

Lecture 17

CS 59000-RL

MDP

- Optimality
- Value iteration
- Policy iteration

Bellman optimality operator is a contraction.

In particular:

For $U, V \in \mathbb{R}^{|X|}$

$$\|TU - TV\|_{\infty} \leq \lambda \|U - V\|_{\infty}$$

Theorem (Banach Fixed-Point) Suppose \mathcal{L} is a Banach space and $\mathcal{L} : \mathcal{L} \rightarrow \mathcal{L}$ is a contraction map under a norm $\|\cdot\|$, and parameter λ . Then

i) There exist a unique V^* in \mathcal{L} such that $\mathcal{L}V^* = V^*$

ii) For an arbitrary $V_0 \in \mathcal{L}$, the sequence

$\{V_n\}$ defined as $V_{n+1} = \mathcal{L}V_n = \mathcal{L}^{n+1}V_0$

converges to V^*

Proof:

$$\begin{aligned}
 \Rightarrow \|v_{n+m} - v_n\| &\leq \sum_{k=0}^{m-1} \|v_{n+k+1} - v_{n+k}\| \\
 &= \sum_{k=0}^{m-1} \|\mathcal{L}^{n+k} v_1 - \mathcal{L}^{n+k} v_0\| \\
 &\leq \sum_{k=0}^{m-1} \lambda^{n+k} \|v_1 - v_0\| = \underbrace{\lambda^n \frac{(1-\lambda^m)}{1-\lambda}}_{\leq \varepsilon} \|v_1 - v_0\|
 \end{aligned}$$

Therefore, $\{v_n\}$ is a Cauchy sequence, i.e. for

any ε , there exist a sufficiently large n

such that $\|v_{n+m} - v_n\| \leq \varepsilon$

Since \mathcal{L} is Banach, then the Cauchy sequence $\{v_n\}$ converges. Let v^+ denote the limit.

Now we show that $\mathcal{L} v^+ = v^+$

For $n \geq 1$

$$\hookrightarrow \pm V_n$$

$$0 \leq \| \mathcal{L} V^* - V^* \|$$

$$\leq \| \mathcal{L} V^* - V_n \| + \| V_n - V^* \|$$

$$\leq \lambda \| \underbrace{V^* - V_{n-1}}_{\nearrow} \| + \| \underbrace{V_n - V^*} \|$$

when taking the limit $n \rightarrow \infty$; $\lim_{n \rightarrow \infty} \| V_n - V^* \| = 0$

$$\text{Ergo } \| \mathcal{L} V^* - V^* \| = 0 \Rightarrow \mathcal{L} V^* = V^*$$

So, there exists a solution

Uniqueness: let's imagine there exists another solution U^* such that

- i) $U^* \neq V^*$
- ii) and $\mathcal{L} U^* = U^*$

$$V^* - U^* = \mathcal{L} V^* - \mathcal{L} U^*$$

$$\hookrightarrow \| \underbrace{V^* - U^*} \| = \| \mathcal{L} V^* - \mathcal{L} U^* \| \leq \lambda \| \underbrace{V^* - U^*} \|$$

$$\Rightarrow \| V^* - U^* \| = 0 \Rightarrow V^* = U^*$$

How we use it for MDPs?

Using the fact that $R^{|X|}$ is a Banach space we have the Bellman optimality equation has a unique solution. $V^* = T V^*$

we need to show $V^* = V^*$.

$$V^*(x) = \max_a \bar{r}(x, a) + \lambda \sum_{x'} P(x'|x, a) V^*(x')$$

$$\geq \bar{r}(x, a) + \lambda \sum_{x'} P(x'|x, a) V^*(x')$$

Therefore:

$$V^* \geq \bar{r}_{\pi_1} + \lambda P_{\pi_1} V^*$$

$$\geq \bar{r}_{\pi_1} + \lambda P_{\pi_1} (\bar{r}_{\pi_2} + \lambda P_{\pi_2} V^*)$$

$$= \bar{r}_{\pi_1} + \lambda P_{\pi_1} \bar{r}_{\pi_2} + \lambda^2 P_{\pi_1} P_{\pi_2} V^*$$

$$\geq \bar{r}_{\pi_1} + \lambda P_{\pi_1} \bar{r}_{\pi_2} + \dots + \lambda^n P_{\pi_1} \dots P_{\pi_n} V^*$$

$$V^* = V^{\pi^*} = \sum_{t=1}^{\infty} \lambda^{t-1} P_{\pi_t}^* \bar{r}_{\pi_t}$$

This inequality holds for any sequence of policies $\bar{\pi}$

$$\Rightarrow V^+ \geq V^{\bar{\pi}} + \underbrace{\lambda^n \bar{P}_{\bar{\pi}_n}^n V^+}_{\text{since } V^{\bar{\pi}} = \sum \lambda^{t+1} \bar{P}_{\bar{\pi}_t}^{t+1} r_{\bar{\pi}_t}} - \underbrace{\sum_{k=n}^{\infty} \lambda^k \bar{P}_{\bar{\pi}_k}^k r_{\bar{\pi}_{k+1}}}_{\text{}} \quad \uparrow \quad \uparrow$$

First: $\|\lambda^n \bar{P}_{\bar{\pi}_n}^n V^+\|_{\infty} \leq \lambda^n \|V^+\|_{\infty}$

second: since $\|\bar{r}(x,a)\| \leq M < \infty$, we have:

vector of all ones. $\underbrace{\frac{-\lambda^n M e}{1-\lambda}}_{\text{}} \leq \underbrace{\sum_{k=n}^{\infty} \lambda^k \bar{P}_{\bar{\pi}_k}^k r_{\bar{\pi}_{k+1}}}_{\text{}}$

Therefore, for any $\varepsilon > 0$, there exists an n such that

$$V^+ - V^{\bar{\pi}} \geq -\varepsilon \rightarrow V^+ \geq V^{\bar{\pi}} \\ \Rightarrow V^+ \geq V^{\pi^*} = V^*$$

we are left with showing $V^+ \leq V^*$

Lemma: If $V \geq 0$, then $\overbrace{(I - \lambda P_\pi)^{-1}}^{\downarrow} V \geq V$

Proof: Using the fact that, $\|I - \lambda P_\pi\| < 1$,

and Neumann series of invertible operators:

$$\underbrace{(I - \lambda P_\pi)^{-1}}_{\text{non negative}} V = V + \underbrace{\lambda P_\pi V + \lambda^2 P_\pi^2 V + \dots}_{\text{non negative}} \geq V \geq 0$$

For Bellman optimality equation, we have $TV = V^*$

Therefore, for any $\epsilon > 0$, there exists a π , such that

$$V^* \leq r_\pi + \lambda P_\pi V^* + \epsilon e$$

$$\Rightarrow \underbrace{r_\pi + \epsilon e - (I - \lambda P_\pi) V^*}_{\geq 0} \geq 0$$

This is positive? Using the lemma we just proved we have

$$V^* \leq \underbrace{(I - \lambda P_\pi)^{-1}}_{\text{non negative}} (r_\pi + \epsilon e) = V^\pi + \underbrace{(I - \lambda P_\pi)^{-1}}_{\text{non negative}} \epsilon e$$

$$V^* \leq V^\pi + (1 - \lambda)^{-1} \epsilon e \quad \text{for any } \epsilon > 0$$

$$\Rightarrow V^+ \leq V^{n+}$$

$$\Rightarrow V^+ = V^*$$

Value iteration:

- Initialize with V_0
- Repeatedly apply $V_{n+1} = T V_n$

we had in the analysis:

$$\underbrace{\|V_{n+1} - V^*\|_\infty}_{\hookrightarrow} = \|T V_n - T V^*\|_\infty \leq \underbrace{\lambda \|V_n - V^*\|_\infty}_{\leq \lambda^{n+1} \|V_0 - V^*\|_\infty}$$

we showed value iteration converges exponentially fast.

Policy iteration:

- start with an arbitrary stationary and memory-less policy π_0 , then solve for V

Policy evaluation step $\rightarrow (I - \gamma P_{\pi_n}) V = r_{\pi_n}$

call the solution V_n (value of Π_n)

— Improve policy $\leftarrow \Pi_{n+1} \in \operatorname{argmax} (r_{\Pi} + \lambda P_{\Pi} V_n)$

Policy improvement step

Is Π_{n+1} better than Π_n ?

we know: $r_{\Pi_{n+1}} + \lambda P_{\Pi_{n+1}} V_n > r_{\Pi_n} + \lambda P_{\Pi_n} V_n = V_n$

$$\Rightarrow r_{\Pi_{n+1}} \geq (I - \lambda P_{\Pi_{n+1}}) V_n$$

Applying the lemma we previously used
(for $V \geq v \rightarrow (I - \lambda P_{\Pi})^{-1} V \geq v$)

we have: applying $(I - \lambda P_{\Pi_{n+1}})^{-1}$

$$V_{n+1} = (I - \lambda P_{\Pi_{n+1}})^{-1} r_{\Pi_{n+1}} \geq V_n$$

Policy iteration algorithm converges to an optimal policy and optimal value.

refer to 6.S MDP Book by
Martin. Puterman.

Example:

For $\lambda' < \lambda$



$$V_{\lambda'}^*$$

$$V_{\lambda}^*$$

after n step

Iteration with $\lambda \rightarrow \lambda^n$

$$\|V_{\lambda}^* - V_{\lambda}^n\|_{\infty} < \lambda^n$$

$$\|V_{\lambda'}^* - V_{\lambda'}^n\|_{\infty} < \lambda'^n$$

If n is small and $\|V_{\lambda'}^* - V_{\lambda'}^n\| \leq \delta$

$$\lambda^n > \delta + \lambda'^n$$