# Policy Iteration

# Objectives

- ☐ Understand the policy improvement theorem

- ☐ Use a value function for a policy to produce a better policy for a given MDP

- ☐ Apply policy iteration to compute optimal policies and optimal value functions

# Policy Improvement

□ We can find the optimal policy by choosing the Greedy action.

□ The Greedy action maximizes the Bellman's optimality equat

**Recall that**    Greedy action

$$\pi_*(s) = \operatorname*{argmax}_a \sum_{s'} \sum_r p(s',r\,|\,s,a)\big[r + \gamma v_*(s')\big]$$

**?** $\quad \operatorname*{argmax}_a \sum_{s'} \sum_r p(s',r\,|\,s,a)\big[r + \gamma v_\pi(s')\big]$

$$\pi(s) = \operatorname*{argmax}_a \sum_{s'} \sum_r p(s',r\,|\,s,a)\big[r + \gamma v_\pi(s')\big] \text{ for all } s \in \mathcal{S}$$

➡ $v_\pi$ obeys the Bellman optimality equation ➡ $\pi$ is optimal

# Policy Improvement

□ The new policy obtained in this way must be a strict improvement on Pi, unless Pi was already optimal.

□ Policy improvement theorem: q Pi: value of a state if you take action A, and then follow policy Pi.

$$q_\pi\big(s, \pi'(s)\big) \geq q_\pi\big(s, \pi(s)\big) \text{ for all } s \in \mathcal{S}$$
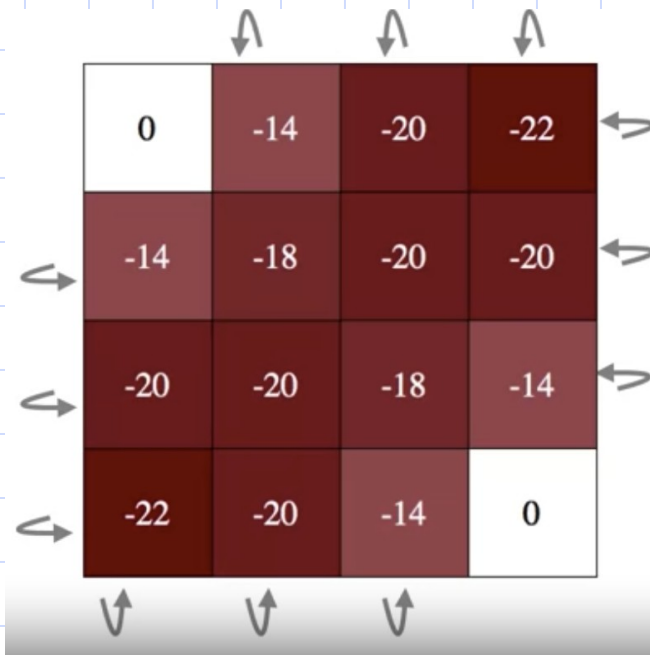
# Policy Improvement

☐ Policy Pi prime is at least as good as Pi if in each state, the value of the action selected by Pi prime is greater than or equal to the value of the action selected by Pi.

☐ Policy pi prime is strictly better if the value is strictly greater and at least one state.

$$q_\pi(s, \pi'(s)) \geq q_\pi(s, \pi(s)) \text{ for all } s \in \mathcal{S} \rightarrow \pi' \geq \pi$$

$$q_\pi(s, \pi'(s)) > q_\pi(s, \pi(s)) \text{ for at least one } s \in \mathcal{S} \rightarrow \pi' > \pi$$

# Policy Improvement

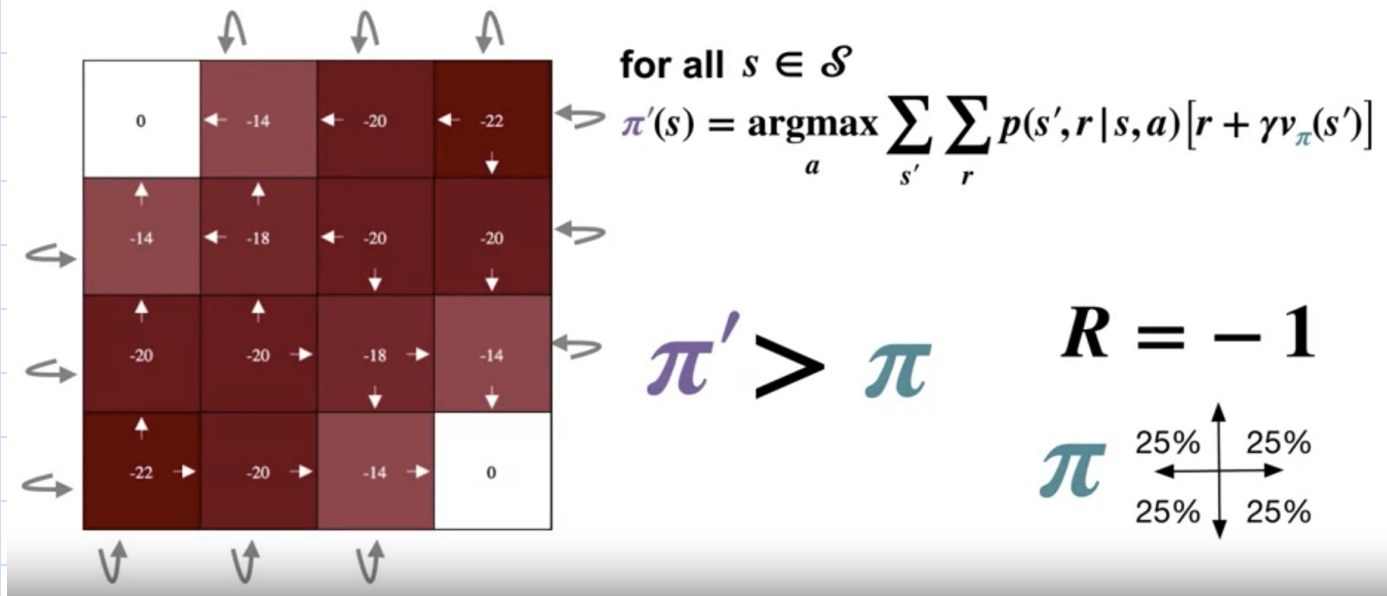□ For example the four-by-four grid. The final value function:



$$R = -1$$

# Policy Improvement

□ The greedy Pi policy



for all $s \in \mathcal{S}$

$$\pi'(s) = \underset{a}{\text{argmax}} \sum_{s'} \sum_{r} p(s',r \mid s,a)\big[r + \gamma v_{\pi}(s')\big]$$

$$\pi' > \pi$$

$$R = -1$$

# Policy Iteration

☐ Policy improvement theorem

Greedy action

$$\pi'(s) \doteq \underset{a}{\mathrm{argmax}} \sum_{s'} \sum_{r} p(s',r\,|\,s,a)\big[r + \gamma v_{\pi}(s')\big]$$

$\pi'$ is strictly better than $\pi$ unless $\pi$ was already optimal

# Policy Iteration

☐ Begin with the policy Pi 1→ evaluate Pi 1 using iterative policy evaluation to obtain the state value, V Pi 1. --> evaluation step.

☐ Using the results of the policy improvement theorem, we can then greedify with respect to v Pi 1 to obtain a better policy, Pi 2--> We call this the improvement step.

☐ We can then compute V Pi 2 and use it to obtain an even better pol[...] f better policies.

Evaluation

Improvement

$$\pi_1 \xrightarrow{E} v_{\pi_1} \xrightarrow{I} \pi_2 \xrightarrow{E} v_{\pi_2} \xrightarrow{I} \pi_3 \xrightarrow{E} \ldots \xrightarrow{I} \pi_*$$

# Policy Iteration

☐ We complete an iteration, and the policy remains unchanged→ found the optimal policy.

☐ Each policy generated in this way is deterministic.

☐ This method of finding an optimal policy is called policy iteration
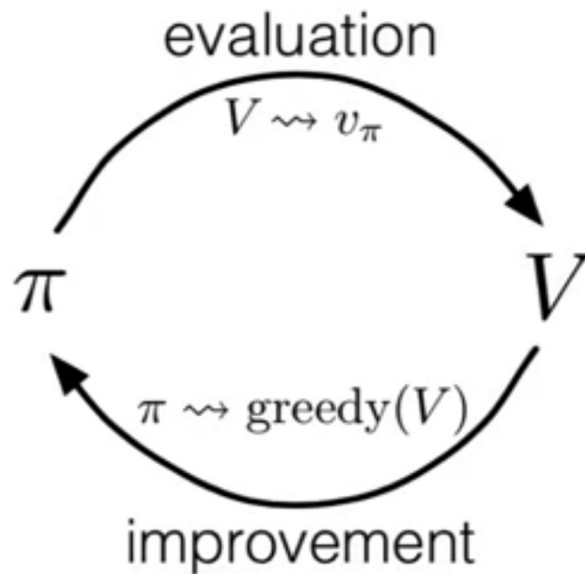
**Evaluation**

**Improvement**

$$\pi_1 \xrightarrow{E} v_{\pi_1} \xrightarrow{I} \pi_2 \xrightarrow{E} v_{\pi_2} \xrightarrow{I} \pi_3 \xrightarrow{} \ldots \xrightarrow{} \pi_* \xrightarrow{E} v_{\pi_*} \xrightarrow{I} \pi_*$$

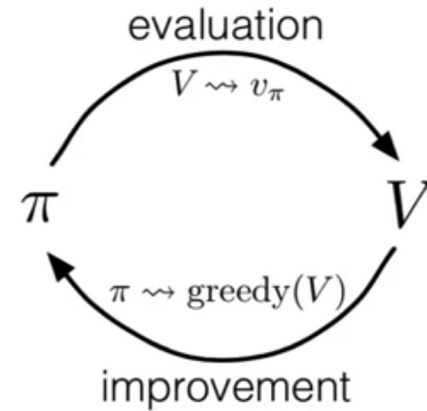**Deterministic**

# Policy Iteration

□ Policy iteration consists of two distinct steps: evaluation and improvement.

# Policy Iteration

☐ We first evaluate our current policy, Pi 1, which gives us a new value function that accurately reflects the value of Pi 1.
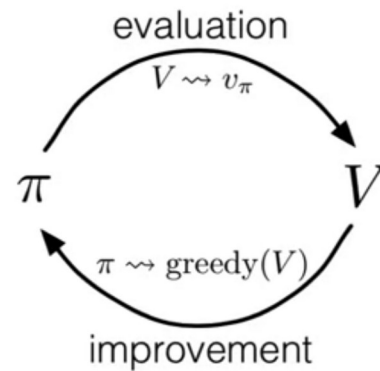
$$\pi_1 \longrightarrow v_{\pi_1}$$



evaluation
$V \rightsquigarrow v_\pi$

$\pi$        $V$

$\pi \rightsquigarrow \text{greedy}(V)$

improvement

# Policy Iteration

☐ The improvement step then uses V Pi 1 to produce a greedy policy Pi 2. At this point, Pi 2 is greedy with respect to the value function of Pi 1, but V Pi 1 no longer accurately reflects the value of Pi 2.
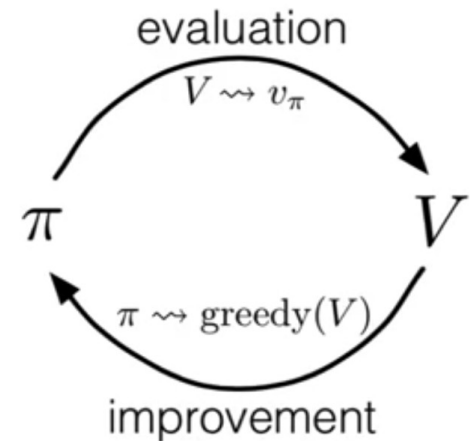
$$\pi_2 \longleftarrow v_{\pi_1}$$



evaluation
$V \rightsquigarrow v_\pi$

$\pi$        $V$

$\pi \rightsquigarrow \text{greedy}(V)$
improvement

13

# Policy Iteration

☐ The policy is greedy and the value function is accurate.

$$\boldsymbol{\pi}_* \longleftrightarrow \boldsymbol{v}_*$$

evaluation

$$V \rightsquigarrow v_\pi$$

$\pi$        $V$

$$\pi \rightsquigarrow \text{greedy}(V)$$

improvement

# Policy Iteration

□ Visualization of policy iteration

# Policy Iteration

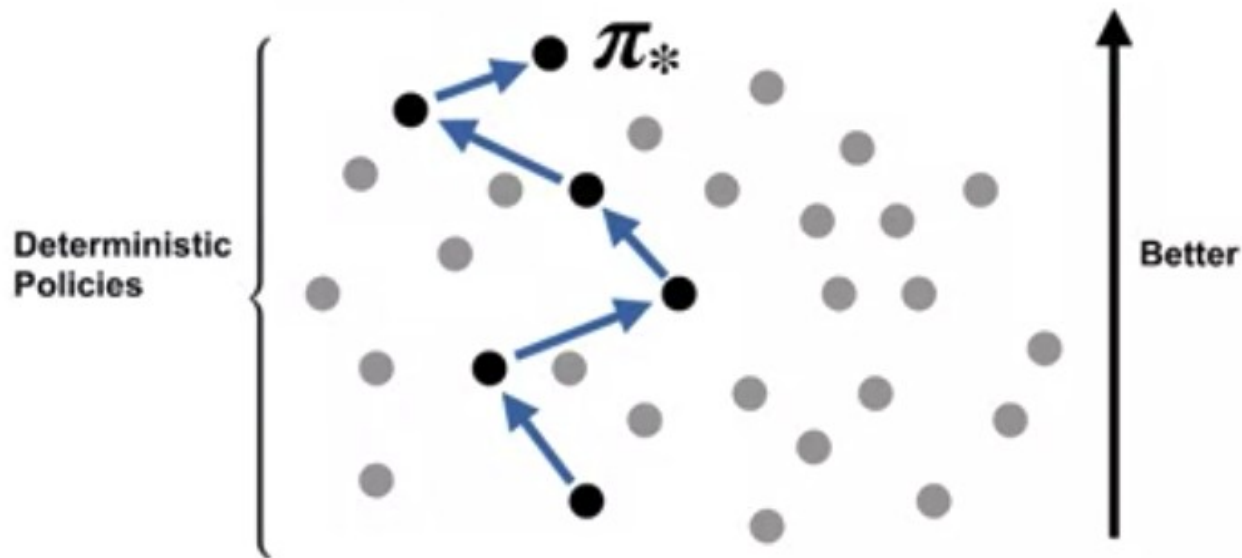☐ Policy Iteration using iterative policy evaluation for estimate pi ~ pi*

1. Initialization
   $V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in \mathcal{S}$

2. Policy Evaluation
   Loop:
       $\Delta \leftarrow 0$
       Loop for each $s \in \mathcal{S}$:
           $v \leftarrow V(s)$
           $V(s) \leftarrow \sum_{s',r} p(s',r|s,\pi(s))\big[r + \gamma V(s')\big]$
           $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
   until $\Delta < \theta$ (a small positive number determining the accuracy of estimation)

3. Policy Improvement
   $policy\text{-}stable \leftarrow true$
   For each $s \in \mathcal{S}$:
       $old\text{-}action \leftarrow \pi(s)$
       $\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s',r} p(s',r|s,a)\big[r + \gamma V(s')\big]$
       If $old\text{-}action \neq \pi(s)$, then $policy\text{-}stable \leftarrow false$
   If $policy\text{-}stable$, then stop and return $V \approx v_*$ and $\pi \approx \pi_*$; else go to 2

# Policy Iteration

□ The power of policy iteration

# Summary

- ☐ Understand the policy improvement theorem

- ☐ Use a value function for a policy to produce a better policy for a given MDP

- ☐ Apply policy iteration to compute optimal policies and optimal value functions

# Q & A