

COMP3211 Tutorial 6: Markov Decision Process

Fengming ZHU

February 17, 2025

Department of CSE
HKUST

© 2024 Fengming Zhu. All rights reserved.

MDP V.S. Search

Value Functions

Bellman Expectation Equation

Bellman Optimality Equation

MDP V.S. Search

Search:

- A set of states \mathcal{S} , initial state I , goal state G
- A set of actions \mathcal{A}
- Deterministic transitions $T : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$
- cost function $c : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$
- Objective: a path p from I to G that minimizes $c(p)$

MDP V.S. Search

Search:

- A set of states \mathcal{S} , initial state I , goal state G
- A set of actions \mathcal{A}
- Deterministic transitions $T : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$
- cost function $c : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$
- Objective: a path p from I to G that minimizes $c(p)$

MDP:

- A set of states \mathcal{S} , a terminating condition $End(s)$
- A set of actions \mathcal{A}
- Stochastic transitions $T : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$
- Reward function $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, with a discount factor γ
- Objective: maximize $\sum_t \gamma^t r_t$

Solution Concept: Policy

Question:

To make sure you can come up with an optimal solution, you'd like to know:

- (A) Only your current state
- (B) All the history from the beginning up until now
- (C) Need to know more

Solution Concept: Policy

Question:

To make sure you can come up with an optimal solution, you'd like to know:

- (A) Only your current state
- (B) All the history from the beginning up until now
- (C) Need to know more

Markov property:

A state S_t is Markovian iff $P[S_{t+1}|S_t, \dots, S_0] = P[S_{t+1}|S_t]$. That is, your current state is already a “sufficient statistic”, also known as the **information state**.

Solution Concept: Policy

Question:

To make sure you can come up with an optimal solution, you'd like to know:

- (A) Only your current state
- (B) All the history from the beginning up until now
- (C) Need to know more

Markov property:

A state S_t is Markovian iff $P[S_{t+1}|S_t, \dots, S_0] = P[S_{t+1}|S_t]$. That is, your current state is already a “sufficient statistic”, also known as the **information state**.

Policy:

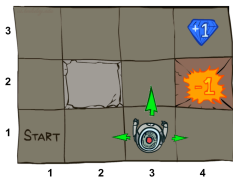
A solution is a policy $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$

Follow-up Question: Maze

Question:

Given a large maze, you (with deterministic actions U/D/L/R) are supposed to find a nice way from the entrance to the exit, which agent you'd like to choose

- (A) State machines with infinite memory
- (B) Agents that can A^* search
- (C) Agents that can compute policies
- (D) None of them

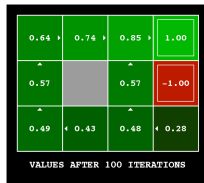
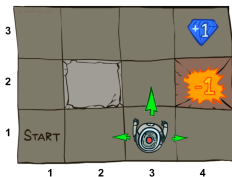


Follow-up Question: Maze

Question:

Given a large maze, you (with deterministic actions U/D/L/R) are supposed to find a nice way from the entrance to the exit, which agent you'd like to choose

- (A) State machines with infinite memory
- (B) Agents that can A^* search
- (C) Agents that can compute policies
- (D) None of them



Value Functions

Notations with Time Index

- **Transition:** $T_{s,s'}^a = \mathbb{P}[S_{t+1} = s' | S_t = s, A_t = a]$
- **Reward:** $R_s^a = \mathbb{E}[R_{t+1} | S_t = s, A_t = a]$
- **Stationary policy:** $\pi(a|s) = \mathbb{P}[A_t = a | S_t = s]$
- **Return:** The return G_t is the total discounted reward from time t ,

$$G_t = R_{t+1} + \gamma R_{t+2} + \cdots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

State-value function:

The state-value function $v_\pi(s)$ for an MDP is the expected return starting from state s , and then following policy π ,

$$v_\pi(s) = \mathbb{E}_\pi[G_t | S_t = s]$$

Action-value function:

The action-value function $q_\pi(s, a)$ for an MDP is the expected return starting from state s , taking action a , and then following policy π ,

$$q_\pi(s, a) = \mathbb{E}_\pi[G_t | S_t = s, A_t = a]$$

Policy iteration:

- Initialize $\pi \leftarrow \pi_0$
- While π still changing:
 - Policy evaluation: iterate until convergence
$$\forall s, v_{\pi}^t(s) = \sum_{s' \in \mathcal{S}} T(s, \pi(s), s') [R(s, \pi(a), s') + \gamma v_{\pi}^{t-1}(s')]$$
 - Policy improvement: greedy update
$$\forall s, \pi_{new}(s) = \arg \max_{a \in A(s)} q_{\pi}(s, a)$$
- $\pi \leftarrow \pi_{new}$

Bellman Expectation Equation

Bellman Expectation Equation

- For policy evaluation, we desire

$$v_{\pi}(s) = \sum_{s' \in S} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma v_{\pi}(s')]$$

- For state-value function,

$$v_{\pi}(s) = E_{\pi}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) | S_t = s]$$

- For action-value function,

$$q_{\pi}(s, a) = E_{\pi}[R_{t+1} + \gamma q_{\pi}(S_{t+1}, A_{t+1}) | S_t = s, A_t = a]$$

Bellman Expectation Equation

Warming-up: Adam's Law

Adam's Law:

For any random variables X and Y ,

$$E[E[X|Y]] = E[X]$$

Adam's Law with Extra Conditioning:

For any random variables X , Y and Z ,

$$E[E[X|Y, Z]|Z] = E[X|Z]$$

Bellman Expectation Equation for $v^\pi(s)$

We first prove the Bellman equation for state-value function.

$$\begin{aligned}v_\pi(S_t = s) &= E_\pi[G_t | S_t = s] \\&= E_\pi[R_{t+1} + \gamma(R_{t+2} + \gamma R_{t+3} + \cdots) | S_t = s] \\&= E_\pi[R_{t+1} + \gamma G_{t+1} | S_t = s]\end{aligned}$$

Since

$$\begin{aligned}E_\pi[G_{t+1} | S_t] &= E_\pi[E_\pi[G_{t+1} | S_{t+1}, S_t] | S_t] \\&= E_\pi[E_\pi[G_{t+1} | S_{t+1}] | S_t] \\&= E_\pi[v_\pi(S_{t+1}) | S_t]\end{aligned}$$

Thus,

$$\begin{aligned}v(S_t = s) &= E_\pi[G_t | S_t = s] \\&= E_\pi[R_{t+1} + \gamma G_{t+1} | S_t = s] \\&= E_\pi[R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s]\end{aligned}$$

Bellman Expectation Equation for $q^\pi(s, a)$

We then prove the Bellman equation for action-state function.

$$\begin{aligned}q_\pi(S_t = s, A_t = a) &= E_\pi[G_t | S_t = s, A_t = a] \\&= E_\pi[R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a]\end{aligned}$$

Since

$$\begin{aligned}E[G_{t+1} | S_t, A_t] &= E[E[G_{t+1} | (S_{t+1}, A_{t+1}), (S_t, A_t)] | (S_t, A_t)] \\&= E[E[G_{t+1} | (S_{t+1}, A_{t+1})] | (S_t, A_t)]\end{aligned}$$

Under policy π , we have

$$\begin{aligned}E_\pi[G_{t+1} | S_t, A_t] &= E_\pi[E_\pi[G_{t+1} | (S_{t+1}, A_{t+1})] | (S_t, A_t)] \\&= E_\pi[q_\pi(S_{t+1}, A_{t+1}) | (S_t, A_t)]\end{aligned}$$

Thus,

$$\begin{aligned}q_\pi(S_t = s, A_t = a) &= E_\pi[R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a] \\&= E_\pi[R_{t+1} + \gamma q_\pi(S_{t+1}, A_{t+1}) | S_t = s, A_t = a]\end{aligned}$$

Bellman Optimality Equation

Optimal Value Function

- For state-value function,

$$v_*(s) = \max_{\pi} v_{\pi}(s)$$

- For action-value function,

$$q_*(s, a) = \max_{\pi} q_{\pi}(s, a)$$

- The optimal value function specifies the best possible performance in the MDP.
- An MDP is “solved” once we know the optimal values.

Define a partial ordering over policies:

$$\pi \geq \pi', \text{ if } v_{\pi}(s) \geq v_{\pi'}(s), \text{ for all } s.$$

Define a partial ordering over policies:

$$\pi \geq \pi', \text{ if } v_{\pi}(s) \geq v_{\pi'}(s), \text{ for all } s.$$

For any MDP:

- There exists an optimal policy π_* that is better than or equal to all other policies, $\pi_* \geq \pi$, for all π .
- All optimal policies achieve the optimal state-value function, $v_{\pi_*}(s) = v_*(s)$, for all s .
- All optimal policies achieve the optimal action-value function, $q_{\pi_*}(s, a) = q_*(s, a)$, for all s, a .

Finding Optimal Policy

Theorem:

- An optimal policy can be found by maximizing over $q_*(s, a)$,

$$\pi_*(a|s) = \begin{cases} 1 & , \text{ if } a = \operatorname{argmax}_{a \in A} q_*(s, a) \\ 0 & , \text{ otherwise} \end{cases}$$

- There is always a **deterministic** optimal policy for any MDP.
- If we know $q_*(s, a)$, we immediately have the optimal policy.

Bellman Optimality Equation

- For state-value function,

$$v_*(s) = \max_a E[R_{t+1} + \gamma v_*(S_{t+1}) | S_t = s, A_t = a]$$

- For action-value function,

$$q_*(s, a) = E[R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') | S_t = s, A_t = a]$$

- However, non-linear thus no closed form solution in general.

Bellman Optimality Equation for v_* and q_*

Under the optimal policy, we first show the relation of v^* and q^* .

- v_* in terms of q_* ,

$$v_*(s) = \max_a q_*(s, a)$$

Bellman Optimality Equation for v_* and q_* (cont'd)

- q_* in terms of v_* ,

$$\begin{aligned} q_*(s, a) &= \max_{\pi} q_{\pi}(s, a) \\ &= R_s^a + \gamma \sum_{s'} T_{ss'}^a \max_{\pi} v_{\pi}(s') \\ &= R_s^a + \gamma \sum_{s'} T_{ss'}^a v_*(s') \\ &= E[R_{t+1} | S_t = s, A_t = a] \\ &\quad + \gamma \sum_{s'} \{ P(S_{t+1} = s' | S_t = s, A_t = a) \\ &\quad \cdot E[v_*(S_{t+1}) | S_{t+1} = s', S_t = s, A_t = a] \} \\ &= E[R_{t+1} | S_t = s, A_t = a] + \gamma E[v_*(S_{t+1}) | S_t = s, A_t = a] \\ &= E[R_{t+1} + \gamma v_*(S_{t+1}) | S_t = s, A_t = a] \end{aligned}$$

Bellman Optimality Equation for V_* and Q_*

Then we show the Bellman optimal equation for v^* ,

$$\begin{aligned} v_*(s) &= \max_a q_*(s, a) \\ &= \max_a E[R_{t+1} + \gamma v_*(S_{t+1}) | S_t = s, A_t = a] \end{aligned}$$

Finally, we show the Bellman equation for q^* ,

$$\begin{aligned} q_*(s, a) &= E[R_{t+1} + \gamma v_*(S_{t+1}) | S_t = s, A_t = a] \\ &= E[R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') | S_t = s, A_t = a] \end{aligned}$$

One Last Remark on PI v.s. VI

Thanks!