

1) Value Iteration

from Q-functions to policies

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

$$Q^t \approx Q^*$$

is $\pi^t(s) = \operatorname{argmax}_{a \in \mathcal{A}} Q^t(s, a)$
a good policy?

Theorem:

$$\forall s \in \mathcal{S} \quad V^{\pi^t}(s) \geq V^*(s) - \frac{2\gamma^t}{1-\gamma} \|Q^0 - Q^*\|_\infty$$

Proof:

assume:

$$\rightarrow V^{\pi^t}(s) - V^*(s) \geq \gamma \mathbb{E}_{s' \sim P(s, \pi^t(s))} [V^{\pi^t}(s') - V^*(s')]$$

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

$$\begin{aligned}
 \underbrace{V^{\pi^t}(s) - V^*(s)} &\geq \gamma \underbrace{\mathbb{E}_{s'} [V^{\pi^t}(s') - V^*(s')]} - 2\gamma^t \|Q^0 - Q^t\|_\infty \\
 &\geq \gamma \mathbb{E}_{s'} \left[\gamma \mathbb{E} [V^{\pi^t}(s') - V^*(s')] - 2\gamma^t \|Q^0 - Q^t\|_\infty \right] \\
 &\quad - 2\gamma^t \|Q^0 - Q^t\|_\infty
 \end{aligned}$$

$$\geq \gamma^k \mathbb{E}_{s^k} [V^{\pi^t}(s^k) - V^*(s^k)] \rightarrow 0$$

$$\gamma < 1 \quad \rightarrow \underbrace{2 \sum_{l=0}^{k-1} \gamma^{t+l}}_{\text{bound}} \|Q^0 - Q^t\|_\infty$$

let $k \rightarrow \infty$

$$\sum_{l=0}^{\infty} \gamma^{t+l} = \gamma^t \sum_{l=0}^{\infty} \gamma^l = \frac{1}{1-\gamma} \cdot \gamma^t$$

proof of assumption

1) write V^{π^t} & V^* in terms of Q^{π^t} & Q^*

2) add & subtract Q-fns.

3) use definition of π^t

2) Policy Iteration

Alg (PI)

initialize $\pi^0: \mathcal{S} \rightarrow \mathcal{A}$

for $t=0, 1, \dots$

Policy Evaluation: $Q^{\pi^t}(s, a) \forall s, a$

Policy Improvement:

$$\forall s \quad \pi^{t+1}(s) = \underset{a \in \mathcal{A}}{\operatorname{argmax}} Q^{\pi^t}(s, a)$$

aside: PE-exact $V = R + \gamma P V$

$$\mathbb{R}^{\mathcal{S}} \rightarrow V^{\pi^t} = (I - \gamma P)^{-1} R$$

$$Q^{\pi^t}(s, a) = r(s, a) + \mathbb{E}_{s' \sim P(s, a)} [V^{\pi^t}(s')]$$

Two key properties

1) Monotonic Improvement

$$Q^{\pi^{t+1}}(s,a) \geq Q^{\pi^t}(s,a) \quad \forall s,a$$

2) convergence

$$\|V^* - V^{\pi^t}\|_{\infty} \leq \gamma^t \|V^* - V^{\pi_0}\|_{\infty}$$

Lemma (Monotonic Improvement)

$$Q^{\pi^{t+1}}(s,a) \geq Q^{\pi^t}(s,a)$$

Proof:

$$\begin{aligned} & Q^{\pi^{t+1}}(s,a) - Q^{\pi^t}(s,a) \\ &= \mathbb{E}_{s' \sim P(s,a)} \left[V^{\pi^{t+1}}(s') - V^{\pi^t}(s') \right] \\ &= \underbrace{Q^{\pi^{t+1}}(s', \pi^{t+1}(s'))}_{\star} - \underbrace{Q^{\pi^t}(s', \pi^t(s'))}_{\star} \end{aligned}$$

$$\begin{aligned} & \rightarrow \sum_{s'} \left(Q^{t+1}(s', \pi^{t+1}(s')) - Q^t(s', \pi^t(s')) \right) \\ & + Q^{t+1}(s', \pi^{t+1}(s')) - Q^t(s', \pi^t(s')) \end{aligned}$$

$$Q^{t+1}(s, a) - Q^t(s, a)$$

$$\sum_{s'} \left[Q^{t+1}(s', \pi^{t+1}(s')) - Q^t(s', \pi^{t+1}(s')) \right]$$

iterate k times

$$\sum 0$$

□

Q: does lemma imply that

$$\underline{V^{\pi^{t+1}}(s) \geq V^{\pi^t}(s) \quad \forall s?}$$

$$Q^{\pi^{t+1}}(s, \pi^{t+1}(s)) \geq \underbrace{Q^{\pi^t}(s, \pi^{t+1}(s))}_{\text{because}}$$

$$\pi^{t+1} = \underset{a}{\operatorname{argmax}} Q^{\pi^t}(s, a)$$

$$\geq \underbrace{Q^{\pi^t}(s, \pi^t(s))}_{V^{\pi^t}(s)}$$

Theorem: Convergence

$$\|V^{\pi^t} - V^*\|_{\infty} \leq \gamma^t \|V^{\pi_0} - V^*\|_{\infty}$$

Proof sketch:

- 1) use Q fn; Bellman opt.
- 2) monotonic improvement
- 3) $\pi^{t+1} = \underset{a}{\operatorname{argmax}} Q^{\pi^t}(s, a)$
- 4) Basic Inequalities

Value & Policy Iteration

- converge exponentially fast
 γ^t

- exact solution?

$\gamma^t \rightarrow 0$ as $t \rightarrow \infty$

But $\gamma^T > 0 \quad \forall T$

- only policy iteration
is guaranteed to
find exact v policy
optimal

in finite # steps

(HW1)

3) Finite Horizon MDP

$$\mathcal{M} = \{ \mathcal{S}, \mathcal{A}, P, r, H, \mu_0 \}$$

$$H \in \mathbb{N}^+ = \{1, 2, 3, \dots\}$$

length of time

$$\mu_0 \in \Delta(\mathcal{S}) \quad \text{initial state distribution}$$

time-varying policies

$$\pi = (\pi_0, \pi_1, \pi_2, \dots, \pi_{H-1})$$

$$\pi_t: \mathcal{S} \rightarrow \mathcal{A}$$

value & Q fn.

$$V_t^\pi(s) = \mathbb{E} \left[\sum_{k=t}^{H-1} r(s_k, a_k) \right]$$

$$\left. \begin{array}{l} s_t = s \\ s_{k+1} \sim P(s_k, a_k) \\ a_k = \pi_k(s_k) \end{array} \right\}$$

$$Q_t^\pi(s, a) = \mathbb{E} \left[\sum_{k=t}^{H-1} r(s_k, a_k) \right]$$

$$\left. \begin{array}{l} s_t = s, a_t = a \\ s_{k+1} \sim P(s_k, a_k) \\ a_k = \pi_k(s_k) \end{array} \right\}$$

Bellman Equation:

$$Q_t^\pi(s, a) = r(s, a) + \mathbb{E} [V_{t+1}^\pi(s')] \\ s' \sim P(s, a)$$

4) Dynamic Programming

to find $\pi^* = (\pi_0^*, \dots, \pi_{H-1}^*)$

start at $H-1$: (note: $V_H^*(s) = 0$)

$$Q_{H-1}^*(s, a) = r(s, a)$$

$$\pi_{H-1}^* = \operatorname{argmax}_{a \in \mathcal{A}} Q_{H-1}^*(s, a)$$

$$V_{H-1}^* = Q_{H-1}^*(s, \pi_{H-1}^*(s))$$

$$\rightarrow Q_t^*(s, a) = r(s, a) + \mathbb{E} [V_{t+1}^*(s')] \\ s' \sim P(s, a)$$

$$\rightarrow \pi_t^*(s) = \operatorname{argmax}_{a \in \mathcal{A}} Q_t^*(s, a) \\ s' \sim P(s, a)$$

Dynamic Programming
terminates in H steps.

exact H^*