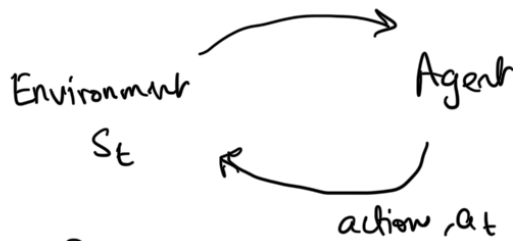


# Reinforcement Learning

MDP  $(S, A, P, R, \gamma, r/H)$

$$r_t \leftarrow s_t, a_t$$



$$s_1 \sim d_0 \quad s_{t+1} \leftarrow P(s_t, a_t)$$

Goal: find policy  $\pi : s \rightarrow a \quad E\left[\sum_{t=1}^{\infty} \gamma^{t-1} r_t \mid \pi, s\right]$

✓ Planning problem: Given state transition + reward models, find optimal policy

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

Bellman optimality equations, optimal  $Q^*, V^*$

$$V^*(s) = \max_{a \in A} [R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) V^*(s')] \equiv T_V V^* \leftarrow$$

$$Q^*(s, a) = R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) \underbrace{V^*(s')}_{\max_{a' \in A} Q^*(s', a')} \equiv T_Q Q^* \leftarrow$$

Challenge: Nonlinear equations (bcz of max)

① Dynamic Programming / Iterative solutions finite S, A

## Policy Iteration

Initialize  $\pi_0$

1) Policy evaluation  $Q_{\pi_{k-1}}$  (linear Bellman eq's)

2) Policy improvement  $\pi_k = \pi_{Q_{\pi_{k-1}}}$   $\leftarrow \Rightarrow A_{\pi_k}(s, \pi_{k-1}) \geq 0$

Guaranteed to improve monotonically. Hence terminate  $Q_{\pi_k} = Q_{\pi_{k-1}}$ .

$\rightarrow$  # policies =  $|A|^{|S|}$

P.I... Improvement Theorem: In policy iteration.  $V_{\pi_k}(s) \geq V_{\pi_{k-1}}(s)$ .

only hypothesis

$$\forall k \geq 1, s \in S, a \in A$$

& strictly positive in atleast 1 state until  $\pi^*$  is found.

Proof: Advantage of action  $a$  in state  $s$  over policy  $\pi$  is defined.

$$as \quad A_{\pi}(s, a) = Q_{\pi}(s, a) - Q_{\pi}(s, \underline{\pi}(s)) = Q_{\pi}(s, a) - V_{\pi}(s)$$

$$\text{Advantage of policy } \pi' \text{ over } \pi \quad A_{\pi}(s, \pi') := A_{\pi}(s, \pi'(s))$$

Performance Difference Lemma For any  $\pi, \pi' \in S$

$$V_{\pi'}(s) - V_{\pi}(s) = \frac{1}{1-\gamma} E_{s' \sim d_{\pi', s}} [A_{\pi}(s', \pi')] \leftarrow$$

↳ normalised discounted occupancy induced by  $\pi'$  starting in state  $s$ .

Proof: Consider a seq<sup>n</sup> of policies  $\{\pi_i\}_{i=0}^{\infty}$

$$\pi_0 = \pi$$

$$\pi_{\infty} = \pi'$$

$\pi_i$  follow  $\pi'$  for first  $i$  steps, then switch to  $\pi$ .

$$\begin{aligned} V_{\pi'}(s) - V_{\pi}(s) &= V_{\pi_{\infty}}(s) - V_{\pi_0}(s) \\ &= \sum_{i=0}^{\infty} (V_{\pi_{i+1}}(s) - V_{\pi_i}(s)) \\ &= \sum_{i=0}^{\infty} (E[\sum_{t=i}^{\infty} \gamma^{t-i} r_t | \pi_{i+1}, s_i = s] \\ &\quad - E[\sum_{t=i}^{\infty} \gamma^{t-i} r_t | \pi_i, s_i = s]) \end{aligned}$$

Note:  $\pi_{i+1} \neq \pi_i$  only deviate at step  $i+1$ .

Same roll-in policy  $\pi'$  for first  $i$  steps

⇒ defines roll-in distribution  $P(s_{i+1} | s_i = s, \pi')$

Same roll-out policy  $\pi$  starting at  $i+2$

⇒ conditioned on  $s_{i+1} = s$ ,  $a_{i+1} = a$  total expected discounted reward picked up in rest of trajectory is

$$\gamma^i Q_{\pi}(s, a) \leftarrow$$

↑ ↑

$$\begin{aligned} \Rightarrow V_{\pi'}(s) - V_{\pi}(s) &= \sum_{i=0}^{\infty} \gamma^i \sum_{s' \in S} P(s_{i+1}=s' | s_i=s, \pi') \\ &\quad \underbrace{(\underbrace{Q_{\pi'}(s', \pi'(s'))}_{\text{new}} - \underbrace{Q_{\pi}(s', \pi(s))}_{\text{old}}))}_{A_{\pi}(s', \pi')} \\ &= \sum_{s' \in S} \underbrace{P(s_{i+1}=s' | s_i=s, \pi')}_{\text{prob}} A_{\pi}(s', \pi') \cdot \frac{1}{1-\gamma} \\ &= E_{s' \sim \mathcal{D}_{\pi', s}} [A_{\pi}(s', \pi')] \frac{1}{1-\gamma} \end{aligned}$$

Proof of Policy improvement thm

Invoke policy diff lemma  $\pi' = \pi_{k+1} \neq \pi = \pi_k$ .

$$\Rightarrow V_{\pi_{k+1}} \geq V_{\pi_k} \quad \because A_{\pi_k}(s, \pi_{k+1}) \geq 0 \quad \forall s.$$

due to policy improvement step

Worst case  $|A|^{(s)}$  but in practice works better.

Often approximate sol<sup>n</sup> suffices.

Thm: Policy iteration converges exponentially in sup-norm

$$\|Q_{\pi_{k+1}} - Q^*\|_{\infty} \leq \gamma \|Q_{\pi_k} - Q^*\|_{\infty}$$

Contraction property of  $T_Q$ :

$T_Q$  is a  $\gamma$ -contraction in sup-norm

$$\|TQ - TQ'\|_{\infty} \leq \gamma \|Q - Q'\|_{\infty} \quad \forall Q, Q'$$

Proof:  $\|TQ - TQ'\|_{\infty} = \gamma \max_{s,a} \left| E_{s' \sim P(\cdot|s,a)} \left[ \underbrace{\max_{a'} Q(s', a')}_{V(s')} \right] \right|$

$$- E_{s' \sim P(\cdot|s,a)} \left[ \max_{a'} Q'(s', a') \right]$$

$$|E| \leq E|1| \leq \gamma \max_{s,a} E_{s' \sim P(\cdot|s,a)} \left| \max_{a'} Q(s', a') - \max_{a'} Q'(s', a') \right|$$

$$|\max| \leq \max|1| \leq \gamma \max_{s,a} E_{s' \sim P(\cdot|s,a)} \max_{a'} |Q(s', a') - Q'(s', a')|$$

$$\leq \gamma \|Q - Q'\|_\infty$$

Proof of exp conv of policy iteration.

$$Q_{\pi_{k+1}} \geq TQ_{\pi_k} \quad **$$

$$\text{Using this, } Q^* - Q_{\pi_{k+1}} \leq TQ^* - TQ_{\pi_k}$$

$$\Rightarrow \|Q^* - Q_{\pi_{k+1}}\| \leq \gamma \|Q^* - Q_{\pi_k}\|_\infty \quad \text{contraction of } T$$

Now lets prove \*\*

$$Q_{\pi_{k+1}}(s,a) = R(s,a) + \gamma E_{s' \sim P(\cdot|s,a)} V_{\pi_{k+1}}(s')$$

$$\geq R(s,a) + \gamma E_{s' \sim P(\cdot|s,a)} V_{\pi_k}(s') \quad \text{Policy improvement thm}$$

$$\geq R(s,a) + \gamma E_{s' \sim P(\cdot|s,a)} \max_{a'} Q_{\pi_k}(s',a') = TQ_{\pi_k}$$

### Value Iteration

Approx Q/V directly using a sep' of Q function (without going between Q &  $\pi$ )

Initialize  $Q_0$

$$Q_t = TQ_{t-1}$$

$$Q_{\pi_{k+1}} \geq TQ_{\pi_k} \quad (\text{policy iteration})$$

Contraction property implies  $Q$  converges to  $Q^*$  exponentially.

Translate to policy value error:

$$V^*(s) - V_{\pi_Q}(s) = Q^*(s, \pi^*(s)) - Q_{\pi_Q}(s, \pi_Q(s)) \quad \text{Note: } Q_{\pi_Q} \neq Q$$

$$= \underbrace{Q^*(s, \pi^*(s)) - Q(s, a)}_1 + \underbrace{Q(s, a) - Q^*(s, a)}_2$$

$$1 \leq Q^*(s, \pi^*(s)) - Q(s, \pi^*(s)) \quad \because a = \arg \max_{a'} Q(s, a')$$

$$\leq \|Q - Q^*\|_\infty$$

$$2 \leq \|Q - Q^*\|_\infty$$

$$3 \leq \gamma E_{s' \sim P(\cdot|s,a)} [V^*(s') - V_{\pi_Q}(s')]$$

$$\leq 2 \|Q - Q^*\|_\infty + \gamma E_{s' \sim P(\cdot|s,a)} [V^*(s') - V_{\pi_Q}(s')]$$

$$\Rightarrow \left( \frac{1}{1-r} \right) \|V^*(s) - V_{\pi_0}(s)\|_{\infty} \leq \frac{2\|Q^* - Q\|_{\infty}}{1-r} \leq \frac{2r^t R_{\max}}{1-r}$$

□

## ② Linear Programming

$$\begin{aligned} \min_V \quad & d_0^T V \\ \text{st. } \quad & V \geq TV \quad \checkmark \quad \quad V = [V(s)]_{s \in S} \end{aligned}$$

Separation constraints (nonlinear)

max  
a

↓  
SxA constraints (linear)