

lower $\gamma \rightarrow$ faster convergence

Last Time

Proved VI converges
Bellman Operator, $B[V]$, is a γ contraction mapping in $(\mathbb{R}^{|S|}, \|\cdot\|_\infty)$
 V^* is unique

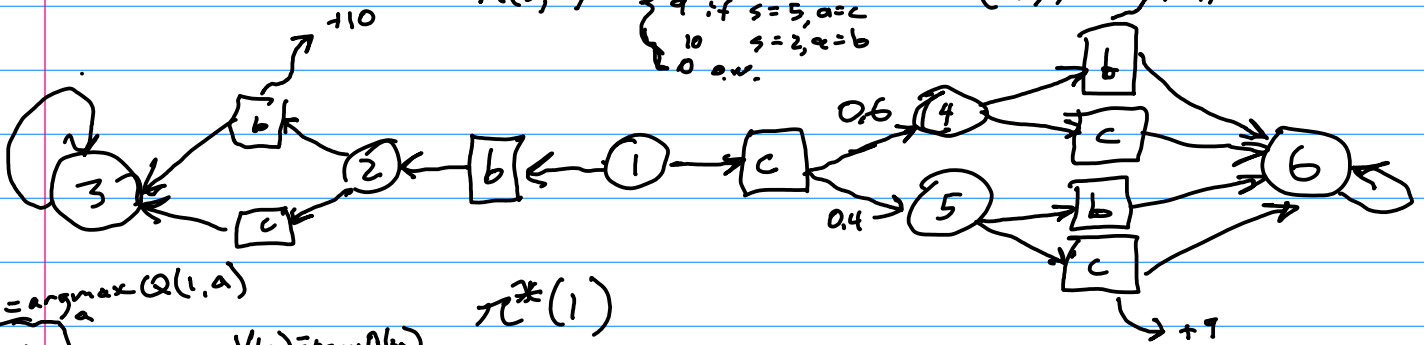
This Time

\rightarrow Debugging and Perf.
 \rightarrow Online Solution Methods

Breakout Rooms

$S = \{1, \dots, 6\}$ $A = \{b, c\}$ $\gamma = 1$

$$R(s, a) = \begin{cases} 11 & \text{if } s=4, a=b \\ 9 & \text{if } s=5, a=c \\ 10 & \text{if } s=3, a=b \\ 0 & \text{o.w.} \end{cases} \quad T(s'|s, a)$$



$$\pi^*(1) = \arg\max_a Q(1, a)$$

$$\pi^*(1)$$

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

$$Q$$

	V	b	c
1	10.2	10	10.2
2	10	10	0
3	0		
4	11	11	0
5	9	0	9
6	0		

$$Q(s, a) = R(s, a) + \sum_{s'} T(s'|s, a) V(s')$$

b	c
$0 + 1.0 \cdot 10$	$0 + 0.6 \cdot 11 + 0.4 \cdot 9$
$10 + 0$	$0 + 0$
$11 + 0$	$0 + 0$
$0 + 0$	$9 + 0$

Offline

P.I. V.I.

Before Execution: find V^*/Q^*

During : $\pi^*(s) = \arg\max_a Q^*(s,a)$

"Policy Extraction"

Online

Before Execution: nothing

During : Everything

Consider Actions and Consequences

$|S|$ too big

Only consider states that are reachable before γ gets too small

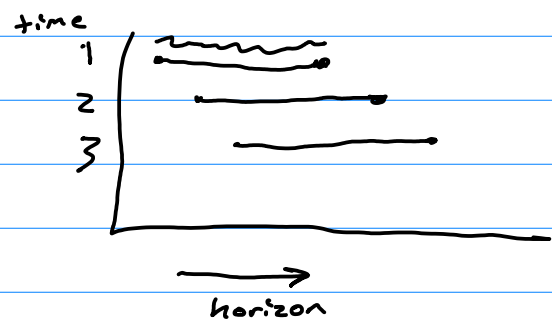
Receding Horizon

loop

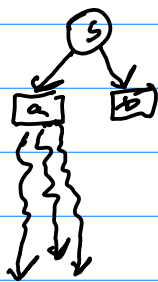
measure state s

Calculate $a = \pi(s)$

take a



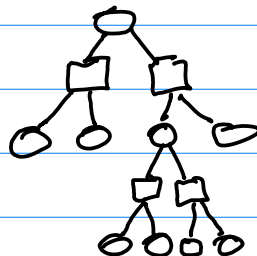
Lookahead with Rollouts



$$R(s,a) + \gamma \hat{V}(s')$$

$\hat{V}(s') = \text{rollout with } \pi_0$

Tree Search



$\leftarrow Q$

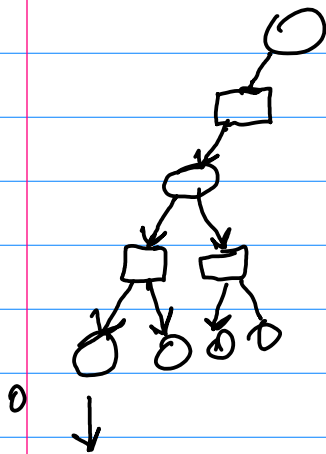
$$\leftarrow V(s) = \max_a Q(s,a)$$

$$\leftarrow Q(s,a) = R(s,a) + \gamma \mathbb{E}_{s' \sim T(s,a)} [V(s')]$$

$\leftarrow V$

$$\sum_{s'} \pi(s'|s,a) V(s')$$

Forward Search



$$O((|S| \times |A|)^d)$$

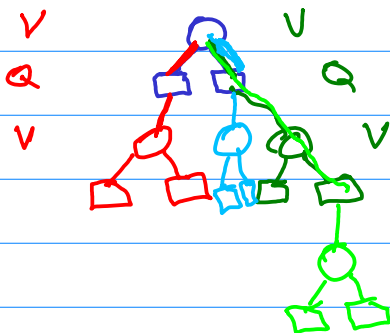
Algorithm 4.6 Forward search

```

1: function SELECTACTION( $s, d$ )
2:   if  $d = 0$ 
3:     return (NIL, 0)
4:   ( $a^*, v^*$ )  $\leftarrow$  (NIL,  $-\infty$ )
5:   for  $a \in A(s)$ 
6:      $v \leftarrow R(s, a)$ 
7:     for  $s' \in S(s, a)$ 
8:       ( $a', v'$ )  $\leftarrow$  SELECTACTION( $s', d - 1$ )
9:        $v \leftarrow v + \gamma T(s' | s, a) v'$ 
10:    if  $v > v^*$ 
11:      ( $a^*, v^*$ )  $\leftarrow$  ( $a, v$ )
12:   return ( $a^*, v^*$ )
  
```

Monte Carlo Tree Search (MCTS) (UCT)

Search Expand Rollout Backup



```

function simulate!( $\pi$ :MonteCarloTreeSearch,  $s, d=\pi.d$ )
  if  $d \leq 0$ 
    return 0.0
  end
   $\mathcal{P}, N, Q, c = \pi.\mathcal{P}, \pi.N, \pi.Q, \pi.c$ 
   $\mathcal{A}, TR, \gamma = \mathcal{P}.\mathcal{A}, \mathcal{P}.TR, \mathcal{P}.\gamma$ 
  if !haskey( $N, (s, first(\mathcal{A}))$ )
    for  $a$  in  $\mathcal{A}$ 
       $N[(s, a)] = 0$ 
       $Q[(s, a)] = 0.0$ 
    end
    return rollout( $\mathcal{P}, s, \pi, \pi, d$ )
  end
   $a = \text{explore}(\pi, s)$ 
   $s', r = TR(s, a)$ 
   $q = r + \gamma \text{simulate}!(\pi, s', d-1)$ 
   $N[(s, a)] += 1$ 
   $Q[(s, a)] += (q - Q[(s, a)]) / N[(s, a)]$ 
  return  $q$ 
end
  
```

$$Q(s, a) + \underbrace{c \sqrt{\frac{\log N(s)}{N(s, a)}}}_{\text{Exploration Bonus}}$$

low $N(s, a) / N(s)$
 ↳ long Exp Bonus

Exploration Bonus