

Dynamic Programming

Assumption

- env. is a finite MDP
- s, A, R
- dynamics : $p(s', r | s, a)$ is given for all $s \in S, a \in A(s), r \in R, s' \in S^+$

Key Idea

$$\begin{aligned} V_*(s) &= \max_a E[R_{t+1} + \gamma V_*(S_{t+1}) | S_t = s, A_t = a] \\ &= \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V_*(s')] \end{aligned}$$

Policy Evaluation (Prediction)

for arbitrary policy π , compute state-value function (evaluate)

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

Strategy : Iteration

$$\begin{aligned} V_{k+1}(s) &= E_\pi[R_{t+1} + \gamma V_k(S_{t+1}) | S_t = s] \\ &= \sum_a \pi(a | s) \sum_{s', r} p(s', r | s, a) [r + \gamma V_k(s')] \end{aligned}$$

iterative policy Iteration

Input : π , to be evaluated

$V_0(s) = 0, k = 0.$

Repeat

for each all $s \in S$

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

$$\Delta V = \max_s |V_{k+1}(s) - V_k(s)|$$

$$k = k + 1$$

until $\Delta V < \epsilon.$

$$V_\pi = V_{k+1}$$

Policy Improvement

How to improve policy? Try new policy and compare.

\therefore If $q_{\pi}(s, \pi'(s)) \geq v_{\pi}(s)$, $\pi(s) \rightarrow \pi'(s)$ is better.

(proof) $v_{\pi}(s) \leq q_{\pi}(s, \pi'(s))$

$$= E[R_{t+1} + \gamma v_{\pi}(S_{t+1}) \mid S_t = s, A_t = \pi'(s)]$$

$$= E_{\pi'}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) \mid S_t = s]$$

$$\leq E_{\pi'}[R_{t+1} + \gamma q_{\pi}(S_{t+1}, \pi'(S_{t+1})) \mid S_t = s]$$

$$= E_{\pi'}[R_{t+1} + \gamma E_{\pi'}[R_{t+2} + \gamma v_{\pi}(S_{t+2}) \mid S_{t+1}]] \mid S_t = s]$$

$$= E_{\pi'}[R_{t+1} + \gamma R_{t+2} + \gamma^2 v_{\pi}(S_{t+2}) \mid S_t = s]$$

$$\leq E_{\pi'}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 v_{\pi}(S_{t+3}) \mid S_t = s]$$

\vdots

$$\leq E_{\pi'}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \mid S_t = s]$$

$$= v_{\pi}(s)$$

Intuitively, if one step try other policy $\pi'(s)$ for all states s better than $\pi(s)$, π' is better for whole process.

$$\pi'(s) = \operatorname{argmax}_a q_{\pi}(s, a) \rightarrow \text{greedy policy (simple solution)}$$

$$= \operatorname{argmax}_a E[R_{t+1} + \gamma v_{\pi}(S_{t+1}) \mid S_t = s, A_t = a]$$

$$= \operatorname{argmax}_a \sum_{s', r} p(s', r \mid s, a) [r + \gamma v_{\pi}(s')]$$

Policy Iteration

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

greedy policy.

1. Initialization : $V(s)$, $\pi(s)$ arbitrarily

2. Policy Evaluation

$$\text{iterate } V_{k+1}(s) = \sum_{s', r} p(s', r \mid s, a = \pi(s)) [r + \gamma V_k(s')]$$

until $\Delta V < \epsilon$

3. Policy Improvement

$$\pi(s) = \operatorname{argmax}_a \sum_{s', r} p(s', r \mid s, a) [r + \gamma V(s')]$$

If $\pi(s) = \text{old policy}$, Stop

$$V \approx V_*, \pi \approx \pi_*$$

Value Iteration

One drawback of policy Iteration

- too many iteration in policy evaluation step
- especially, initial iterations meaningless because policy itself is meaningless

policy evaluation + policy improvement = value Iteration

$$V_{k+1}(s) = \max_a E[R_{t+1} + \gamma V_k(S_{t+1}) \mid S_t = s, A_t = a]$$
$$= \max_a \sum_{s', r} p(s', r \mid s, a) [r + \gamma V_k(s')]$$

Algorithm

$$V_0(s) = 0, \quad k = 0.$$

Repeat

For each $s \in S$:

$$V_{k+1}(s) = \max_a \sum_{s', r} p(s', r \mid s, a) [r + \gamma V_k(s')]$$

$$\Delta V = \max_s |V_{k+1}(s) - V_k(s)|, \quad k = k+1,$$

until $\Delta V < \epsilon$.

~~return~~

$$\pi_k(s) = \arg \max_a \sum_{s', r} p(s', r \mid s, a) [r + \gamma V_k(s')].$$