cs 4/5789

2 Feb 2012 Prot Sarah Dean

Lecture 4: Value & Policy Heration, Dynamic Programming

1) Value Heration

Recall algorithm, contraction proof from last lecture. Setting of the JQt results in approximately optimal a function! optimal

10t-0*115 8t 1100-0*112

From a functions to policies

We know $T^*(s) = \underset{\alpha}{\text{argmax}} Q^*(s, \alpha)$

since Qt(s,a) & Qt(s,a) during value iteration,

 $TT(s) = argmax Q^t(s,a)$

a good choice?

Theorem: The quality of Tit is bounded below:

 $V^{Ht}(s) \geq V^{*}(s) - \frac{28^{t}}{1-8} \|Q^{0} - Q^{*}\|_{\infty} \quad \forall s \in S$

'Proof:

Assume the following claim is true:

 $V^{\pi t}(s) - V^{*}(s) \ge Y \mathbb{E} \left[V^{\pi t}(s') - V^{*}(s') \right] - 2Y^{t} \| Q^{\circ} - Q^{*} \|_{\infty}$

Then recursing k times,

$$V^{\dagger t}(s) - V^{*}(s) \geq 8^{k} \left[\sum_{s \in P(S)} v^{\dagger t}(s') - V^{*}(s) \right] - 2 \sum_{k=0}^{k} 8^{k+t} ||Q^{0} - Q^{*}||_{\infty}$$

White thing $k \rightarrow \infty$,

$$V^{\dagger t}(s) - V^{*}(s) \geq -28^{t} \sum_{k=0}^{\infty} 8^{k} ||Q^{0} - Q^{*}||_{\infty}$$

$$= -28^{t} ||Q^{0} - Q^{*}||_{\infty}$$

$$= -\frac{2}{1-8} \|Q^{0} - Q^{*}\|_{\infty}$$

$$V^{\text{rt}}(S) - V^{*}(S) = Q^{\text{tt}}(S, \Pi^{t}(S)) - Q^{*}(S, \Pi^{*}(S))$$

$$- Q^{*}(S, \Pi^{t}(S)) + Q^{*}(S, \Pi^{t}(S))$$

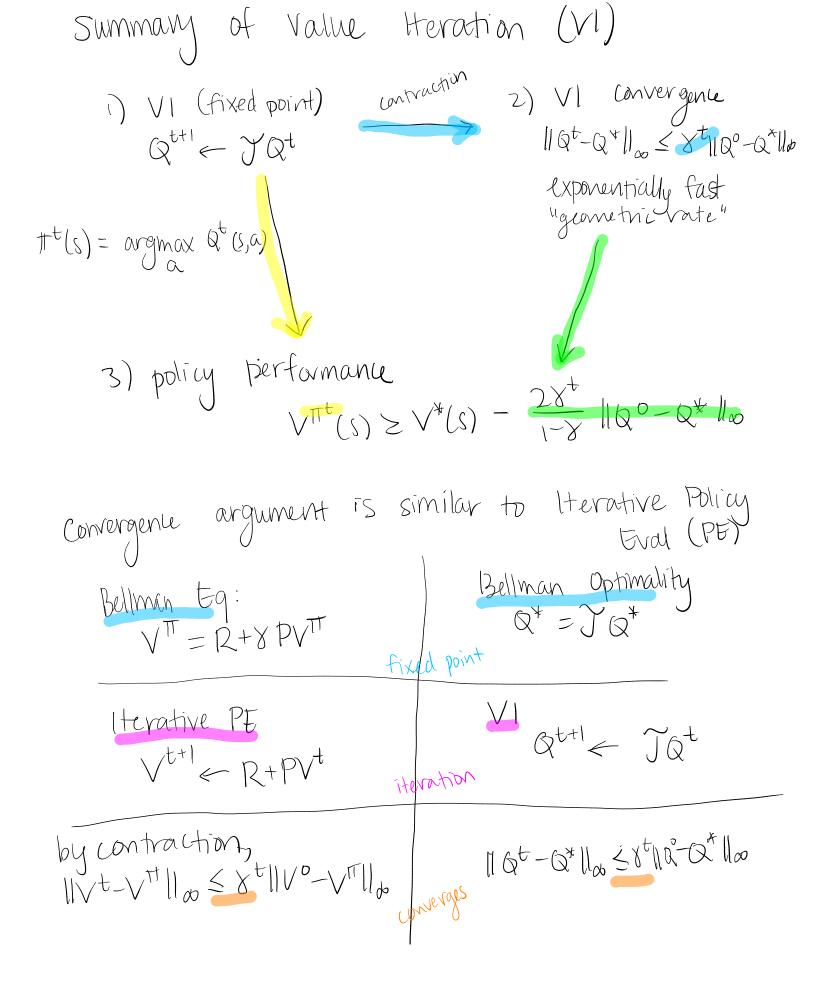
$$= X \mathbb{E} \left[V^{\text{tt}}(S) - V^{*}(S) \right] + Q^{*}(S, \Pi^{t}(S)) - Q^{*}(S, \Pi^{t}(S))$$

$$= S^{\text{re}}(S, \Pi^{t}(S))$$

$$= \begin{cases} \sum_{s \neq l, |\pi^{t}(s)|} - \sum_{s \neq l, |\pi^{t}(s)|} \\ = \sum_{s \neq l, |\pi^{t}(s)|} + \sum_{s \neq l, |\pi^{t}(s)|} \\ = \sum_{s \neq l, |\pi^{t}(s)|} + \sum_{s \neq l, |\pi^{t}(s)|$$

$$\geq 8 \mathbb{E} \left[\sqrt{\pi^{t}(s')}, \sqrt{4(s')} \right] - \frac{1}{100} \mathbb{E} \left[\sqrt{\pi^{t}(s')}, \sqrt{4(s')}, \sqrt{4(s')}, \sqrt{4(s')} \right] - \frac{1}{100} \mathbb{E} \left[\sqrt{\pi^{t}(s')}, \sqrt{4(s')}, \sqrt{$$

$$\geq X \mathbb{E}[V^{\pi^{t}}(S'), V^{x}(S')] - 2X^{t} ||Q^{0} - Q^{x}||_{\infty}$$
 (convergence S'~P(S, $\pi^{t}(S)$)



2) Yoliy Iteration Another iterative Algorithm for approximating the optimal policy to* White Value iteration updates a function at each timestep (and then at the very end we transform at into It), policy iteration updates both a policy and a a function at each timestep. Algoritmi Policy Heration Initialize TO: S-D(A) for t=0,1, ---Polity Evaluation Qtt (s, a) \ts,a) Policy Improvement: It to (s) = argmax Qt (s, a) 45 In each iteration, we first use policy evaluation to compute the Q function associated with the current policy. Then, we "argmax" that à function to generate a new policy, aka, policy improvement. Aside: How do we get QTT from policy evaluation?

VT = (I-8P)-'R

Then $Q^{\text{fit}}(S,a) = r(S,a) + \mathbb{E}[V^{\text{fit}}(S')] \forall S, a$

we will prove two key properties of policy iteration.

n) Monotonic Improvement

 $Q^{Tt+1}(S, \alpha) \geq Q^{Tt+1}(S, \alpha) \forall S, \alpha$

2) convergence

NY*-VTT N= 8 t NV*-VTO No

```
Lemma (Monotonic Improvement):
                                                                                                                              For policy iteration, QTt+1 (s,a) > QTt+1 (s
                  Proof:

\frac{\partial^{+}(s,\alpha)}{\partial s} = \frac{\partial s}{\partial s} + \frac{\partial s}{\partial s} = \frac{\partial s}
                                                                                                                                                                                                                                                                                                                                                           = \chi \mathbb{E}\left[Q^{\dagger tr'}(S') - Q^{\dagger t}(S', \dagger t'(S'))\right]
= \chi \mathbb{E}\left[Q^{\dagger tr'}(S', \dagger t'(S')) - Q^{\dagger t}(S', \dagger t'(S'))\right]
                expectation of ord
                                                                                                                                                                                                                                                                               = \times \mathbb{E} \left[ Q^{\mathsf{T}^{\mathsf{t}^{\mathsf{t}}}}(S', \mathsf{T}^{\mathsf{t}^{\mathsf{t}^{\mathsf{t}}}}(S')) - Q^{\mathsf{T}^{\mathsf{t}^{\mathsf{t}}}}(S', \mathsf{T}^{\mathsf{t}^{\mathsf{t}^{\mathsf{t}}}}(S') \right]
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            +Q^{Tt}(S', Tt^{tt}(S')) - Q^{Tt}(S', Tt^{t}(S')) > 0
 Tto is
defined as
                                                                                                                                                                                                                                                                       \geq \chi \left[ \left[ Q^{Ttt'}(S') - Q^{Tt}(S') \right] - Q^{Tt}(S') \right]
                                                                                                                                                                                                                                                                                                                                                    5'~P(5,a)
         (iterate)
                                                                                                                                                                                                                                                                \geq \chi^{2} \left[ Q^{\dagger t+1}(S'', T^{t+1}(S'')) - Q^{\dagger t}(S'', T^{t+1}(S'')) \right]
                                                                                                                                                                                                                                                                                                                                      S "~ P(s',TT+1(s'))

\[
\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\ti}}}}}} \ext{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\tinite\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\text{\
                                        Does this imediately imply that V^{Tt+1}(S) \geq V^{Tt}(S)? (see proof below)
```

```
Theorem (convergence):
        For policy iteration, \|V^{\dagger t} - V^*\|_{\infty} \leq \sqrt[4]{\|V^{\dagger \dagger} - V^*\|_{\infty}}
   Proot
      \frac{1}{V^*(s)} - V^{\pi t+1}(s) = \max_{\alpha} \left[ r(s,\alpha) + \chi \mathbb{E} V^*(s') \right] - \left[ (s,\pi t+1) + \chi \mathbb{E} V^*(s') \right] 
                     Aside: by Lemma, Q(s,a) \ge Q(s,a) + s, a
                          setting a= tttl(s),
                                        QHttl(s)) > QT(S,TT(s)) YS
                          Recall that mtti(s) is defined to maximize

Out (s, ·). Therefore

Ht
                                definition VTTt1 (S) ZQ(S, a) 4S, a
                                       choosing \alpha = T^{t}(S), V^{TT}(S) \geq V^{TT}(S) + S.
       V^{*}(s) - V^{(s)} \leq \max_{\alpha} \left[ v(s_{|\alpha}) + \gamma \notin V(s_{|\alpha}) \right] - \left[ v(s_{|\pi} t_{|s|}^{t+1}) + \gamma \notin V(s_{|\alpha}) \right]
  (of The)
                          = max [r(s,a)+8 EV(s')] - max [r(s,a) +8 EV(s)]
 (mox E(x) mox dx) ch)
                          < max [risiat + 8 FV*(s') - (risia) + 8 FV*(s')]
(Fig. 7 not Erry)
                          \leq \max_{a, s'} \chi(V^*(s') - V^{\mathsf{Tt}}(s')) = \chi ||V^* - V^{\mathsf{Tt}}||_{\infty}
            11Vt+1-V*1100 = 811Vt-V*1100
                          => 11Vt -V*1100 < 8t 11V0-V*1100
```

Both value iteration and policy iteration have geometric/exponential convergence Value It.

Noting It.

Noting It.

Noting It.

Noting It.

Noting It. while this is a very fast convergence rate, for any finite , its not equal to 0, r.e. it is not exact. In HWZ, you will see that in fact Policy iteration is guaranteed to exactly converge to the optimal policy in a finite number of steps. (the same 4's not true for value Iteration)

3) Finite Hovizon MDP

M= {S, A, P, r, H, Mos

states S, actions 86, transitions P, rewards r as before Horizon HENT (length of time) Initial state distribution $Mo \in \Delta(S)$ $S_o \sim M_o$

The task starts from an initial distribution and lasts for H steps (common in robotics) $\max_{t} \mathbb{E}\left[\sum_{t=0}^{H-1} r(a_{t}, s_{t})\right]$

Stop P(Stat), So~Mo,] $\alpha_t = \Pi_t(S_t)$

(deterministic

In general, we consider time-vary Policies $T = (TO_1, TI_1, -, TI_{H-1})$ The value and Q function are $V_{t}^{\#}(S) = \mathbb{E} \left[\sum_{k=t}^{H-1} r(S_{k}, a_{k}) \middle| S_{t} = S, \ a_{k} = TI_{k}(S_{k}), \ S_{k+1} \sim P(S_{k}, a_{k}) \right]$ $Q_{t}^{\#}(S_{t}a) = \mathbb{E} \left[\sum_{k=t}^{H-1} r(S_{k}, a_{k}) \middle| (S_{t}, a_{t}) = (S_{t}a), \ a_{k} = TI_{k}(S_{k}), \ S_{k+1} \sim P(S_{k}, a_{k}) \right]$ time-varying: $S_{k+1} \sim P(S_{k}, a_{k})$

Bellman Equation:

$$Q_{t}^{t}(s,a) = r(s,a) + \mathbb{E}\left[V_{t+1}^{TT}(s')\right]$$

$$S'\sim P(s,a)$$

$$V_{t}^{t}(s')$$

$$V_{t}^{t}(s')$$

$$V_{t}^{t}(s')$$

Because the horizon is finite, the recursion implied by the Bellman equation is finite and we can compute the optimal policy backwards through time.

4) Dynamic Programming
to find TT*= (TT*, TT+1)

Start with H-1 (note VH(S)=0 since H is past horizon)

$$Q_{H-1}^{*}(s, \alpha) = r(s, \alpha)$$
 $T_{H+1}^{*}(s) = \underset{\alpha}{\operatorname{argmax}} Q_{H-1}^{*}(s, \alpha)$

Bellman by
$$V_{H-1}^{*}(S) = Q_{H-1}^{*}(S) = Q_{H-1}^{*}(S)$$