

LECTURE 3 - DYNAMIC PROGRAMMING

GOAL: Given full dynamics of the MDP, find an optimal policy π^* .

POLICY EVALUATION.

Given a policy π and full knowledge of MDP Dynamics, we want to compute $V_\pi(s)$.

From Lecture 2, we have the Bellman Expectation

$$\text{Equation } V_\pi(s) = \mathbb{E} [R_{t+1} + \gamma V_\pi(s_{t+1})] =$$

$$\sum_a \pi(a|s) \left[R_s^a + \gamma \sum_{s'} p(s'|s,a) V_\pi(s') \right].$$

These form a system of $|S|$ linear equations. Existence/Uniqueness is guaranteed if either

$$\begin{cases} \gamma < 1 \text{ or} \\ \text{termination is guaranteed from all states by following } \pi. \end{cases}$$

the system is defined by $V^\# = R^\pi + \gamma P^\pi V^\pi$

where $V^\# \in \mathbb{R}^{|S|}$

$$R^\# \in \mathbb{R}^{|S|}, \quad R^\#_s = \mathbb{E}[R_{t+1}|s_t=s] = \sum_a \pi(a|s) R_s^a$$

$$P^\# \in \mathbb{R}^{|S| \times |S|}, \quad P^\#_{ss'} = \mathbb{P}(s_{t+1}=s'|s_t=s) = \sum_a \pi(a|s) P_{ss'}$$

the solution $V^\#$ is a fixed point of the function $T(V) = R^\# + \gamma P^\# V : \mathbb{R}^{|S|} \rightarrow \mathbb{R}^{|S|}$.

We can find it by iterative methods. We define the sequence $V_\pi^{(k)} = T(V_\pi^{(k-1)})$

$$V_\pi^{(0)} = 0$$

Under the assumptions we mentioned,

$$V_\pi^{(k)} \xrightarrow{k} V_\pi.$$

$$V_{\pi}^k \rightarrow V_{\pi}.$$

$$V_{\pi}^k(s) = \sum_a \pi(a|s) \left[R_s + \gamma \sum_{s'} P_{ss'}^a \cdot V_{\pi}^{k-1}(s') \right]$$

- WHEN IMPLEMENTING THE METHOD, WE CAN USE 2 ARRAYS; ONE FOR OLD VALUE FUNCTIONS AND ONE FOR NEW VALUES, OR WE CAN DO IT IN PLACE USING ONLY ONE ARRAY. THIS STILL CONVERGES.

- COMMON STOPPING CRITERIA: $\max |V_k(s) - V_{k+1}(s)| < \epsilon$.

POLICY IMPROVEMENT

GOAL: Given a policy π w/ value V_{π} , can we define a new policy π' st $V_{\pi'} \geq V_{\pi}$?

POLICY IMPROVEMENT THEOREM: Given 2 policies π, π' if $Q_{\pi}(s, \pi'(s)) \geq V_{\pi}(s) \quad \forall s \in S \Rightarrow V_{\pi'}(s) \geq V_{\pi}(s) \quad \forall s \in S$.

COROLLARY: Given π , consider DETERMINISTIC π' defined by $\pi'(s) = \operatorname{argmax}_{\pi} Q_{\pi}(s, \pi) = \operatorname{argmax}_{\pi} \mathbb{E}[R_{t+1} + \gamma V_{\pi}(S_{t+1}) \mid S_t = s, A_t = a]$. satisfies THM because given any state $s \in S$,

$$Q_{\pi}(s, \pi'(s)) = Q_{\pi}(s, \operatorname{argmax}_{\pi} Q_{\pi}(s, \pi)) =$$

$$\max_{\pi} Q_{\pi}(s, \pi) \geq Q_{\pi}(s, \pi'(s)) = V_{\pi}(s).$$

$$\text{So } \pi'(s) = \text{Greedy } (\pi) \geq \pi.$$

PROOF THM Given any state $s, \dots \Rightarrow V_{\pi}(s) \leq V_{\pi'}(s)$

$$V_{\pi}(s) = Q_{\pi}(s, \pi(s)) \leq Q_{\pi}(s, \pi'(s)) =$$

$$\mathbb{E}_{\pi} \left[R_{t+1} + \gamma V_{\pi}(S_{t+1}) \mid S_t = s \right] \leq$$

$$\mathbb{E}_{\pi'} \left[R_{t+1} + \gamma Q_{\pi'}(S_{t+1}, \pi'(S_{t+1})) \mid S_t = s \right] =$$

$$\mathbb{E}_{\pi^*} [R_{T+2} + \gamma \mathbb{E}_{\pi^*} [R_{T+2} + \gamma V_{\pi^*}(S_{T+2}) \mid S_T = s]] =$$

$$\mathbb{E}_{\pi^*} [R_{T+1} + \gamma R_{T+2} + \gamma^2 V_{\pi^*}(S_{T+2}) \mid S_T = s]$$

$$= \dots = \mathbb{E}_{\pi^*} [R_{T+1} + \gamma R_{T+2} + \gamma^2 R_{T+3} + \dots \mid S_T = s]$$

$$= V_{\pi^*}(s) \quad \square$$

- Stopping Condition? What if $\pi' = \pi$?

Then: If $\pi' = \text{Greedy}(\pi) = \pi \Rightarrow \pi' = \pi = \pi^*$.

Proof: $\pi' = \pi^* \Leftrightarrow \pi'$ satisfies Bellman optimality Eq.

We want to prove that:

$$V_{\pi'}(s) = \max_a \mathbb{E} [R_{T+2} + \gamma V_{\pi'}(S_{T+2}) \mid S_T = s, A_T = a]$$

$$\bullet V_{\pi'}(s) = Q_{\pi'}(s, \pi'(s)) = \mathbb{E} [R_{T+2} + \gamma V_{\pi'}(S_{T+2}) \mid S_T = s, A_T = \pi'(s)]$$

$$\stackrel{\downarrow}{=} \mathbb{E} [R_{T+2} + \gamma V_{\pi'}(S_{T+2}) \mid S_T = s, A_T = \pi'(s)] = Q_{\pi'}(s, \pi'(s)) =$$

$$\begin{aligned} \pi &= \pi' \Rightarrow \\ V_{\pi} &= V_{\pi'} \end{aligned}$$

$$\max_a Q_{\pi}(s, a) = \max_{a=\pi'} Q_{\pi'}(s, a) =$$

$$\max_a \mathbb{E} [R_{T+2} + \gamma V_{\pi'}(S_{T+2}) \mid S_T = s, A_T = a] \quad \square$$

POLICY ITERATION.

$$\pi_0 \xrightarrow{\mathbb{E}} V_{\pi_0} \xrightarrow{\mathbb{E}} \pi_1 \xrightarrow{\mathbb{E}} V_{\pi_1} \xrightarrow{\mathbb{E}} \pi_2 \xrightarrow{\dots} \xrightarrow{\mathbb{E}} \pi_* \xrightarrow{\mathbb{E}} V_*$$

VALUE ITERATION.

- ONE DRAWBACK OF POLICY ITERATION IS the Amount of COMPUTATION Done Every time we EVALUATE a POLICY.
- CONVERGENCE CAN STILL HOLD EVEN IF POLY-EVAL RUNS FOR ONLY K ITERATIONS
- $K=1$ IS CALLED VALUE ITERATION.

- When we perform policy evaluation for the k^{th} time,
 $V_k(s) = \mathbb{E}[R_{t+1} + \gamma V_{k-1}(s_{t+1}) | s_t = s]$ using the
 Bellman Eq. Ex. V_{k-1} will correspond to the VALUE
 FUNCTION S.T $\pi_k = \text{SReedy}(V_{k-1})$.

$$\pi_0 \xrightarrow{\mathbb{E}} v_0 \xrightarrow{\mathbb{E}} \pi_1 \xrightarrow{\mathbb{E}} v_1$$

↳ SReedy(v_0)

$$\Rightarrow V_k(s) = \mathbb{E}[R_{t+1} + \gamma V_{k-1}(s_{t+1}) | s_t = s, A_t = \pi_k(s)]$$

$$= \max_a \mathbb{E}[R_{t+1} + \gamma V_{k-1}(s_{t+1}) | s_t = s, A_t = a]$$

$$\pi_k(s) = \underset{a}{\max} Q(s, a) = \underset{Q}{\max} \mathbb{E}[R_{t+1} + \gamma V_{k-1}(s_{t+1}) | s_t = s, A_t = a]$$

$$\therefore V_k(s) = \max_a \mathbb{E}[R_{t+1} + \gamma V_{k-1}(s_{t+1}) | s_t = s, A_t = a]$$

- Another way to view this update, is as iterating over the Bellman optimality equation in order to find the fixed point.

EXAMPLE 4.3) the GAMBLERS PROBLEM.

- Gambler makes a STAKE and coin is FLIPPED.
 - $H \Rightarrow$ He wins as many \$ he Bet on that FLIP
 - $T \Rightarrow$ Loses everything " "
- $\begin{cases} \text{He Reaches \$100} \\ \text{Runs out of money} \end{cases} \Rightarrow \text{some ENDS.}$

FORMULATE THIS PROBLEM AS MDP

- $S = \text{Gambler's capital } \in \{1, \dots, 99\} \quad S' = \{0, 100\}$

$A = \text{amount of money he Bets on FLIP} = \{0, 1, \dots, \min\{S, 100-S\}\}$

- S = Gambler's Capital $\in \{0, 1, \dots, 100\}$
- A_t = Amount of money he bets on flip = $\{0, 1, \dots, \min\{S, 100-S\}\}$
- $P = P(S'|R|S, A)$ OR $P(S'|S, A)$ OR
 $\mathbb{E}[R_{t+1} | S_t=S, A_t=A]$.

Rewards: 0 for all transitions, except for 1 when we reach the terminal state $S=100$.

$V_H(S) = \mathbb{E} \left[\underbrace{R_{t+1} + \gamma R_{t+2} + \dots}_{\text{Bell(P)}} | S_t=S \right] = \text{PROBABILITY OF WINNING FROM STATE } S$

π^* : Policy that maximizes the probability of win from every state

If $P_H = \text{PROB HEADS IS known}$, we can solve the MDP w/ Dynamic Programming -

ASYNCHRONOUS DYNAMIC PROGRAMMING.

ONE drawback to DP methods comes when $|S|$ is huge. For ex, before 1 step of value iteration, we need to perform a sweep over the entire state space S .

ONE IDEA IS TO IMPLEMENT IN-PLACE ALGORITHMS.
 WE UPDATE THE VALUE OF A STATE AS SOON AS IT IS READY.
 WE CAN ALSO CHOOSE WHICH STATES TO UPDATE THE MOST ACCORDING TO THEIR IMPORTANCE.