

Review  $(\pi)^T V \leftarrow \text{actions}$   $\pi = \text{actions} \quad (\pi)^T V \leftarrow \text{actions}$

\* Markov decision process (MDP)

$$\{ \mathcal{S}, \mathcal{A}, P(s'|s, a), R(s), \gamma \}$$

states, actions, transitions, rewards, discount factor

\* Policy mapping  $\pi: \mathcal{S} \rightarrow \mathcal{A}$  of states to actions

\* Value functions not in responses yet

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

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

\* Planning - given parameters of MDP what can we compute?

1. Policy evaluation - how to compute  $V^\pi(s)$ ?

Solve linear equations:

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

2. Policy improvement - how to improve on  $\pi(s)$ ?

Greedy policy  $\pi'(s) = \arg \max Q^\pi(s, a)$

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

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

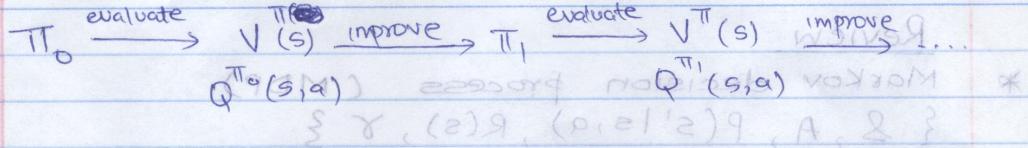
Algorithm [NOT the algorithm in HW]

(1) initialize policy at random

(2) repeat until convergence

- compute state & action value function of current policy.

- derive greedy policy from action value function.



\* Is policy iteration guaranteed to converge? Yes  
why? 1) # policies is finite  $|A|^{|\mathcal{S}|}$

2) cannot cycle b/c policy improvements at every iteration.

Typically converges in far less steps than  $|A|^{|\mathcal{S}|}$

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

\* thm: suppose  $\pi(s) = \pi^*(s)$  for all states  $s$ , or even more generally, that  $V^\pi(s) = V^*(s)$ . Then  $V^\pi(s) = V^*(s)$  w.r.t. no-tautous policy.

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

\* Proof strategy: w.r.t.  $\pi^*$   
1) Derive "Bellman optimality eqn" satisfied by  $V^\pi(s)$  at convergence.  
2) Show that  $V^\pi(s) \geq V^*(s)$  for all states  $s$  and other policies  $\pi$ . Hence  $V^\pi(s) = V^*(s)$

(2) \*  $\pi$  shows w.r.t.  $\pi^*$

From Bellman eqn for  $\pi^*(s)$  initial (1)  

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

To obtain  $\pi$  shows  $\pi^*(s)$  -

By assumption,  $V^{\pi}(s) = V^*(s)$  at convergence

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

By assumption,  $\pi^*(s)$  is greedy w.r.t.  $V^\pi(s)$   
Hence:

$$V^\pi(s) = R(s) + \gamma \max \sum p(s'|s, a) V^\pi(s')$$

at step  $t$ ,  $V^\pi(s)$  no longer s.t.  $s' \in S$ ,  $a' \in A$  & different than linear Bellman equation.

"Bellman optimality equation" to ENH equation.  
(set of  $n$  non-linear eqns for  $s = 1, 2, \dots, n$ )

(2) non-linear b/c of max operation

Step 2

Iterate RHS b/c  $V^\pi(s)$

$$V^\pi(s) = R(s) + \gamma \max \left( \sum p(s'|s, a) [R(s') + \right.$$

$$(2)^* V = (2)^* a \quad x_{\text{opt}}^{s'} = (2)^* V$$

$$\left. \gamma \max P(s''|s', a') V^\pi(s'') \right]$$

(note:  $a'$  for action  $a$ )

( $\epsilon$ -optimal vs. non-optimal  $V^\pi(s')$ )

Iterate again and again:

$$V^\pi(s) = R(s) + \gamma \max \sum_a \sum_{s'} p(s'|s, a) \left[ R(s') + \gamma \max \sum_{s''} p(s''|s', a') V^\pi(s'') \right]$$

$$(2)^* V = (2)^* a \quad x_{\text{opt}}^{s'} = (2)^* V$$

Now show that this iterated expression (taken out to infinity terms) implies optimality.

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

From Bellman's eqn:  $V^\pi(s) = R(s) + \gamma \sum_{s'} p(s'|s, \pi(s)) V^\pi(s')$

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

$$\leq R(s) + \gamma \max_a \sum_{s'} p(s'|s, a) V^{\tilde{\pi}}(s')$$

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

$$\leq R(s) + \gamma \max_a \sum_{s'} p(s'|s, a) [R(s') + \gamma \max_{a'} \sum_{s''} p(s''|s', a') V^{\tilde{\pi}}(s'')]$$

be greedy

(e)  $V^{\pi}$  considers upper bound  $\leq$  on (a)  $V^{\pi}(s)$  from iterating above  $t$  times.

(e) compare to (a) equality for  $V^{\pi}(s)$  after iterating  $t$  times.

As  $t \rightarrow \infty$ , RHS on upper bound on  $V^{\pi}(s)$  converges to RHS of "equality for  $V^{\pi}(s)$ "

Thus as  $t \rightarrow \infty$ ,

$$V^{\pi}(s) \leq \lim_{t \rightarrow \infty} [q_0 \cdot x_0] = \lim_{t \rightarrow \infty} [p/d] = V^{\pi}(s)$$

(e)

thus for all policies  $\pi$  and states  $s$

$$+ (e) we have  $V^{\pi}(s) \leq p_m \pi(s) \quad (e) \Rightarrow (e) V^{\pi}(s)$$$

$$\text{or: } V^{\pi}(s) = \max_a V^{\pi}(s) = V^*(s)$$

$(e) V^{\pi}(s) \leq V^*(s)$

Pros / cons of policy iteration:

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

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

: more time per iteration

(b) Value iteration - how to compute  $V^*(s)$  directly?

$$V^*(s) = \max_a Q^*(s, a)$$

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

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

\*  $n$  non-linear eqns for  $s = 1, 2, \dots, n$

$n$  unknowns  $V^*(s)$

How to solve?

$$(e) V^*(s) \geq R(s) + \gamma \max_a Q^*(s, a) \geq$$

$$[(e) V^*(s) \geq R(s) + \gamma \max_a Q^*(s, a)] \geq [(e) V^*(s) \geq R(s) + \gamma \max_a Q^*(s, a)]$$

$$[(e) V^*(s) \geq R(s) + \gamma \max_a Q^*(s, a)] \geq [(e) V^*(s) \geq R(s) + \gamma \max_a Q^*(s, a)]$$

iterate

Algorithm : value iteration

(1) initialize  $V_0(s) = 0$  for all states  $s$

(2) iterate

$$V_{k+1}(s) = R(s) + \gamma \max_a \left[ \sum_{s'} P(s'|s,a) V_k(s') \right]$$

for all  $s = 1, 2, \dots, n$

Note : this algorithm works directly on value functions  
not policies.

But incremental policies can be computed from :

$$\pi_{k+1}(s) = \text{greedy} [ V_k(s) ] = \arg \max_a \left[ \sum_{s'} P(s'|s,a) V_k(s') \right]$$

(3) suppose this converges :

$$\lim_{k \rightarrow \infty} V_k(s) \rightarrow V^*(s)$$

$$\text{compute } \pi^*(s) = \arg \max_a \left[ \sum_{s'} P(s'|s,a) V^*(s') \right]$$

Does algorithm converge ? Yes.