

CS 59006-RL MDPs

Agenda:

- Tabular Fixed Horizon MDPs
- Infinite Horizon Discounted reward MDP

Tabular fixed horizon MDPs with finite state and action:

$$M := (\mathcal{X}, \mathcal{A}, \underset{\uparrow}{P}, \underset{\uparrow}{R}, P_1, H)$$

$$|\mathcal{X}| < \infty \quad |\mathcal{A}| < \infty$$

at time t , state X_t , taking action A_t
result in reward $r_t \sim R_t(X_t, A_t)$

$$X_{t+1} \sim P_t(X_t, A_t)$$

simplicity $P_t(X_{t+1} | X_t, A_t)$ to denote
the probability of next state.

(It is fine when we deal with finite
state and action spaces)

$$\text{Let } \bar{r}_b(x_b, A_b) = E[r_b | \sigma(x_b, A_b)]$$

Also, for the last state, we have

$$r_H \sim R_H(x_H) \quad \text{and} \quad \bar{r}_H \text{ is its expected value.}$$

starting from state x , and following π , we have

$$\begin{aligned} V_1^\pi(x) &= E^\pi \left[\sum_{k=1}^{H-1} r_k + r_H \mid \sigma(x) \right] (x) \\ &= E^\pi \left[\sum_{k=1}^{H-1} \bar{r}_k(x_k, A_k) + r_H(x_H) \mid \sigma(x) \right] (x) \end{aligned}$$

similarly for any $1 < b < H$

$$V_b^\pi(h_b) = E^\pi \left[\sum_{k=b}^{H-1} \bar{r}_k(x_k, A_k) + \bar{r}_H(x_H) \mid \sigma(h_b) \right] (h_b)$$

For simplicity we use the following notation

$$V_b^\pi(h_b) = E^\pi \left[\sum_{k=b}^{H-1} \bar{r}_k(x_k, A_k) + \bar{r}_H(x_H) \mid h_b \right]$$

This is the value of h_b following π

note that

$$v_t^\pi(h_t) = E^\pi \left[\bar{r}_t(x_t, A_t) \mid h_t \right] + E^\pi \left[E^\pi \left[\sum_{k=t+1}^{H-1} \bar{r}_k(x_k, A_k) + \bar{r}_H(x_H) \mid \mathcal{I}_{t+1} \right] \mid h_t \right]$$

$\leftarrow v_{t+1}^\pi(h_{t+1})$

$$v_t^\pi = E^\pi \left[\bar{r}_t(x_t, A_t) \mid h_t \right] + E^\pi \left[v_{t+1}^\pi(h_{t+1}) \mid h_t \right]$$

Let's look at the first term:

$$E^\pi \left[\bar{r}_t(x_t, A_t) \mid h_t \right] = \sum_{a \in \mathcal{A}} \pi(a; h_t) \bar{r}_t(x_t, a)$$

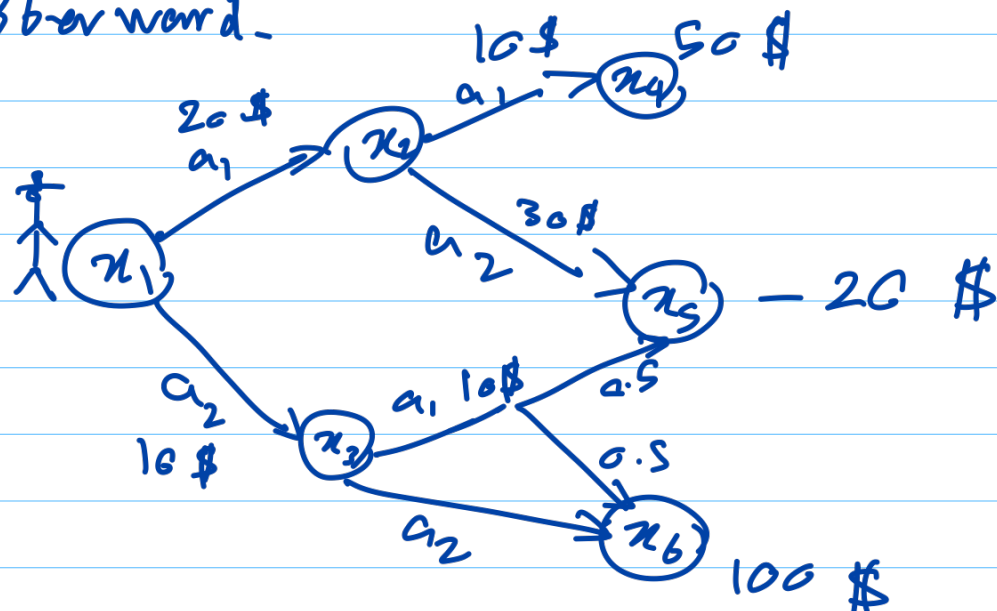
And the second term:

$$E^\pi \left[v_{t+1}^\pi(h_{t+1}) \mid h_t \right] = \sum_{x \in \mathcal{X}} \sum_{a \in \mathcal{A}} \pi(a; h_t) P_{t+1}^{\pi}(x_{t+1} \mid x_t, a) v_{t+1}^\pi(h_t, a, x)$$

most likely chapter 2, 4, 5, 6, 7 of
MDP book by Martin Puterman

Monday, October 12, 2020

Therefore to compute V_b^π , we look at what is the immediate reward, and what we expect afterward.



Consider this policy π : At x_1 , choose a_1 , and at any other state choose a_1, a_2 uniformly at random.

$$\left. \begin{aligned} V_3^\pi(x_4) &= 50 \\ V_3^\pi(x_5) &= -20 \\ V_3^\pi(x_6) &= 100 \end{aligned} \right\} \begin{aligned} V_2^\pi(x_2) &= \frac{1}{2} 10 + \frac{1}{2} 30 \\ &\quad + \frac{1}{2} 50 + \frac{1}{2} (-20) = 35 \end{aligned}$$

$$V_1(x_1) = 20 + 35 = 55$$

The thing we have done is called
Dynamic programming

Optimality:

$$V_1^*(x) = \max_{\pi \in \Pi^{\text{HR}}} V_1^\pi(x) \quad x \in \mathcal{X}$$

Let's define $u_\#(h_\#) = \bar{r}_\#(x_\#)$

$$\text{and } u_b(h_b) = \max_{a \in \mathcal{A}} \left(\bar{r}_b(x_b, a) + \sum_{x_{b+1}} P(x_{b+1} | x_b, a) u_{b+1}(h_{b+1}, a, x_{b+1}) \right)$$

now we show $u_b = V_b^*$ and there a deterministic optimal policy.

Proof:

— we know that $u_\# = V_\#^*$ by definition

— we can write $u_b(h_b) \geq V_b^*(h_b)$ for all h_b .

— Now let's do some induction:

$$\text{Let } u_{b+1}(h_{b+1}) \geq V_{b+1}^*(h_{b+1});$$

then we show it results in

$$u_b(h_b) \geq V_b^*(h_b)$$

We have

$$\begin{aligned}
 u_b(h_t) &= \max_{a \in A} \left(\bar{r}_b(x_t, a) + \sum_{x \in \mathcal{X}} p_t(x_{t+1} = x | x_t, a) u_{t+1}(h_{t+1}, a, x) \right) \\
 &\geq \max_{a \in A} \left(\bar{r}_b(x_t, a) + \sum_{x \in \mathcal{X}} p_t(x_{t+1} = x | x_t, a) V_{t+1}^\pi(h_{t+1}, a, x) \right) \\
 &\geq \sum_{a \in A} \pi(a | h_t) \left[\bar{r}_b(x_t, a) + \sum_{x \in \mathcal{X}} p_b(x_{t+1} = x | x_t, a) V_{t+1}^\pi(h_{t+1}, a, x) \right] \\
 &= V_t^\pi(h_t)
 \end{aligned}$$

what we showed is

$$\begin{aligned}
 u_b(h_t) &\geq V_t^\pi(h_t) \text{ for all } h_t \text{ and } \pi \in \Pi^{\text{PR}} \\
 \Rightarrow u_b(h_t) &\geq V_t^*(h_t).
 \end{aligned}$$

we just proved that following a deterministic policy which gives us u_b has value at least as good as the V_b^* . \Rightarrow Therefore, this policy is optimal and $V_b^* = u_b$

Let's show $V_t^*(h_t)$ depends just on x_t , not the whole h_t . \Rightarrow we have

$$V_t^*(h_t) = \max_{a \in \mathcal{A}} \left(\bar{r}_t(x_t, a) + \sum_{x \in \mathcal{X}} P_{t+1}^t(x | x_t, a) V_{t+1}^*(h_{t+1}) \right)$$

we also know that $V_{t+1}^*(h_{t+1}) = \bar{r}_{t+1}(x_{t+1})$ is a function of h_{t+1} through x_{t+1} .

$$\text{so we write } V_{t+1}^*(h_{t+1}) = V_{t+1}^*(x_{t+1})$$

Induction:

$$\text{Let for all } h_{t+1}, V_{t+1}^*(h_{t+1}) = V_{t+1}^*(x_{t+1})$$

Then we show, same is true for V_t^*

$$\underline{V_t^*(h_t) = \max_{a \in \mathcal{A}} \left(\bar{r}_t(x_t, a) + \sum_{x \in \mathcal{X}} P_{t+1}^t(x | x_t, a) V_{t+1}^*(x) \right)}$$

The right hand side just depends on x_t in h_t .

Therefore, V_t^* is a function of x_t .
no matter what is h_t as a whole.

The action we choose at x_t is

$$a_t \in \operatorname{argmax}_{a \in A} \left(\bar{r}_t(x_t, a) + \sum_{x' \in X} P(x' = x_{t+1} | x_t, a) \bar{V}_{t+1}^+(x') \right)$$

which is a function of just x_t .

Therefore, there exist an optimal policy

which is $\left\{ \begin{array}{l} \text{Deterministic} \\ \text{Markovian} \end{array} \right\} \in \Pi^{\text{MD}}$

(Note: These awesome results do not hold for general state-action spaces. In general the optimal policy may not even exist.)