

CS 839: PROBABILITY & LEARNING IN HIGH DIMENSION
LECTURE 18: MARKOV DECISION PROCESSES

Yudong Chen
UW-Madison CS
April 6, 2022

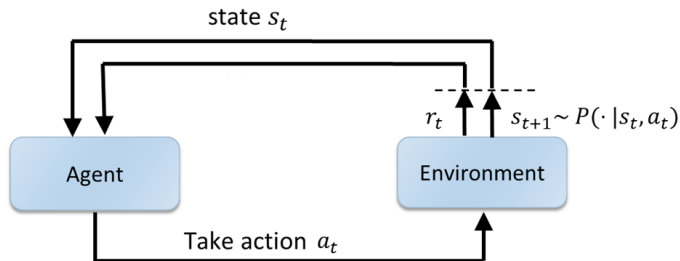
- Section 1.1–1.3.3 in Reinforcement Learning: Theory and Algorithms, by Alekh Agarwal, Nan Jiang, Sham M. Kakade, Wen Sun, 2021.
(“**AJKS**” book)
<https://rltheorybook.github.io>

Discounted Markov Decision Process (MDP)

- Interaction between agent and environment is described by an MDP $M = (\mathcal{S}, \mathcal{A}, \mathbb{P}, r, \gamma)$.
- \mathcal{S} : *state space*; finite
- \mathcal{A} : *action space*; finite
- $\mathbb{P} : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ is the *transition kernel*, where $\Delta(\mathcal{S}) :=$ space of probability distributions over \mathcal{S} .
 - $\mathbb{P}(s'|s, a) =$ probability of transitioning to next state s' given the current state is s and action a is taken
- $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the *reward function*
 - $r(s, a) =$ one-step (deterministic) reward when current state is s and action is a .
- $\gamma \in [0, 1)$: *discount factor*.

Policy and MDP Dynamics

- Agent adopts a *stochastic*/randomized policy $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$
 - $\pi(a|s)$ = probability of taking action a at state s
- Initial state $s_0 \sim \mu \in \Delta(\mathcal{S})$
- At each time t , agent observes state s_t , takes action $a_t \sim \pi(\cdot|s_t)$, and receives reward $r_t = r(s_t, a_t)$. System then transitions to state $s_{t+1} \sim \mathbb{P}(\cdot|s_t, a_t)$ at time $t + 1$



Value and Q-Functions

- Fix a policy π
- The **value function** $V^\pi : \mathcal{S} \rightarrow \mathbb{R}$ is defined to be

$$V^\pi(s) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s \right]$$

- The **action-value (Q-value) function** $Q^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is defined as

$$Q^\pi(s, a) := \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, a_0 = a \right]$$

- Relationship:

$$V^\pi(s) := \mathbb{E}_{a \sim \pi(\cdot \mid s)} [Q^\pi(s, a)] \tag{1}$$

- The objective is to find a policy π that maximizes $V^\pi(s)$

Optimal Value Function and Policy

- Let Π be the set of all randomized policies. Define the **optimal value and Q functions** as

$$V^*(s) := \sup_{\pi \in \Pi} V^\pi(s), \quad \text{for each } s \in \mathcal{S},$$

$$Q^*(s, a) := \sup_{\pi \in \Pi} Q^\pi(s, a), \quad \text{for each } s \in \mathcal{S}, a \in \mathcal{A}.$$

THEOREM 1 (PUTERMAN '94, THM 6.2.7)

Assume \mathcal{S} and \mathcal{A} are finite. There exists a (deterministic) policy π^* such that

$$V^{\pi^*}(s) = V^*(s), \quad \forall s \in \mathcal{S},$$

$$Q^{\pi^*}(s, a) = Q^*(s, a), \quad \forall s \in \mathcal{S}, a \in \mathcal{A}$$

π^* is called an **optimal policy**.

- π^* maximizes $V^\pi(s)$ *simultaneously* for all $s \in \mathcal{S}$.

Bellman Equations

THEOREM 2 (BELLMAN EQUATIONS FOR π)

For each policy $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$, V^π and Q^π are the unique functions that satisfy

$$V^\pi(s) = \mathbb{E}_{a \sim \pi(\cdot|s), s' \sim \mathbb{P}(\cdot|s,a)} [r(s, a) + \gamma V^\pi(s')] \quad (2)$$

$$Q^\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim \mathbb{P}(\cdot|s,a), a' \sim \pi(\cdot|s')} [Q^\pi(s', a')]. \quad (3)$$

THEOREM 3 (BELLMAN OPTIMALITY EQUATIONS)

V^* and Q^* are the unique functions that satisfy

$$V^*(s) = \max_{a \in \mathcal{A}} \left[r(s, a) + \gamma \mathbb{E}_{s' \sim \mathbb{P}(\cdot|s,a)} V^*(s') \right], \quad \forall s \in \mathcal{S}. \quad (4)$$

$$Q^*(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim \mathbb{P}(\cdot|s,a)} \left[\max_{a' \in \mathcal{A}} Q^*(s', a') \right], \quad \forall s \in \mathcal{S}, a \in \mathcal{A}. \quad (5)$$

Value Iteration

- When r and \mathbb{P} are known, Q^* can be computed using **Value Iteration**:
 - Specify an initial function $Q^{(0)} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$
 - For $k = 0, 1, \dots$, compute

$$Q^{(k+1)}(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim \mathbb{P}(\cdot | s, a)} \left[\max_{a' \in \mathcal{A}} Q^{(k)}(s', a') \right], \quad \text{for each } s, a. \quad (6)$$

- Note: Q^* is a fixed point of (6) by Bellman Optimality Equation
- Value Iteration converges to Q^* geometrically:

$$\|Q^{(k+1)} - Q^*\|_{\infty} \leq \gamma^k \|Q^{(0)} - Q^*\|_{\infty}.$$

- Given Q^* , an optimal deterministic policy can be computed:

$$\pi^*(s) = \operatorname{argmax}_{a \in \mathcal{A}} Q^*(s, a), \quad s \in \mathcal{S}.$$

Finite-Horizon MDP

- $M = (\mathcal{S}, \mathcal{A}, \mathbb{P}, r, \mathbf{H})$.
- \mathcal{S} : state space; \mathcal{A} : action space
- \mathbf{H} : **horizon** (number of steps in an **episode**)
- $\mathbb{P} = (\mathbb{P}_h : h = 1, \dots, H)$ are the *transition kernels*
 - $\mathbb{P}_h(s'|s, a) =$ probability of next state s' when current state is s and action a is taken **at step h**
- $r = (r_h : h = 1, \dots, H)$ are the *reward functions*
 - $r_h(s, a) =$ one-step reward when current state-action is (s, a) **at step h**
- Policy $\pi = (\pi_h : h = 1, \dots, H)$
 - $\pi_h(a|s) =$ probability of taking action a at state s and **step h**
- Value and Q functions: for $h = 1, \dots, H$

$$V_h^\pi(s) := \mathbb{E}_\pi \left[\sum_{t=h}^H r_t(s_t, a_t) | s_h = s \right], \quad Q_h^\pi(s, a) := \mathbb{E}_\pi \left[\sum_{t=h}^H r_t(s_t, a_t) | s_h = s, a_h = a \right]$$

- π^*, V_h^*, Q_h^* defined analogously