# Lecture 16

CS 59000_RL

MDP

Agenda

- Bellman equation for Fixed Horizon MDPs
- Discounted infinite horizon MDPs.

---

For a policy $\pi \in \Pi^{MR}$ , $\forall_t, \forall_x$

$$\underbrace{V_t^\pi(x_t)} = \sum_a \pi(a|x_t)\, \bar{r}_t(x_t, a)$$

$$V_t^\pi = T^\pi V_{t+1}^\pi \implies \qquad + \sum_a \sum_x P(x_{t+1} = x | x_t, a)\, \pi(a; x_t)\, \underbrace{V_{t+1}^\pi(x)}$$

This is known as Bellman equation

or Bellman consistency equation.

$T^\pi$ is Bellman operator.

$$V_t^*(x_t) = \max_{a \in \mathcal{A}} \left( \bar{r}(x_t, a) + \sum_x P(x_{t+1} = x | x_t, a)\, V_{t+1}^*(x) \right)$$

is known as Bellman optimality equation.

and $\qquad V_b^* = T^*\left(V_{b+1}\right)$

where $T^*$ is Bellman optimality operator.

Q function: For a policy $\pi \in \Pi^{MR}$

$$Q_t^\pi(x, a) = E^\pi \left[ \sum_{k=t}^A r_k \mid X_t = x, A_t = a \right]$$

$$= \bar{r}_t(x, a) + \sum_{x'} P(X_{t+1} = x' \mid x, a) V_{t+1}^\pi(x)$$

and for the optimal policy.

$$Q_t^*(x, a) = \bar{r}_t(x, a) + \sum_{x'} P(X_{t+1} = x' \mid x, a) V_{t+1}^*(x')$$

$$\hookrightarrow$$

$$\pi^* : \quad \arg\max_{a \in t} Q_t^*(x, a)$$

---

Tabular discounted infinite horizon MDPs
(Finite state-action spaces

Consider stationary setting where we have
$$P(x' \mid x, a) \qquad R(x, a) \qquad \text{independent of time}$$

For each state $x$, we have:

$$\implies V^\pi(x) = E \left[ \sum_{t=1}^\infty \lambda^{t-1} r_t \mid X_1 = x \right]$$

Optimal: $V^*(x) = \max_{\pi \in \Pi^{HR}} V^\pi(x)$

Assume $|\bar{r}(x,a)| \leq M < \infty$.

✳ With somewhat different argument,
There is a <u>stationary</u> <u>memory-less</u>
Markovian policy $\pi^*$ that is optimal (Also deterministic)

Therefore we focus on this policy class.
For the proof: Chapter 6 of MDP book]

Let's use vector and matrix notation

$$\gamma_\pi \in R^{|X|} \quad , \quad \gamma_{\pi,i} = \sum_a \bar{r}(x=i, a) \, \pi(a; x=i)$$

$$P_\pi \in R^{|X| \times |X|} \quad , \quad (P_\pi)_{i,j} = \sum_a P(j|x,a) \, \pi(a; i)$$

Now

$$V^\pi = \sum_{t=1}^{\infty} \lambda^{t-1} P_\pi^{t-1} r_1$$

$$= \gamma_\pi + [\lambda P_\pi r_\pi + \lambda^2 P_\pi^2 r_\pi + \cdots]$$

$$= \gamma_\pi + \lambda P_\pi V_\pi$$

$$= T^\pi(V^\pi) \qquad \text{which is Bellman operator.}$$

Cool thing: $T^\pi(V) = V \implies$ has
a unique solution, and $V^\pi$ is that one.

Let's play with this equation:

$$v^{\pi} = T^{\pi}(v^{\pi}) = T^{\pi} v^{\pi} = r_{\pi} + \lambda \underbrace{P_0 v^{\pi}}$$

$$\underbrace{(I - \lambda P_{\pi})} v^{\pi} = r$$

If $(I - \lambda P_{\pi})$ is full rank, $\Rightarrow v^{\pi}$ is the solution

to $T^{\pi}(v) = v$

proof:

Note that $\|P_{\pi}\|_{\infty} = 1$, and $\|\lambda P_{\pi}\|_{\infty} = \lambda < 1$

Therefor, $(I - \lambda P_{\pi})$ is full rank and $(I - \lambda P_{\pi})^{-1}$

exists.

Another way of proving it. For $\alpha \in R^{|x|} \neq 0$

$$\underbrace{\|(I - \lambda P^{\pi}) \alpha\|_{\infty}} = \|\alpha - \lambda P_{\pi} \alpha\|_{\infty}$$

Triangle inequality $\geq \|\alpha\|_{\infty} - \lambda \|P_{\pi}\alpha\|_{\infty}$

$$\geq \|\alpha\|_{\infty} - \lambda \|\alpha\|_{\infty}$$

$$= (1 - \lambda)\|\alpha\|_{\infty} > 0$$

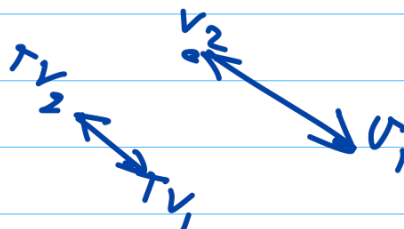Now consider the Bellman optimality equation

$$V = \max_{\Pi} \left( r_{\Pi} + \lambda P_{\Pi} V \right)$$

$$= T V$$

Theorem: There exist a unique $V$ that satisfies the Bellman optimality equation and it is the optimal $V^*$.

---

In the first step we need to show a solution exists. Then show it is unique.
    Ten, show, such solution is $V^*$.

Lemma: The $T$ operator is contraction.
under $\|\cdot\|_{\infty}$ norm.



proof: Consider $U, V \in \mathbb{R}^{|x|}$. For a state $x$,

$$a_x^* \in \underset{a \in A}{\arg\max} \left( \bar{r}(x,a) + \lambda \sum_{x'} P(x'|x,a) V(x') \right)$$

an optimal action if $V$ was the value.

Now: Assume $T V(x) \geqslant T U(x)$, then

$$0 \leq TV(x) - TU(x)$$

$$\leq \bar{r}(x, a_x^*) + \lambda \sum_{x' \in \mathcal{X}} P(x'|x, a_x^*) V(x')$$

$$- \left( \bar{r}(x, a_x^*) + \lambda \sum_{x' \in \mathcal{X}} P(x', x, a_x^*) U(x') \right)$$

$$\leq \lambda \sum_{x' \in \mathcal{X}} P(x'|x, a_x^*) \left( V(x') - U(x') \right)$$

$$\leq \lambda \sum_{x' \in \mathcal{X}} P(x'|x, a_x^*) \| V - U \|_\infty = \lambda \| V - U \|_\infty$$

we can do the same when $Tv(x) \leq TU(x)$

$$\Rightarrow$$

$$0 \leq TU(x) - TV(x) \leq \lambda \| V - U \|_\infty$$

Making the same argument for all $x \in \mathcal{X}$

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

Theorem: (Banach Fixed-Point)