# Reinforcement Learning

## Learning and Value-based Methods

- An RL agent may include one or more of these components
    - Policy : agent's behavior func:on
    - Value func : how good is a state and/or ac:on
    - Model: agent's representa:on of the environment
- Markov Decision Processes $$ \langle\mathcal{S}, \mathscr{A}, \mathscr{P}, \mathscr{R}, \gamma\rangle $$ ： RL 的学习环境
    - $ \mathcal{S} $ -a set of states
    - $\mathscr{A}$- a set of ac:ons
    - $\mathscr{P}$ - transi:on probability function, $\mathscr{P}_{s s^{\prime}}^{a}=\mathbb{P}\left[S_{t+1}=s^{\prime} \mid S_{t}=s, A_{t}=a\right]$
    - $\mathscr{R}$ - reward function, $\mathscr{R}_{s}^{a}=\mathbb{E}\left[R_{t+1} \mid S_{t}=s, A_{t}=a\right]$
    - $\gamma$ - discounting factor for future reward
    - Polity $\pi$ : $\pi(a \mid s)=\mathbb{P}\left[A_{t}=a \mid S_{t}=s\right]$
- **Value Function** :

    $\begin{aligned} v_{\pi}(s) &=\mathbb{E}_{\pi}\left[G_{t} \mid S_{t}=s\right] \\ &=\mathbb{E}_{\pi}\left[R_{t+1}+\gamma R_{t+2}+\gamma^{2} R_{t+3}+\ldots \mid S_{t}=s\right] \\ &=\mathbb{E}_{\pi}\left[R_{t+1}+\gamma\left(R_{t+2}+\gamma R_{t+3}+\ldots\right) \mid S_{t}=s\right] \\ &=\mathbb{E}_{\pi}\left[R_{t+1}+\gamma G_{t+1} \mid S_{t}=s\right] \\ &=\mathbb{E}_{\pi}\left[R_{t+1}+\gamma v_{\pi}\left(S_{t+1}\right) \mid S_{t}=s\right] \end{aligned}$

- $v_{\pi}(s)=\sum_{a \in \mathscr{A}} \pi(a \mid s) q_{\pi}(s, a)$

- $q_{\pi}(s, a)=\mathscr{R}_{s}^{a}+\gamma \sum_{s^{\prime} \in \mathcal{S}} \mathscr{P}_{s s^{\prime}}^{\alpha} v_{\pi}\left(s^{\prime}\right)$

- $v_{\pi}(s)=\sum_{a \in \mathscr{A}} \pi(a \mid s)\left(\mathscr{R}_{s}^{a}+\gamma \sum_{s^{\prime} \in \mathcal{S}} \mathscr{P}_{s s^{\prime}}^{a} v_{\pi}\left(s^{\prime}\right)\right)$

- $q_{\pi}(s, a)=\mathscr{R}_{s}^{a}+\gamma \sum_{s^{\prime} \in \delta} \mathscr{P}_{s s^{\prime}}^{a} \sum_{a^{\prime} \in \mathscr{A}} \pi\left(a^{\prime} \mid s^{\prime}\right) q_{\pi}\left(s^{\prime}, a^{\prime}\right)$

- **Optimal Value Function**
    - $v_{*}(s)=\max _{\pi} v_{\pi}(s)$
    - $q_{*}(s, a)=\max _{\pi} q_{\pi}(s, a)$
    - $v_{*}(s)=\max _{a} q_{*}(s, a)$
    - $q_{*}(s, a)=\mathscr{R}_{s}^{a}+\gamma \sum_{s^{\prime} \in \delta} \mathscr{P}_{s s^{\prime}}^{a} v_{*}\left(s^{\prime}\right)$
    - $v_{*}(s)=\max _{a}\left(\mathscr{R}_{s}^{a}+\gamma \sum_{s^{\prime} \in \delta} \mathscr{P}_{s s^{\prime}}^{a} v_{*}\left(s^{\prime}\right)\right)$

- $q_{*}(s, a)=\mathscr{R}_{s}^{a}+\gamma \sum_{s^{\prime} \in \mathcal{S}} \mathscr{P}_{s s^{\prime}}^{a} \max _{a^{\prime}} q_{*}\left(s^{\prime}, a^{\prime}\right)$

- Dynamic Programming：假设我们已知 MDP

  - Policy Evaluation：利用迭代法计算给定 policy $\pi$的 value function
    $v^{k+1}=\mathscr{R}^{\pi}+\gamma \mathscr{P}^{\pi} v^{k}$

  - Policy Improvement：利用贪心法

    $\pi^{\prime}(s)=\underset{a \in \mathscr{A}}{\operatorname{argmax}} q_{\pi}(s, a)$

    $q_{\pi}\left(s, \pi^{\prime}(s)\right)=\max _{a \in \mathscr{A}} q_{\pi}(s, a) \geq q_{\pi}(s, \pi(s))=v_{\pi}(s)$

  - value iteration: update policy every iteration
    $v_{*}(s) \leftarrow \max _{a \in \mathscr{A}}\left(\mathscr{R}_{s}^{a}+\gamma \sum_{s^{\prime} \in \mathcal{S}} \mathscr{P}_{s s^{\prime}}^{a} v_{*}\left(s^{\prime}\right)\right)$

  - 可以证明，这样迭代下去，最后 value function 会以$\gamma$的比例线性收敛到最优解

  - 每个迭代复杂度是$O(m n ^2)$，m 是action 数目，n 是 state 数目

- Monte-Carlo Methods

  - 给定