nabla Q_w(s, a)$

        $\theta \leftarrow \theta + \alpha_\theta \cdot \delta \cdot \nabla \ln \pi_\theta(s, a)$

        $s \leftarrow s'; a \leftarrow a'$

- Replace G with TD estimate -> more stable -> converge
- Critic: feedback on actions

# 题目