

RL Week - 3

L1 Bellman Equations

- Value of state \rightarrow expected long-term reward starting from that state, depends on the agent's policy.

- Bellman Eqⁿ for policy π

$$\begin{aligned}
 v_{\pi}(s) &= E_{\pi}[G_t \mid S_t = s] \\
 &= E_{\pi}[R_{t+1} + \gamma G_{t+1} \mid S_t = s] \\
 &= \sum_a \pi(a \mid s) \sum_{s'} \sum_a p(s', a \mid s, a) \\
 &\quad [r + \gamma E_{\pi}[G_{t+1} \mid S_{t+1} = s']] \\
 &= \sum_a \pi(a \mid s) \sum_{s', a} p(s', a \mid s, a) \\
 &\quad [r + \gamma v_{\pi}(s')] \quad \forall s \in S
 \end{aligned}$$

\rightarrow linear eqⁿ in $|S|$ variables

\rightarrow unique solⁿ exists

\rightarrow for the first reward, we calculate the probability, but from next step onwards, take the expected value into account.

→ we assume things markov & stationary

Eg: Ac[~] → N, S, W, E . deterministic . If would take agent off the grid . No movement but reward = 0. Other ac[~] produce reward = 0 , except sp. ac[~].

State value func[~] for equiprobable random policy . $\gamma = 0.9$

L2 Bellman Optimality Equa[~]s

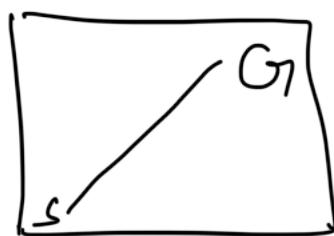
- for finite MDPs , policies can be partially ordered .
 $\pi \geq \pi'$ iff . $v_{\pi}(s) \geq v_{\pi'}(s) \quad \forall s \in S$
- There is always at least 1 (or many) policies that is better than or equal to all the others . This will be our optimal policy . That is π^* . ($\pi^* \geq \pi$)
- Optimal policies share the same v^*
Optimal state - value func[~].

$$v^*(s) = \max_{\pi} v_{\pi}(s) \quad \forall s \in S$$

$$q^*(s, a) = \max_{\pi} q_{\pi}(s, a) \quad \forall s \in S, a \in A$$

- there is 1 deterministic policy (at least 1) out of all policies.

Eg. Grid world setting



$$M = \langle S, A, p, r \rangle$$

N
W E
S

- many policies, but at least 1 optimal policy.

- Bellman Optimality eq^n for v_*

$$v_*(s) = \max_{a \in A(s)} q_{\pi_*}(s, a)$$

$$= \max_a E_{\pi_*} [G_t | s_t = s, a_t = a]$$

$$= \max_a E_{\pi_*} [R_{t+1} + \underbrace{\gamma G_{t+1}}_{\gamma v_*(s_{t+1})} | s_t = s, a_t = a]$$

$$v_*(s) = \max_a \sum_{s', a'} p(s', a | s, a) [a + \gamma v_*(s')]$$

- The value of a state under an optimal policy must equal the expected return for the best action from that state.

$$q_*(s, a) = E[R_{t+1} + \gamma \max_{a'} q_*(s_{t+1}, a') | s_t = s, a_t = a]$$

$$= \sum_{s', a} p(s', a | s, a) \{ u + \gamma \max_{a'} q^*(s', a') \}$$

- Any policy that is greedy w.r.t v^* is an optimal policy. \therefore , given v^* , 1-step look ahead search produces the long term optimal actions.

grid world v^*

$$\begin{array}{cc} 22.0 & 22.4 \\ 19.8 & - \end{array} \Rightarrow \begin{array}{c} \cdot \\ \cdot \\ \cdot \end{array}$$

- given v^* , the agent does not have to do a 1-step ahead search.

$$\pi^*(s) = \arg \max_{a \in A(s)} q^*(s, a)$$

L3 DP Policy Iteration

- Stochastic DP is the sole "method of choice" for MDPs (greedy, conquer & divide rule)

1. req. comp. knowledge of system dynamics (transit matrix & reward)
2. computationally expensive

"value" v_{π} q_{π}

control v^* q^*

3. curse of dimensionality
 4. guaranteed to converge

- Given a policy π , compute the state value funcⁿ v_π .

$$v_\pi(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma v_\pi(s')]$$

→ we iteratively solve it. You will have $|S|$ size of variables. $|S|$ linear equations. There will be a single unique soln. So, these are simultaneous linear equations.

$$v_0 = 0 \leftarrow s$$

$$v_{k+1}(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma v_k(s')] \leftarrow$$

similarly, calculate v_{k+2} using v_{k+1} . This iteration solves the problem gives approximate results.

Input π

$\theta > 0 \rightarrow$ algo param

$V(s)$ initally, $V(\text{terminal}) = 0$

$\Delta = 0$

(works in place)

loop $s \in S$
 $v \leftarrow V(s)$

$$V(s) \leftarrow \sum_a \pi(a|s) \sum_{s', a} p(s'|a|s, a)$$

$$[v + \gamma V(s')]$$

$$\Delta \leftarrow \max(D, |v - V(s)|)$$

until $\Delta < 0$

L4 Policy Improvement

- Given a π , we can now find v_π .
- For a state s , is it possible to choose $a^* \neq \pi(s)$

Value of a in state s -

$$q_\pi(s, a) = E[R_{t+1} + \gamma v_\pi(s_{t+1}) | \begin{matrix} S_t = s, \\ A_t = a \end{matrix}]$$
$$= \sum_{s', a} p(s'|a|s, a) [r + \gamma v_\pi(s')]$$

- So, it is better to switch to a^* for state s iff. $q_\pi(s, a^*) > v_\pi(s)$

- Now we do this for all states, that gives new policy π' (greedy w.r.t v_{π})

$$\pi'(s) = \arg \max_a q_{\pi}(s, a)$$

$$= \sum_{s', a} p(s', a | s, a) [a + \gamma v_{\pi}(s')]$$

$\therefore v_{\pi'} \geq v_{\pi} \Rightarrow$ Policy Improvement
 $\forall s \in S$

- What if $v_{\pi'} = v_{\pi}$, then

$$v_{\pi'}(s) = \max_a \sum_{s', a} p(s', a | s, a) [a + \gamma v_{\pi}(s')]$$

↳ Bellman optimality eqⁿ again

$\therefore v_{\pi'} = v_{*}$ & both π & π' are the optimal policies.

$$\pi_0 \xrightarrow{E} v_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} v_{\pi_1} \xrightarrow{I} \pi_2 \dots$$

E → evaluate, I → improve
 ↳ quodlibet

Policy Iteration