

Planning in MDPs

Sham Kakade and Kianté Brantley

CS 2824: Foundations of Reinforcement Learning

Announcements

HW0 is **due** Mon Feb. 2nd

First reading assignment **due** Wed. Feb 4th

Waitlist

Recap: Infinite Horizon MDPs

$$\mathcal{M} = \{S, A, P, r, \gamma\}$$

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \rightarrow [0,1], \quad \gamma \in [0,1)$$

Stationary Policy $\pi : S \mapsto \Delta(A)$

$$\text{Value function } V^\pi(s) = \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid s_0 = s, a_h \sim \pi(s_h), s_{h+1} \sim P(\cdot \mid s_h, a_h) \right]$$

$$\text{Q function } Q^\pi(s, a) = \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid (s_0, a_0) = (s, a), a_h \sim \pi(s_h), s_{h+1} \sim P(\cdot \mid s_h, a_h) \right]$$

Recap: Bellman Optimality

$$\mathcal{M} = \{S, A, P, r, \gamma\}$$

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \rightarrow [0,1], \quad \gamma \in [0,1)$$

Theorem 1: Bellman Optimality (Q-version)

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

Main Question for Today:

Given an MDP $\mathcal{M} = (S, A, P, r, \gamma)$, How to find π^\star (stationary & deterministic)

Outline

1. Bellman optimality — property of V^\star
2. Optimal planning: Value Iteration

Bellman Optimality

Theorem 2:

For any $V : S \rightarrow \mathbb{R}$, if $V(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V(s') \right]$ for all s ,
then $V(s) = V^*(s), \forall s$

$$\begin{aligned} |V(s) - V^*(s)| &= \left| \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s')) - \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s')) \right| \\ &\leq \max_a \left| (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s')) - (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s')) \right| \\ &\leq \max_a \gamma \mathbb{E}_{s' \sim P(s, a)} |V(s') - V^*(s')| \\ &\leq \max_a \gamma \mathbb{E}_{s' \sim P(s, a)} \left(\max_{a'} \gamma \mathbb{E}_{s'' \sim P(s', a')} |V(s'') - V^*(s'')| \right) \\ &\leq \max_{a_1, a_2, \dots, a_{k-1}} \gamma^k \mathbb{E}_{s_k} |V(s_k) - V^*(s_k)| \end{aligned}$$

Bellman Optimality for Q^\star

What about Q^\star ?

We should have:

For any $Q : S \times A \rightarrow \mathbb{R}$, if $Q(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a'} Q(s', a')$
for all s , then $Q(s, a) = Q^\star(s, a), \forall s, a$

Outline

1. Bellman optimality — property of V^\star

2. Optimal planning: Value Iteration

Define Bellman Operator \mathcal{T} :

Given a function $f : S \times A \mapsto \mathbb{R}$,

$$\mathcal{T}f : S \times A \mapsto \mathbb{R},$$

$$(\mathcal{T}f)(s, a) := r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a' \in A} f(s', a'), \forall s, a \in S \times A$$

Q: what is $\mathcal{T}Q^\star$?

Value Iteration Algorithm:

1. Initialization: $Q^0 : \|Q^0\|_\infty \in (0, \frac{1}{1-\gamma})$
2. Iterate until convergence: $Q^{t+1} = \mathcal{T} Q^t$

Intuition:

Via Bellman optimality theorem:

$$Q^{\star} = \mathcal{T} Q^{\star}$$

i.e., Q^{\star} is the fixed point solution of $f = \mathcal{T}f$

Consider the simple problem: finding fixed point solution $x^{\star} = \ell(x^{\star})$

$$x_0, x_{t+1} = \ell(x_t), t = 0, \dots,$$

$$|x_t - x^{\star}| = |\ell(x_{t-1}) - \ell(x^{\star})| \leq L |x_{t-1} - x^{\star}|$$

If $L < 1$ (i.e., contraction), then it converges exponentially fast

Convergence of Value Iteration:

Lemma [contraction]: Given any Q, Q' , we have:

$$\|\mathcal{T}Q - \mathcal{T}Q'\|_\infty \leq \gamma \|Q - Q'\|_\infty$$

Proof:

$$\begin{aligned} |\mathcal{T}Q(s, a) - \mathcal{T}Q'(s, a)| &= \left| r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} \max_{a'} Q(s', a') - \left(r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} \max_{a'} Q'(s', a') \right) \right| \\ &\leq \gamma \sum_{s'} P(s' | s, a) \left| \left(\max_{a'} Q(s', a') - \max_{a'} Q'(s', a') \right) \right| \\ &\leq \gamma \sum_{s'} P(s' | s, a) \max_{a'} \left| (Q(s', a') - Q'(s', a')) \right| \\ &\leq \gamma \max_{s'} \max_{a'} \left| (Q(s', a') - Q'(s', a')) \right| = \gamma \|Q - Q'\|_\infty \end{aligned}$$

Convergence of Value Iteration:

Lemma [Convergence]: Given Q^0 , we have:

$$\|Q^t - Q^\star\|_\infty \leq \gamma^t \|Q^0 - Q^\star\|_\infty$$

Proof:

$$\|Q^{t+1} - Q^\star\|_\infty = \|\mathcal{T}Q^t - \mathcal{T}Q^\star\|_\infty \leq \gamma \|Q^t - Q^\star\|_\infty$$

$$\dots \leq \gamma^{t+1} \|Q^0 - Q^\star\|_\infty$$

Final Quality of the Policy:

$$\pi^t : \pi^t(s) = \arg \max_a Q^t(s, a)$$

$$\textbf{Theorem: } V^{\pi^t}(s) \geq V^\star(s) - \frac{2\gamma^t}{1-\gamma} \|Q^0 - Q^\star\|_\infty \forall s \in S$$

Proof:

$$\begin{aligned} V^{\pi^t}(s) - V^\star(s) &= Q^{\pi^t}(s, \pi^t(s)) - Q^\star(s, \pi^\star(s)) \\ &= Q^{\pi^t}(s, \pi^t(s)) - Q^\star(s, \pi^t(s)) + Q^\star(s, \pi^t(s)) - Q^\star(s, \pi^\star(s)) \\ &= \gamma \mathbb{E}_{s' \sim P(s, \pi^t(s))} \left(V^{\pi^t}(s') - V^\star(s') \right) + Q^\star(s, \pi^t(s)) - Q^\star(s, \pi^\star(s)) \\ &\geq \gamma \mathbb{E}_{s' \sim P(s, \pi^t(s))} \left(V^{\pi^t}(s') - V^\star(s') \right) + Q^\star(s, \pi^t(s)) - Q^t(s, \pi^t(s)) + Q^t(s, \pi^\star(s)) - Q^\star(s, \pi^\star(s)) \\ &\geq \gamma \mathbb{E}_{s' \sim P(s, \pi^t(s))} \left(V^{\pi^t}(s') - V^\star(s') \right) - 2\gamma^t \|Q^0 - Q^\star\|_\infty \quad \dots \text{Recursion} \end{aligned}$$

Outline



1. Bellman optimality — property of V^*

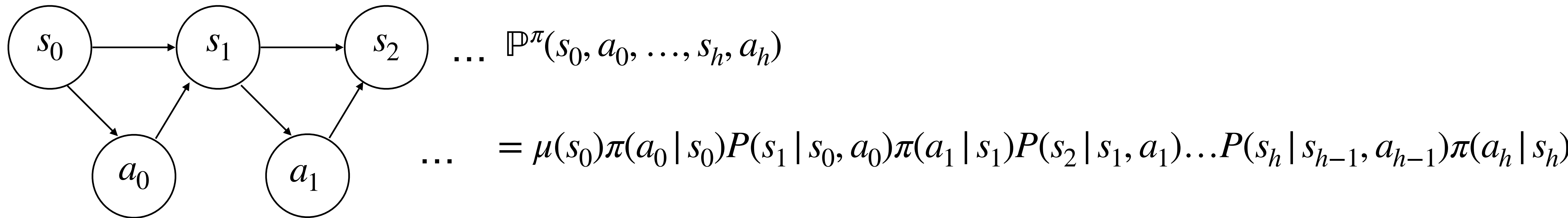


2. Optimal planning: Value Iteration

3. State-action distribution

Trajectory distribution and state-action distribution

Q: what is the probability of π generating trajectory $\tau = \{s_0, a_0, s_1, a_1, \dots, s_h, a_h\}$?



Q: what's the probability of π visiting state (s,a) at time step h ?

$$\mathbb{P}_h^\pi(s, a) = \sum_{s_0, a_0, s_1, a_1, \dots, s_{h-1}, a_{h-1}} \mathbb{P}^\pi(s_0, a_0, \dots, s_{h-1}, a_{h-1}, s_h = s, a_h = a)$$

Averaged state action occupancy measure

$\mathbb{P}_h^\pi(s, a)$: probability of π visiting (s, a) at time step $h \in \mathbb{N}$

$$d^\pi(s, a) = (1 - \gamma) \sum_{h=0}^{\infty} \gamma^h \mathbb{P}_h^\pi(s, a)$$

$$\mathbb{E}_{s_0 \sim \mu} V^\pi(s_0) = \frac{1}{1 - \gamma} \sum_{s, a} d^\pi(s, a) r(s, a)$$

Summary for today

Planning algorithm (no learning so far):

VI: fixed point iteration $Q^{t+1} = \mathcal{T} Q^t$

1. Bellman operator is a contraction map
2. $\|Q^t - Q^\star\|_\infty$ being small implies V^{π^t} & V^\star are close