

1/18

Review

* Markov decision process

$$MDP = \{ \mathcal{S}, \mathcal{A}, P(s'|s, a), R(s) \}$$

* Policy = deterministic mapping $\pi: \mathcal{S} \rightarrow \mathcal{A}$

* Value functions

$$V^\pi(s) = E^\pi \left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \mid S_0 = s \right]$$

$$Q^\pi(s, a) = E^\pi \left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \mid S_0 = s, a_0 = a \right]$$

* Policy evaluation

$$\text{Solve linear equations: } \sum_{s'} [I(s, s') - \gamma P(s'|s, \pi(s))] V^\pi(s') = R(s)$$

* Policy improvement

$$\text{Greedy policy } \pi'(s) = \underset{a}{\operatorname{argmax}} Q^\pi(s, a)$$

Thm: $V^{\pi'}(s) \geq V^\pi(s)$ for all states s

* Policy iteration

$$\pi_0 \xrightarrow{\text{evaluate}} \begin{matrix} V^{\pi_0}(s) \\ Q^{\pi_0}(s) \end{matrix} \xrightarrow{\text{improve}} \pi_1 \longrightarrow \dots$$

Thm: if $\pi'(s) = \underset{a}{\operatorname{argmax}} Q^\pi(s, a)$ and $V^{\pi'}(s) = V^\pi(s)$ for all s , then $V^{\pi'}(s) = V^*(s)$.

Proof: 1) From last time:

$$V^\pi(s) = R(s) + \gamma \max_a \sum_{s'} P(s'|s, a) V^\pi(s') \quad \text{: Bellman optimality equation}$$

2) Iterate RHS:

$$V^\pi(s) = R(s) + \gamma \max_a \sum_{s'} P(s'|s, a) [R(s') + \gamma \max_{a'} \sum_{s''} P(s''|s', a') V^\pi(s'')] \quad \text{: (B)}$$

Imagine iterating above t times.

Let's show that this iterated expression implies optimality.

Let $\hat{\pi}(s)$ be any other policy:

$$V^{\hat{\pi}}(s) = R(s) + \gamma \sum_{s'} P(s'|s, \hat{\pi}(s)) V^{\hat{\pi}}(s') \quad \text{: Bellman equation}$$

$$\leq R(s) + \gamma \max_a \sum_{s'} P(s'|s, a) V^{\hat{\pi}}(s') \quad \text{: greedy}$$

$$= R(s) + \gamma \max_a \sum_{s'} P(s'|s, a) [R(s') + \gamma \sum_{s''} P(s''|s', \hat{\pi}(s')) V^{\hat{\pi}}(s'')] \quad \text{: iterate}$$

$$V^{\hat{\pi}}(s) \leq R(s) + \gamma \max_a \sum_{s'} P(s'|s, a) [R(s') + \gamma \max_{a'} \sum_{s''} P(s''|s', a') V^{\hat{\pi}}(s'')] \quad \text{: be greedy}$$

$$= \dots \leq \dots = \dots \quad \text{: (A)}$$

Consider upper bound on $V^{\hat{\pi}}(s)$ by iterating t times. (being greedy then applying Bellman equation.)

Compare to equality after t iterations for $V^\pi(s)$.

As $t \rightarrow \infty$, RHS on upper bound on $V^{\hat{\pi}}(s)$ converges to RHS of equality for $V^\pi(s)$.

Thus, as $t \rightarrow \infty$: $V^{\hat{\pi}}(s) \leq \lim_{t \rightarrow \infty} [\textcircled{A}] = \lim_{t \rightarrow \infty} [\textcircled{B}] = V^{\pi}(s)$

Thus, for all policies $\hat{\pi}(s)$ and states s : $V^{\pi}(s) \geq V^{\hat{\pi}}(s)$

$$V^{\pi}(s) = \max_{\hat{\pi}} V^{\hat{\pi}}(s) \rightarrow V^{\pi}(s) = V^*(s)$$

Finally, once you have $V^*(s)$,

$$\begin{aligned} \pi^*(s) &= \operatorname{argmax}_a Q^*(s, a) \\ &= \operatorname{argmax}_a \sum_{s'} P(s'|s, a) V^*(s') \end{aligned}$$

* Pros & Cons of policy iteration.

(+) converges quickly in finite # steps

(-) requires policy evaluation $O(n^3)$ at each iteration.

Value Iteration

* How to compute $V^*(s)$ directly?

* Bellman optimality equation.

$$\begin{aligned} V^*(s) &= \max_a Q^*(s, a) \\ &= \max_a [R(s) + \gamma \sum_{s'} P(s'|s, a) V^*(s')] \\ &= R(s) + \gamma \max_a \sum_{s'} P(s'|s, a) V^*(s') \quad (*) \end{aligned}$$

Before we showed that if $V^{\pi}(s)$ satisfied (*), then π is optimal.

Here, we show that $V^*(s)$ must satisfy (*).

Q: How to solve this set of N nonlinear equations for N unknowns?

* Algorithm

1) initialize: $V_0(s) = 0$ for all s .

2) iterate: $V_{k+1}(s) = R(s) + \gamma \max_a [\sum_{s'} P(s'|s, a) V_k(s')] \quad \leftarrow \text{estimate at } k^{\text{th}} \text{ iteration}$ for all $s=1, 2, \dots, n$

Note: this algorithm works on value functions, not policies.

But, incremental policies can be computed as:

$$\pi_{k+1}(s) = \operatorname{greedy}[V_k(s)] = \operatorname{argmax}_a [\sum_{s'} P(s'|s, a) V_k(s')]$$

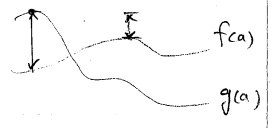
3) after convergence:

$$\pi^*(s) = \operatorname{argmax}_a [Q^*(s, a)] = \operatorname{argmax}_a [\sum_{s'} P(s'|s, a) V^*(s')]$$

Does algorithm converge? $V^*(s)$ is obviously a fixed point.

But are there others? Does it always reach $V^*(s)$?

* Lemma: $|\max_a [f(a)] - \max_a [g(a)]| \leq \max_a |f(a) - g(a)|$



Proof of lemma:

For all a : $f(a) - \max_{a'} g(a') \leq f(a) - g(a)$

Max over a : $\max_a f(a) - \max_{a'} g(a') \leq \max_a [f(a) - g(a)] \leq \max_a |f(a) - g(a)|$

By symmetry: $\max_a g(a) - \max_{a'} f(a') \leq \max_a |f(a) - g(a)|$

Combining last two bounds proves lemma.

* Theorem: value iteration converges.

$$\lim_{k \rightarrow \infty} [V_k(s)] \rightarrow V^*(s) \text{ for all states } s.$$

* Proof: let $\Delta_k = \max_s |V_k(s) - V^*(s)|$ error at k^{th} iteration.

$$\begin{aligned} \Delta_{k+1} &= \max_s |V_{k+1}(s) - V^*(s)| \\ &= \max_s \left| [R(s) + \gamma \max_a \sum_{s'} P(s'|s,a) V_k(s')] \right. \\ &\quad \left. - [R(s) + \gamma \max_a \sum_{s'} P(s'|s,a) V^*(s')] \right| \\ &= \gamma \max_s \left| \underbrace{\max_a \sum_{s'} P(s'|s,a) V_k(s')}_{f(a)} - \underbrace{\max_a \sum_{s'} P(s'|s,a) V^*(s')}_{g(a)} \right| \end{aligned}$$

Apply lemma

$$\begin{aligned} \Delta_{k+1} &\leq \gamma \max_s \max_a \left| \sum_{s'} P(s'|s,a) V_k(s') - \sum_{s'} P(s'|s,a) V^*(s') \right| \\ &= \gamma \max_s \max_a \left| \sum_{s'} P(s'|s,a) [V_k(s') - V^*(s')] \right| \\ &\leq \gamma \max_s \max_a \left| \sum_{s'} P(s'|s,a) \max_{s''} |V_k(s'') - V^*(s'')| \right| \\ &= \gamma \max_s \max_a \left| \sum_{s'} P(s'|s,a) \Delta_k \right| \\ &= \gamma \Delta_k \max_s \max_a (1) \end{aligned}$$

Hence: $\Delta_{k+1} \leq \gamma \Delta_k$ with $\gamma < 1$. so-called "contraction" mapping.

By iteration: $\Delta_k \leq \gamma^k \Delta_0 \rightarrow 0$ as $k \rightarrow \infty$.

$$\begin{aligned} \text{Assume rewards are bounded: } \Delta_0 &= \max_s |V_0(s) - V^*(s)| \\ &= \max_s |V^*(s)| \\ &\leq \max_s |R(s)| (1 + \gamma + \gamma^2 + \dots) \\ &= \max_s |R(s)| \cdot \left(\frac{1}{1-\gamma} \right) \end{aligned}$$

More iterations required as $\gamma \rightarrow 1$.