

AML, 22 Oct 2019

## Markov Decision Processes

State dependent rewards

$$p(s', r | s, a)$$

Assume states,  
rewards are finite

At  $s$ , action  $a$   
leads to  $s'$  with  
reward  $r$

Policy - "strategy to choose action"

Given a policy  $\pi$ ,  $v_{\pi}(s)$  - expected long term value of state  $s$

Over a path, we discount rewards

$$G_t = R_{t+1} + \gamma G_{t+1}$$

$$v_{\pi}(s) = E_{\pi} [G_t | S_t = s]$$

$$= \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma v_{\pi}(s')]$$

$$q_{\pi}(s, a) = \dots$$

## Optimal policy

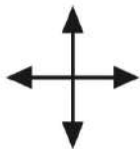
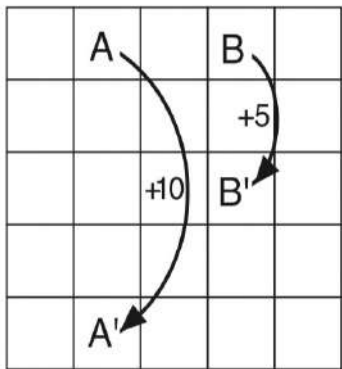
$$\pi_1 \geq \pi_2 \quad \text{if } \forall s. \quad v_{\pi_1}(s) \geq v_{\pi_2}(s)$$

## Optimal value function

$$v_*(s) = \max_a E [R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a]$$

$$= \max_a \sum_{s', r} p(s', r \mid s, a) [r + \gamma v_*(s')]$$

$$q_*(s, a) = \dots$$



Actions

N, S, E, W

$\sqrt{V}\pi$

3.3	8.8	4.4	5.3	1.5
1.5	3.0	2.3	1.9	0.5
0.1	0.7	0.7	0.4	-0.4
-1.0	-0.4	-0.4	-0.6	-1.2
-1.9	-1.3	-1.2	-1.4	-2.0

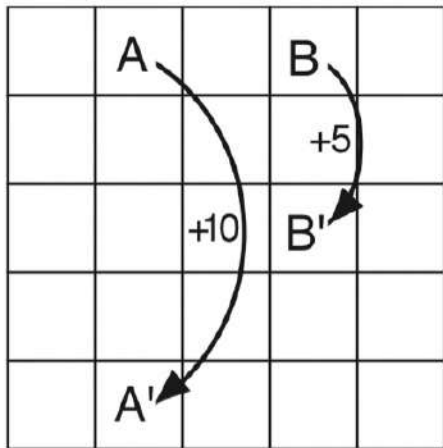
Uniform  $\uparrow$   
policy  $\pi$

At A, B all actions move as shown

All other rewards are 0

except, -1 for hitting edge

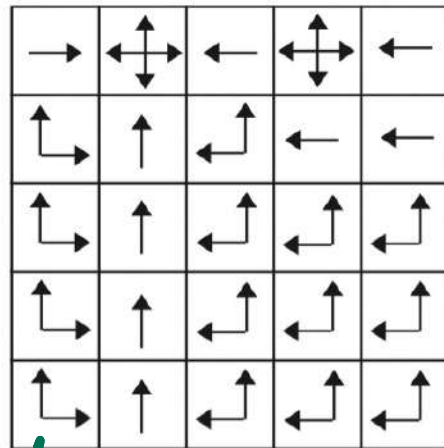
$\pi \rightarrow \sqrt{V}\pi$



Gridworld

22.0	24.4	22.0	19.4	17.5
19.8	22.0	19.8	17.8	16.0
17.8	19.8	17.8	16.0	14.4
16.0	17.8	16.0	14.4	13.0
14.4	16.0	14.4	13.0	11.7

$v_*$



$\pi_*$

Same Example

Choose  $\uparrow \rightarrow$   
with any probability

0 prob for all nonoptimal dir.

Computing  $V_\pi$  from  $\pi$

$$V_\pi(s) = \sum_a \pi(a|s) \sum_{s', r} P(s', r | s, a) [r + \gamma V_\pi(s')]$$

Bellman Eqn

Provided  $\gamma < 1$ , this set of eqns has a unique solution

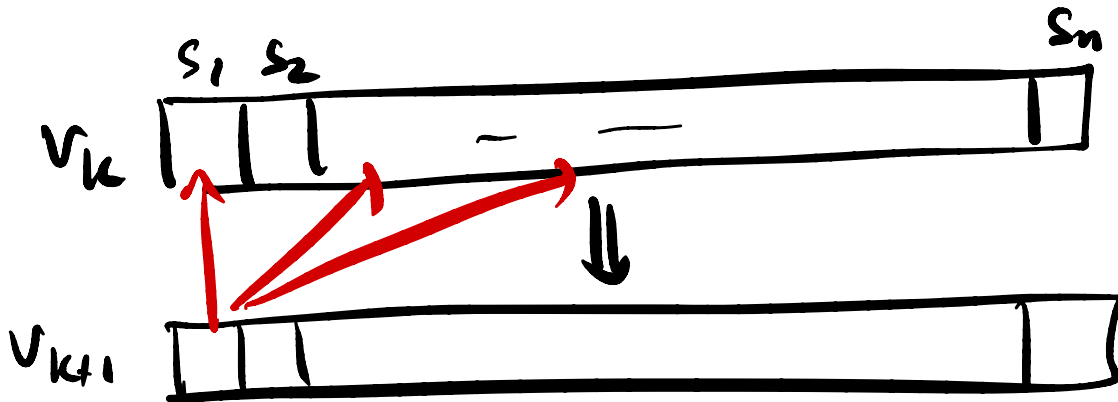
If we have episodic system with finite paths, assume a terminal state with value 0  
S original states,  $S^+$  with added terminal state

Iterative soln

Treat equation as update rule

$$V_{k+1}(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma V_k(s')]$$

Update all  $V_{k+1}(s)$  in parallel



Repeat until convergence ( $\Delta \leq \epsilon$ )

Each update:  $\Delta = \max_s (V_{k+1}(s) - V_k(s))$

$\pi \rightarrow V_\pi$  : Policy evaluation

Convergence guaranteed

In practice - need not update in parallel

Sequentially update  $V_{k+1}(s_1), V_{k+1}(s_2) \dots$



Goal is to reach  $V_*$ ,  $\pi_*$

Policy update

$$\pi \rightarrow V_\pi$$

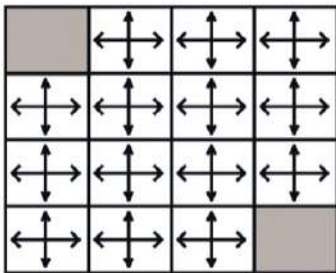
At  $s$ , choose action  $\pi'(s)$  instead of  $\pi(s)$  s.t.

$$q(s, \pi'(s)) \geq V_\pi(s)$$

$$\Rightarrow \forall s \quad V_{\pi'}(s) \geq V_\pi(s)$$

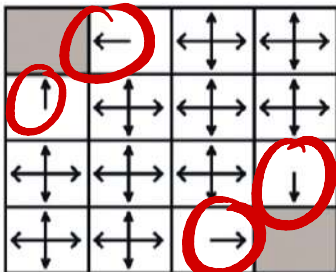
Greedy policy update "works"

0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0



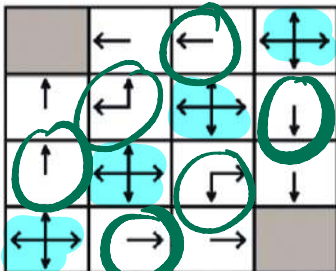
← random policy

0.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	0.0



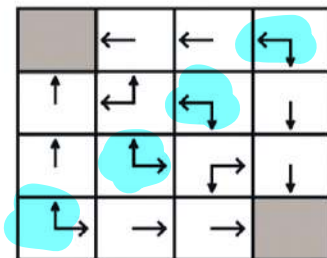
*greedy update.*

0.0	-1.7	-2.0	-2.0
-1.7	-2.0	-2.0	-2.0
-2.0	-2.0	-2.0	-1.7
-2.0	-2.0	-1.7	0.0



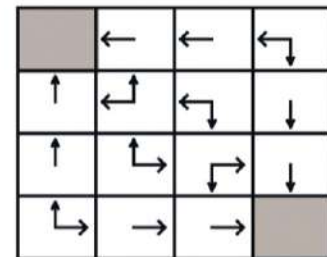
$k = 3$

0.0	-2.4	-2.9	-3.0
-2.4	-2.9	-3.0	-2.9
-2.9	-3.0	-2.9	-2.4
-3.0	-2.9	-2.4	0.0



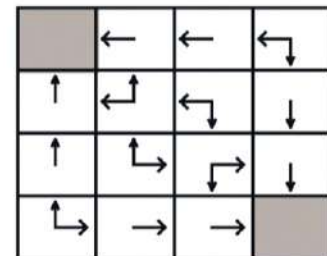
$k = 10$

0.0	-6.1	-8.4	-9.0
-6.1	-7.7	-8.4	-8.4
-8.4	-8.4	-7.7	-6.1
-9.0	-8.4	-6.1	0.0



$k = \infty$

0.0	-14.	-20.	-22.
-14.	-18.	-20.	-20.
-20.	-20.	-18.	-14.
-22.	-20.	-14.	0.0



optimal policy

$$\begin{aligned}
v_\pi(s) &\leq q_\pi(s, \pi'(s)) \\
&= \mathbb{E}[R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s, A_t = \pi'(s)] \\
&= \mathbb{E}_{\pi'}[R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s] \\
&\leq \mathbb{E}_{\pi'}[R_{t+1} + \gamma q_\pi(S_{t+1}, \pi'(S_{t+1})) \mid S_t = s] \\
&= \mathbb{E}_{\pi'}[R_{t+1} + \gamma \mathbb{E}[R_{t+2} + \gamma v_\pi(S_{t+2}) \mid S_{t+1}, A_{t+1} = \pi'(S_{t+1})] \mid S_t = s] \\
&= \mathbb{E}_{\pi'}[R_{t+1} + \gamma R_{t+2} + \gamma^2 v_\pi(S_{t+2}) \mid S_t = s] \\
&\leq \mathbb{E}_{\pi'}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 v_\pi(S_{t+3}) \mid S_t = s] \\
&\vdots \\
&\leq \mathbb{E}_{\pi'}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots \mid S_t = s] \\
&= v_{\pi'}(s).
\end{aligned}$$

Proof that greedy policy update  
is sound

# Policy iteration

Initial policy

$\pi_0$

policy  
evaluation

$V_{\pi_0}$

greedy  
update

$\pi_1$

$V_{\pi_1}$

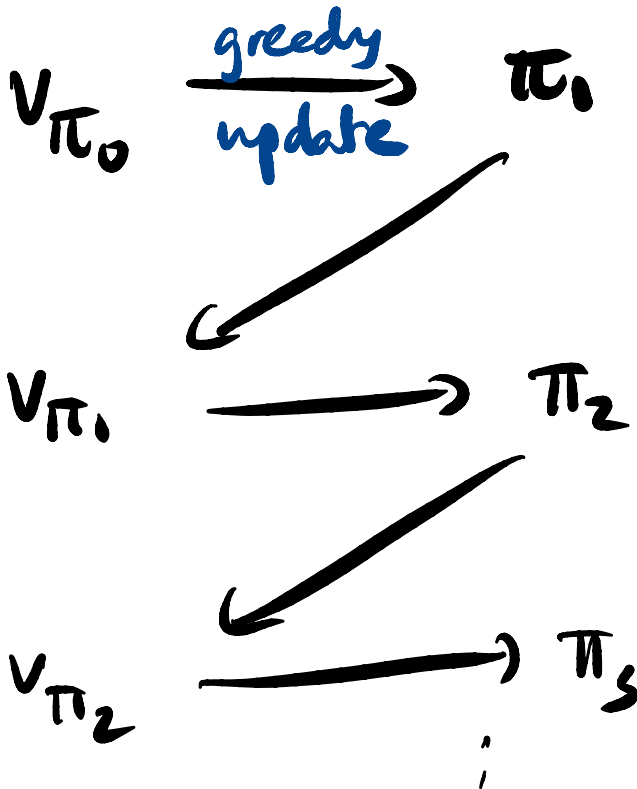
$\pi_2$

$V_{\pi_2}$

$\pi_3$

iterative  
evaluation

Expensive, but theoretically  
feasible



Treating Bellman equations as update rules

## Dynamic Programming

Policy evaluation - repeatedly sweep across all states

## Simplifying Policy Iteration

Don't need to fully evaluate  $v_{\pi}$  for each  $\pi$  in the iteration

Suffices to compute

$$V_{\pi}^0 \rightarrow V_{\pi}^1 \quad (1 \text{ sweep})$$

Collapse  $V_{\pi} \rightarrow \pi' \rightarrow V_{\pi'} [1 \text{ step}]$

$$V_{k+1}(s) = \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V_k(s)]$$

Value Iteration

In practice

Asynchronous D.P.

Need not update  $V(s)$   $\forall s$  in each pass

"Fairness"

In an infinite sequence of updates,  
every state is updated infinitely often