

Markov Decision Process (MDP)

$$M = (S, A, T, R, \gamma)$$

state \uparrow action \uparrow reward $R(s, a, s')$
 discount γ
 transition $T(s, a, s') = \Pr[s' | s, a]$
 total: $\sum_t \gamma^t r_t$
 max: $\frac{r_{\max}}{1-\gamma}$

Value function $V_t^\pi(s)$:

expected total reward with ① policy π ② begin in s ③ have T steps remaining

- Finite horizon optimal: $V_t^*(s) = \max_{a \in A} [R(s, a) + \gamma \sum_{s'} P(s' | s, a) V_{t-1}^*(s')]$.
 $V_0^* = 0$

- Infinite horizon optimal: $V^*(s) = \max_{a \in A} [R(s, a) + \gamma \sum_{s'} P(s' | s, a) V^*(s')]$.

Policy iteration:

① Policy evaluation: $\begin{cases} \text{Value iteration} = V_t^\pi(s) \leftarrow R(s, \pi(s)) + \gamma \sum_{s'} P(s' | s, \pi(s)) V_{t-1}^\pi(s) \\ V_0^\pi(s) \leftarrow 0 \\ \text{Solve linear system (Bellman equation)} \end{cases}$

② Policy improvement: For state s , $\pi(s) \leftarrow \operatorname{argmax}_{a \in A} [R(s, a) + \gamma \sum_{s'} P(s' | s, a) V^\pi(s')]$.

Reinforcement Learning Transition T / Reward R unknown!

Monte-Carlo method:

① Monte-Carlo prediction $\begin{cases} \text{Average} = \text{average}(G_t^\pi(s)) \quad \leftarrow \begin{matrix} \text{first-visit} \\ \text{every-visit} \end{matrix} \\ \text{Incremental} = V^\pi(s_t) \leftarrow V^\pi(s_t) + \alpha [G_t^\pi - V^\pi(s_t)] \end{cases}$

- Functional approximation: $\theta_i \leftarrow \theta_i + \alpha [U_j(s) - \hat{V}_\theta(s)] \frac{\partial \hat{V}_\theta(s)}{\partial \theta_i}$.

② Monte-Carlo control $\begin{cases} \text{Naive: } \pi(s) \leftarrow \operatorname{argmax}_a Q(s, a) \text{ in each episode} \\ \text{on-policy } \leftarrow \begin{cases} \epsilon\text{-soft: } \pi(a|s) \leftarrow \begin{cases} 1-\epsilon + \epsilon/|A| & \text{if } a = a^* \\ \epsilon/|A(s)| & \text{if } a \neq a^* \end{cases} \end{cases} \end{cases}$

Temporal difference method

① TD prediction: $V(s_t) \leftarrow V(s_t) + \alpha [R_t + \gamma V(s_{t+1}) - V(s_t)]$.

- Functional approximation: $\theta_i \leftarrow \theta_i + \alpha [R_t + \gamma \hat{V}_\theta(s') - \hat{V}_\theta(s)] \frac{\partial \hat{V}_\theta(s)}{\partial \theta_i}$.

② TD control $\begin{cases} \text{SARSA: } Q(s, A) \leftarrow Q(s, A) + \alpha [R + \gamma Q(s', A) - Q(s, A)] \\ \text{Q-learning: } Q(s, A) \leftarrow Q(s, A) + \alpha [R + \gamma \max_{a'} Q(s', a') - Q(s, A)] \end{cases}$

- Functional approximation: $\theta_i \leftarrow \theta_i + \alpha [R + \gamma \max_{a'} \hat{Q}(s', a') - \hat{Q}(s, a)] \frac{\partial \hat{Q}(s, a)}{\partial \theta_i}$

Monte-Carlo policy gradient: $\theta_{j+1} \leftarrow \theta_j + \alpha U_j \nabla_\theta \ln \pi_\theta(s, a_j)$

Partially Observable Markov Decision Process (POMDP) $M = (S, A, E, T, O, R, \gamma)$

Belief state $b(s)$

evidences/observations \downarrow
 observation function $O(s, e) = P(e|s)$

- Update: $b'(s') = \alpha \Pr[e' | s'] \sum_s \Pr[s' | s, a] b(s)$.

Value iteration: $\alpha_p(s) = \sum_{s'} P_r[s' | s, a] [R(s, a, s') + \gamma \sum_{e'} \Pr[e' | s'] \alpha_{p, e'}(s')]$