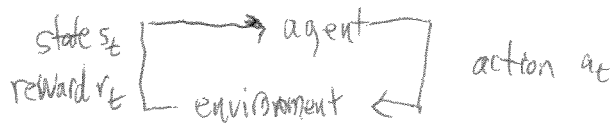


# Review

## • Reinforcement Learning



- Markov Decision Process (MDP) : assume a complete model of fully observed, stochastic environment

$$\text{MDP} = \{S, A, P(s'|s, a), R(s)\}$$

↑            ↓            ↓            ↘  
States    actions    transitions    rewards

- Policy:  $\pi: S \rightarrow A$

- Long-term discounted return =  $\sum_{t=0}^{\infty} \gamma^t r_t$

- State value function

$$V^{\pi}(s) = \mathbb{E}^{\pi} \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t) \mid s_0 = s \right]$$

$$V^{\pi}(s) = R(s) + \gamma \sum_{s'} P(s'|s, \pi(s)) V^{\pi}(s') \quad \star \text{Bellman eq.}$$

- Action value function

$$Q^{\pi}(s, a) = \mathbb{E}^{\pi} \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t) \mid s_0 = s, a_0 = a \right]$$

$$= R(s) + \gamma \sum_{s'} P(s'|s, a) V^{\pi}(s')$$

- Value functions  $\leftrightarrow$  policies

Given  $\pi^*(s)$ :

$$V^*(s) = R(s) + \gamma \sum_{s'} P(s'|s, \pi^*(s)) V^*(s') = \max_a Q^*(s, a) \quad \star$$

$$Q^*(s, a) = R(s) + \gamma \sum_{s'} P(s'|s, a) \cdot V^*(s')$$

Given  $V^*(s)$ ,  $Q^*(s, a)$ :

$$\begin{aligned}\pi^*(s) &= \operatorname{argmax}_a Q^*(s, a) \\ &= \operatorname{argmax}_a [R(s) + \gamma \sum_{s'} P(s'|s, a) \cdot V^*(s')] \\ &= \operatorname{argmax}_a \left[ \sum_{s'} P(s'|s, a) \cdot V^*(s') \right]\end{aligned}$$

• Planning under uncertainty

How to compute  $\pi^*(s)$ , or equivalently,  $V^*(s)$  or  $Q^*(s, a)$ ?

- 1) Policy evaluation: compute  $V^\pi(s)$
- 2) Policy improvement: compute  $\pi'$  s.t.  $V^{\pi'}(s) \geq V^\pi(s) \quad \forall s$
- 3) Policy iteration: compute  $\pi^*$
- 4) Value iteration: compute  $V^*(s)$  directly

1) Policy evaluation

System of  $n$  linear equations for  $n$  unknowns

$$\sum_{s'} [I(s, s') - \gamma P(s'|s, \pi(s))] \cdot V^\pi(s') = R(s) \quad s=1 \dots n$$

$$\rightarrow (I - \gamma P)V = R \quad \text{ex} \quad \left[ \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} - \gamma \begin{pmatrix} P_{00} & P_{01} \\ P_{10} & P_{11} \end{pmatrix} \right] \begin{bmatrix} v_0 \\ v_1 \end{bmatrix} = \begin{bmatrix} r_0 \\ r_1 \end{bmatrix}$$

$$\Rightarrow V^\pi = (I - \gamma P^\pi)^{-1} R$$

## 2) Policy improvement

- How to compute  $\pi'$  such that  $V^{\pi'}(s) \geq V^{\pi}(s)$  for all states  $s$ ?
- Recall  $Q^{\pi}(s,a)$  = "expected return from state  $s$ , follow action  $a$ , then follow policy  $\pi$ "

How to compute  $Q^{\pi}(s,a)$ ?

- Evaluate policy to get  $V^{\pi}(s)$
- $Q^{\pi}(s,a) = R(s) + \gamma \sum_{s'} P(s'|s,a) V^{\pi}(s')$

Define "greedy policy":

$$\begin{aligned}\pi'(s) &= \underset{a}{\operatorname{argmax}} Q^{\pi}(s,a) \\ &= \underset{a}{\operatorname{argmax}} R(s) + \gamma \sum_{s'} P(s'|s,a) V^{\pi}(s') \\ &= \underset{a}{\operatorname{argmax}} \left[ \sum_{s'} P(s'|s,a) V^{\pi}(s') \right]\end{aligned}$$

Theorem: greedy policy  $\pi'$  everywhere performs better or equal to original policy  $\pi$ .

$$V^{\pi'}(s) \geq V^{\pi}(s) \quad \text{for all } s.$$

Intuition: ~~greedy policy  $\pi'$  everywhere performs better or equal to original policy  $\pi$~~

if better to choose action  $a$  in state  $s$ , then follow  $\pi$ ,  
it's always better to choose action  $a$  in state  $s$ .

Proof: 
$$\begin{aligned}
 V^\pi(s) &= Q^\pi(s, \pi(s)) \\
 &\leq \max_a Q^\pi(s, a) \\
 &= Q^\pi(s, \pi'(s)) \quad \text{by definition of greedy policy} \\
 &= R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s)) \cdot V^\pi(s')
 \end{aligned}$$

So far, it is better to take one step under  $\pi'$ , then revert to  $\pi$ , than to follow  $\pi$ .

"one-step inequality": 
$$V^\pi(s) \leq R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s)) V^\pi(s')$$

Apply inequality to  $V^\pi(s')$  on RHS:  $\swarrow V^\pi(s')$

$$V^\pi(s) \leq R(s) + \gamma \sum_{s'} P(s'|s, \pi'(s)) \cdot \left[ R(s') + \gamma \sum_{s''} P(s''|s', \pi'(s')) \cdot V^\pi(s'') \right]$$

So, better to take two steps under  $\pi'$ , then revert to  $\pi$ , than to always follow  $\pi$ .

In general, apply "one-step inequality"  $t$  times, we will show:  
 better to take  $t+1$  steps under  $\pi'$ , then follow  $\pi$ ,  
 than to always follow  $\pi$ .

Let  $t \rightarrow \infty$ , it's always better to follow  $\pi'(s)$  than  $\pi(s)$ .

$$\Rightarrow V^\pi(s) \leq V^{\pi'}(s) \quad \text{since RHS converges to } V^{\pi'}(s) \text{ for } \gamma < 1.$$

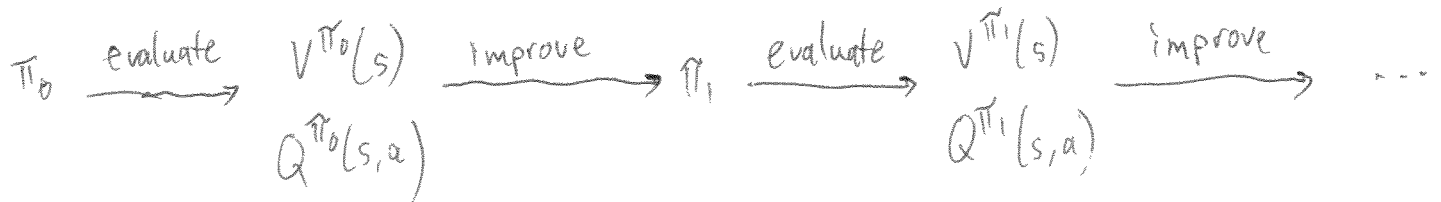
3) Policy iteration: how to compute  $\pi^*(s)$ ?

Algorithm:

(1) initialize policy  $\pi_0$  at random

(2) repeat until convergence:

- compute state + action value function of current policy (policy evaluation)
- derive greedy policy from action value function (policy improvement)



• Two questions: (i) does policy iteration converge at all?  
(ii) if so, what does it converge to?

(i) is policy iteration guaranteed to converge? Yes.

why? 1) cannot cycle b/c policy improve ~~at~~ at every iteration  
2) # policies is finite:  $|A|^{|S|}$

Typically converges in far fewer steps than  $|A|^{|S|}$ .

(ii) Does it always converge to an optimal policy  $\pi^*(s)$ ? Yes.

Thm: Suppose that we've converged at iteration  $k$ :  $V^{\pi_{k+1}}(s) = V^{\pi_k}(s)$   
Then  $V^{\pi_k}(s) = V^*(s)$ .

(Note: optimal value function is unique, even if there are many optimal policies.)

• Proof strategy:

1) Derive "Bellman optimality eqn" satisfied by  $V^{\pi_k}(s)$  at convergence.

2) Show that  $V^{\pi_k}(s) \geq V^{\tilde{\pi}}(s)$  for all states  $s$  and other policies  $\tilde{\pi}$ .

Hence,  $V^{\pi_k}(s) = V^*(s)$ .

## Step 1

From Bellman eqn for  $\pi_{k+1}(s)$ :

$$V^{\pi_{k+1}}(s) = R(s) + \gamma \sum_{s'} P(s'|s, \pi_{k+1}(s)) V^{\pi_{k+1}}(s')$$

By assumption:  $V^{\pi_{k+1}}(s) = V^{\pi_k}(s)$  at convergence.

$$V^{\pi_k}(s) = R(s) + \gamma \sum_{s'} P(s'|s, \pi_{k+1}(s)) V^{\pi_k}(s')$$

By assumption:  $\pi_{k+1}(s)$  is greedy w.r.t.  $V^{\pi_k}(s)$

Hence: 
$$\underline{V^{\pi_k}(s) = R(s) + \gamma \max_a \sum_{s'} P(s'|s, a) V^{\pi_k}(s')} \quad \text{"Bellman optimality eqn"}$$

(set of  $n$  non-linear eqns for  $s=1, 2, \dots, n$ ;  
non-linear b/c of max operation)

## Step 2

Iterate RHS

$\swarrow V^{\pi_k}(s')$

$$V^{\pi_k}(s) = R(s) + \gamma \max_a \sum_{s'} P(s'|s, a) \left[ R(s') + \gamma \max_{a'} \sum_{s''} P(s''|s', a') V^{\pi_k}(s'') \right]$$

denote  $\textcircled{\text{RHS}}$  for RHS

Iterate again + again.

Now show that this iterated expression (taken out to an infinite # of terms) implies optimality.

Let  $\tilde{\pi}(s)$  be any other policy.

From Bellman eqn:

$$\begin{aligned}
 V^{\tilde{\pi}}(s) &= R(s) + \gamma \sum_{s'} P(s'|s, \tilde{\pi}(s)) V^{\tilde{\pi}}(s) \\
 &\leq R(s) + \gamma \max_a \sum_{s'} P(s'|s, a) V^{\tilde{\pi}}(s) \quad \left. \begin{array}{l} \text{be greedy} \\ \text{use Bellman eq.} \end{array} \right\} \\
 &= R(s) + \gamma \max_a \sum_{s'} P(s'|s, a) \left[ R(s') + \gamma \sum_{s''} P(s''|s', \tilde{\pi}(s')) V^{\tilde{\pi}}(s'') \right] \quad \left. \begin{array}{l} \text{be greedy} \\ \text{use Bellman eq.} \end{array} \right\} \\
 &\leq R(s) + \gamma \max_a \sum_{s'} P(s'|s, a) \left[ R(s') + \gamma \max_{a'} \sum_{s''} P(s''|s', a') V^{\tilde{\pi}}(s'') \right] \\
 &\quad \text{denote } \boxed{\$} \text{ for RHS}
 \end{aligned}$$

Consider upper bound on  $V^{\tilde{\pi}}(s)$  from iterating above  $t$  times.

Compare to equality for  $V^{\pi^k}(s)$  after iterating  $t$  times.

As  $t \rightarrow \infty$ ,  $\boxed{\$}$  converges to  $\boxed{@}$ .

upper bound  
on  $V^{\tilde{\pi}}(s)$

RHS of equality for  $V^{\pi^k}(s)$

Thus, as  $t \rightarrow \infty$ :

$$V^{\tilde{\pi}}(s) \leq \lim_{t \rightarrow \infty} \boxed{\$} = \lim_{t \rightarrow \infty} \boxed{@} = V^{\pi^k}(s).$$

Thus for all policies  $\tilde{\pi}$  and states  $s$ , we have  $V^{\tilde{\pi}}(s) \leq V^{\pi^k}(s)$ ,

$$\text{or: } V^{\pi^k}(s) = \max_{\tilde{\pi}} V^{\tilde{\pi}}(s) = V^*(s).$$

To compute  $\pi^*$ :  $\pi^*(s) = \operatorname{argmax}_a Q^*(s, a) = \operatorname{argmax}_a \sum_{s'} P(s'|s, a) V^*(s)$ .

Pros/cons of policy iteration:

(+): converges quickly (in handful of steps often)

(-): each iteration requires  $O(n^3)$  policy evaluation