

# **535514: Reinforcement Learning**

## **Lecture 3 – MDP and Optimal Policies**

Ping-Chun Hsieh

February 26, 2024

# Announcement

- No class this Thursday (2/29)
- To make up for the lecture:
  - **Plan A:** 10min extension on 3/4 (Mon.), 3/11 (Mon.), 3/18 (Mon.), 3/25 (Mon.), 4/1 (Mon.)
  - **Plan B:** Make-up lecture on 3/7 (Thursday), 4:30pm-5:20pm

# Human Intelligence vs Machine Intelligence?



Robot PR1

Is this robot impressive to you?

(A video made by Pieter Abbeel back in 2013)

(Source: <https://www.youtube.com/watch?v=qBZPSTR96N4>)

# Human Intelligence vs Machine Intelligence?



Pieter Abbeel  
14.3K subscribers

Subscribe

Like 36



Share

Download



6,299 views May 29, 2013

This video showcases PR1, a robot developed by Keenan Wyrobek, Eric Berger, HFM Van der Loos, and Ken Salisbury at Stanford. It is the predecessor of the Willow Garage PR2. [The PR1 was tele-operated during this entire video.](#)

Reference: K. Wyrobek, E. Berger, H.F.M. Van der Loos, and K. Salisbury. Towards a personal robotics development platform: Rationale and design of an intrinsically safe personal robot. In the Proceedings of the IEEE International Conference on Robotics and Automation (ICRA), 2008.

Show less

Actually not that impressive  
(as PR1 was tele-operated by a human lol)

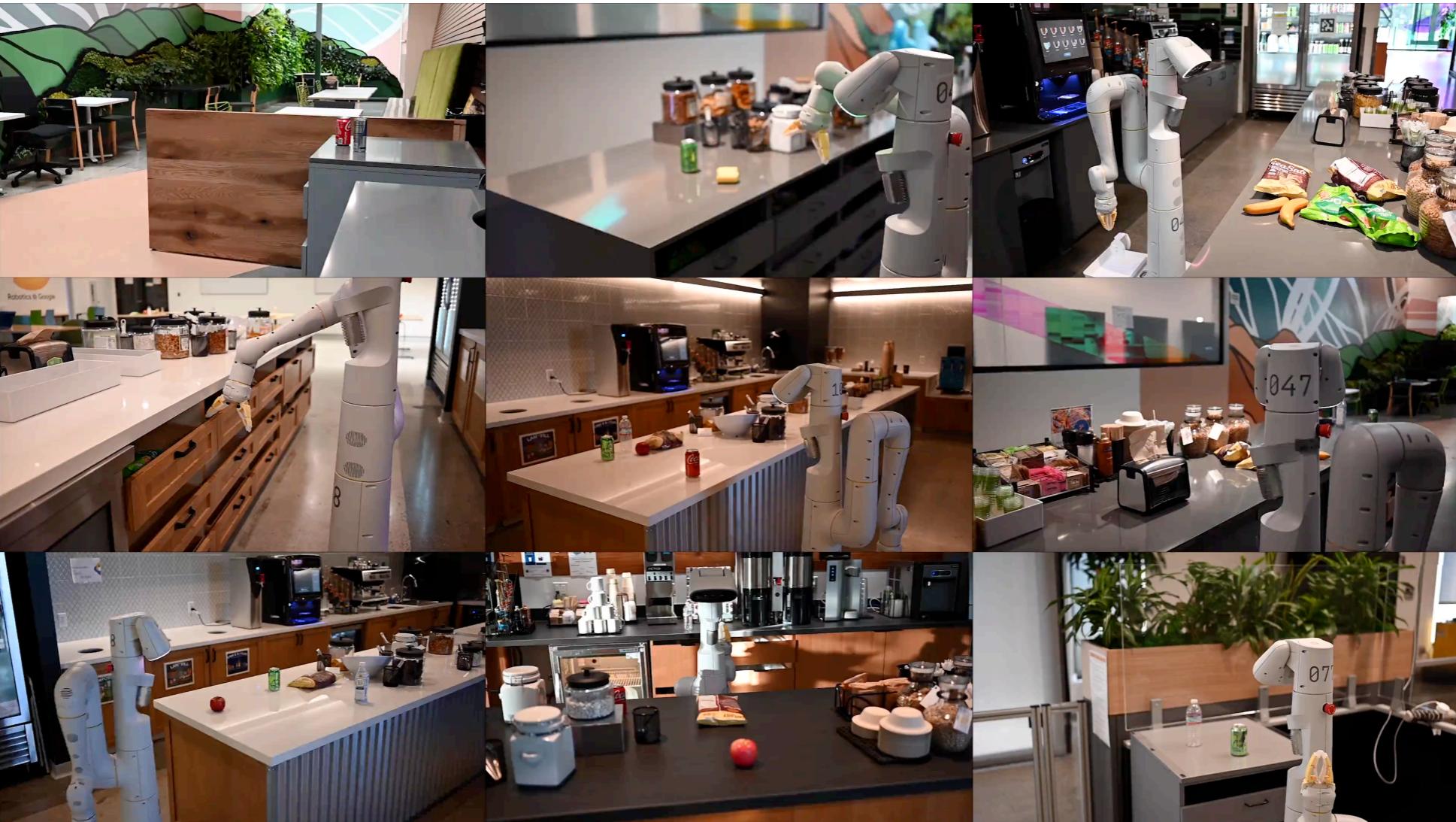
## Takeaways:

1. Humans and machines are NOT very different in “physical capability”
2. The difference lies in the “mind”, rather than the “body”

That is why we love RL research! (To extend human mind!)

# How About Now?

## RT-1 (Robotics Transformer) by Google



(Source: <https://blog.research.google/2022/12/rt-1-robotics-transformer-for-real.html>)

# Review: Markov Decision Process

- ▶ **Markov Decision Process (MDP)**: An MDP  $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$  is specified by

## Underlying Dynamics

1. State space  $\mathcal{S}$  (assumed finite)
2. Action space  $\mathcal{A}$  (assumed finite)
3. Transition matrix  $P = [P_{ss'}^a]$  with  $P_{ss'}^a = \mathbb{P}[s_{t+1} = s' | s_t = s, a_t = a]$

## Task / Goal

4. Reward function  $R_{s,a} = \mathbb{E}[r_{t+1} | s_t = s, a_t = a]$
5. Discount factor  $\gamma \in [0,1]$

- ▶ In this lecture, we shall assume the model parameters  $P$  and  $R_s^a$  are **known** (i.e. no learning)

# Review: How to Specify a Policy?

History-dependent policies  
 $\Pi_t(a_t=a \mid s_1, s_2, \dots, s_t)$

- ▶ **Idea:** “policy” is a lookup table specifying the action taken at any given state
- ▶ **(Randomized) Policy:** A policy  $\pi$  is a conditional distribution over possible actions given state  $s$ , i.e for any  $s \in \mathcal{S}, a \in \mathcal{A}$

$$\underline{\pi(a \mid s)} := \mathbb{P}(a_t = a \mid s_t = s)$$

- ▶ Let's use  $\Pi$  to denote the set of all stationary policies
- ▶ **Remark:** Here we focus on stationary policies, i.e.,  $\pi$  does not depend on time  $t$  [Puterman, 1994]
- ▶ **Question:** What's the intuition behind using stationary policies?

Markov property

# Connection Between MDP and MRP

- Idea: Fix a policy  $\pi(a | s)$  for an MDP  $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$ :

1. What is the probability of  $s \rightarrow s'$  under  $\pi$ ?

$$P_{ss'}^\pi = \sum_{a \in \mathcal{A}} \pi(a | s) P_{ss'}^a \quad (\text{total probability theorem})$$


2. What is the expected reward of begin in  $s$  under  $\pi$ ?

$$R_s^\pi = \sum_{a \in \mathcal{A}} \pi(a | s) R_{s,a}$$

- Under a fixed  $\pi(a | s)$ , we get an ( $\pi$ -induced) MRP  $(\mathcal{S}, P^\pi, R^\pi, \gamma)$

# Goals, Return, and State-Value Function of MDPs

- ▶ **Goal:** Given  $P$  and  $R$ , find a policy  $\pi$  that maximizes the expected cumulative reward (Question: this formulation can be viewed as optimal control, model-based RL, or model-free RL?)
- ▶ **Return  $G_t$ :** Cumulative discounted rewards over a single trajectory from  $t$  onwards (random)

$$G_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k r_{t+k}$$

- ▶ **State-value function  $V^\pi(s)$ :** Expected return if we start from state  $s$

$$V^\pi(s) = \mathbb{E}[G_t | s_t = s; \pi]$$

- ▶ **Question:** The expectation above is taken w.r.t. randomness of ?

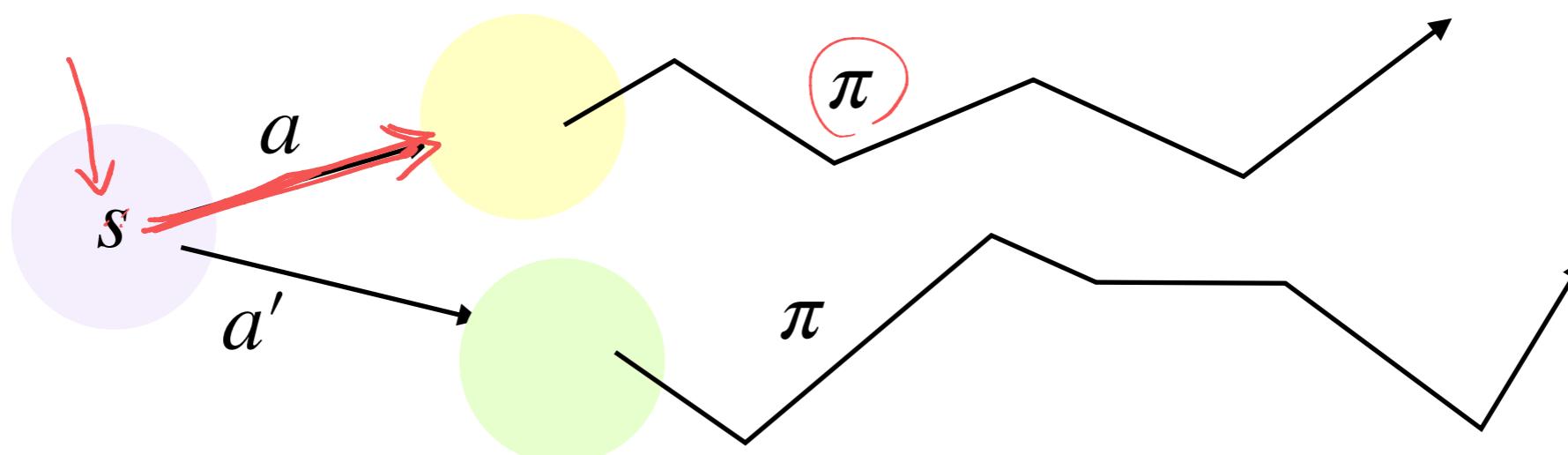
①  $\pi$     ② State transition    ③ Reward

# Action-Value Function $Q^\pi(s, a)$

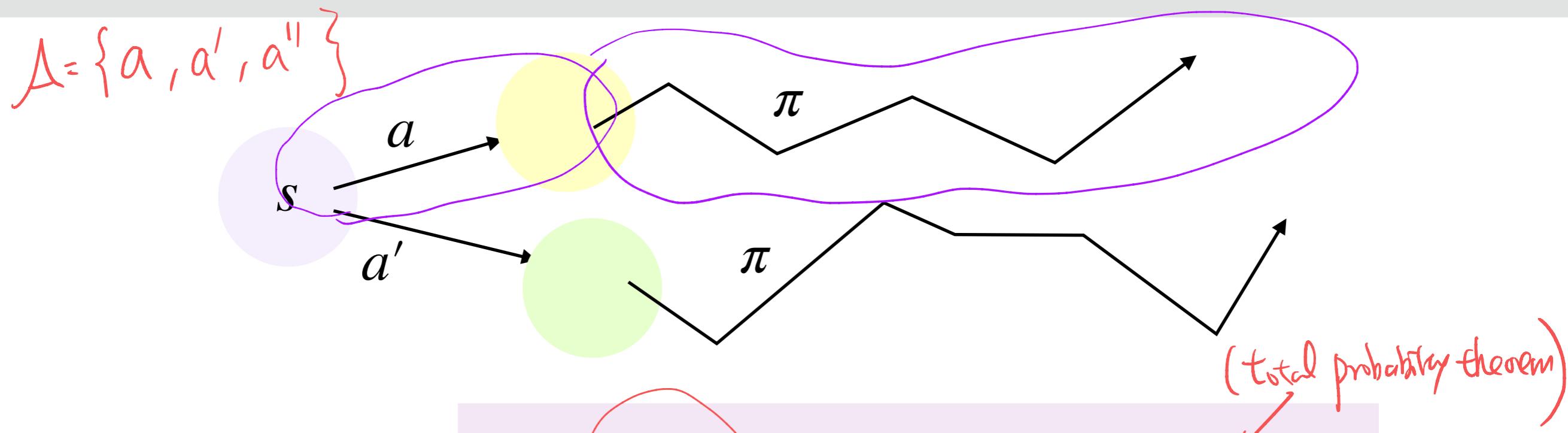
(Q-function)

- ▶ Action-value function  $Q^\pi(s, a)$ : Expected return if we start from state  $s$  and take action  $a$ , and then follow policy  $\pi$

$$Q^\pi(s, a) = \mathbb{E}[G_t | s_t = s, a_t = a; \pi]$$



# Natural Connection Between $V^\pi(s)$ and $Q^\pi(s, a)$



(1)  $V$  written in  $Q$

$$V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a | s) Q^\pi(s, a)$$

(total probability theorem)

(2)  $Q$  written in  $V$

$$Q^\pi(s, a) = \underbrace{R_{s,a}} + \gamma \sum_{s' \in \mathcal{S}} \underbrace{P_{ss'}^a} \underbrace{V^\pi(s')}$$

# Example: Value Functions

## Deep Sea Treasure

Initial state  $s_0 = (0,0)$



- ▶ **State:** Coordinate  $(x, y)$
- ▶ **Actions:** Right or Down
- ▶ **Reward:** Values on the treasures
- ▶ Submarine stays at the same place after hitting a wall
- ▶ The episode ends after finding a treasure
- ▶ **Policy:**  $\pi(\text{Right}|s) = \pi(\text{Down}|s) = 0.5$ , for all states  $s$
- ▶ **Question:**  $Q^\pi(s_0, \text{Down})=?$   $Q^\pi(s_0, \text{Right})=?$   $V^\pi(s_0)=?$

$$Q^\pi(s_0, \text{Down}) = 0.1$$

$$Q^\pi(s_0, \text{Right}) = \sum_i \left( \text{treasure}_i \right) \cdot (0.5)^{\# \text{ steps to reach treasure } i} \cdot (\# \text{ of possible paths to reach treasure } i)$$

$$(0.5)^2 \cdot (0.5)^4 \cdot C_1^2 C_1^1$$

$$V^\pi(s_0) = \underbrace{\pi(\downarrow|s_0)}_{0.5} Q^\pi(s_0, \text{Down}) + \underbrace{\pi(\rightarrow|s_0)}_{0.5} Q^\pi(s_0, \text{Right})$$

# Recursions for Computing $V^\pi(s)$ and $Q^\pi(s, a)$

(1) V written in Q

$$V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a | s) Q^\pi(s, a)$$

(2) Q written in V

$$Q^\pi(s, a) = R_{s,a} + \gamma \sum_{s' \in \mathcal{S}} P_{ss'}^a V^\pi(s')$$

(3) V written in V

$$V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a | s) \left( R_{s,a} + \gamma \cdot \sum_{s' \in \mathcal{S}} P_{ss'}^a V^\pi(s') \right)$$

(4) Q written in Q

$$Q^\pi(s, a) = R_{s,a} + \gamma \cdot \sum_{s' \in \mathcal{S}} P_{ss'}^a \cdot \left( \sum_{a' \in \mathcal{A}} \pi(a' | s) \cdot Q^\pi(s', a') \right)$$

Bellman equations

# (Non-Iterative) MDP Policy Evaluation

For  $V^\pi(s)$ :

$$V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a | s) \left( R_{s,a} + \gamma \sum_{s' \in \mathcal{S}} P_{ss'}^a V^\pi(s') \right)$$

Consider  $\pi$ -induced  
MRP  $(\mathcal{S}, P^\pi, R^\pi, \gamma)$ :

$$\begin{bmatrix} V^\pi(s^{(1)}) \\ \vdots \\ V^\pi(s^{(100)}) \end{bmatrix}$$

Matrix form:

$$R^\pi = \sum_{a \in \mathcal{A}} \pi(a | s) R_{s,a}$$

$$P_{ss'}^\pi = \sum_{a \in \mathcal{A}} \pi(a | s) P_{ss'}^a$$

$$V^\pi = R^\pi + \gamma P^\pi V^\pi$$

From Lecture 2

Solution of  $V^\pi$ :

$$V^\pi = (I - \gamma P^\pi)^{-1} R^\pi$$

# Iterative MDP Policy Evaluation (IPE)

We know:

$$V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a | s) \left( R_{s,a} + \gamma \sum_{s' \in \mathcal{S}} P_{ss'}^a V^\pi(s') \right)$$

- ▶ Iterative policy evaluation for a fixed policy  $\pi$ :

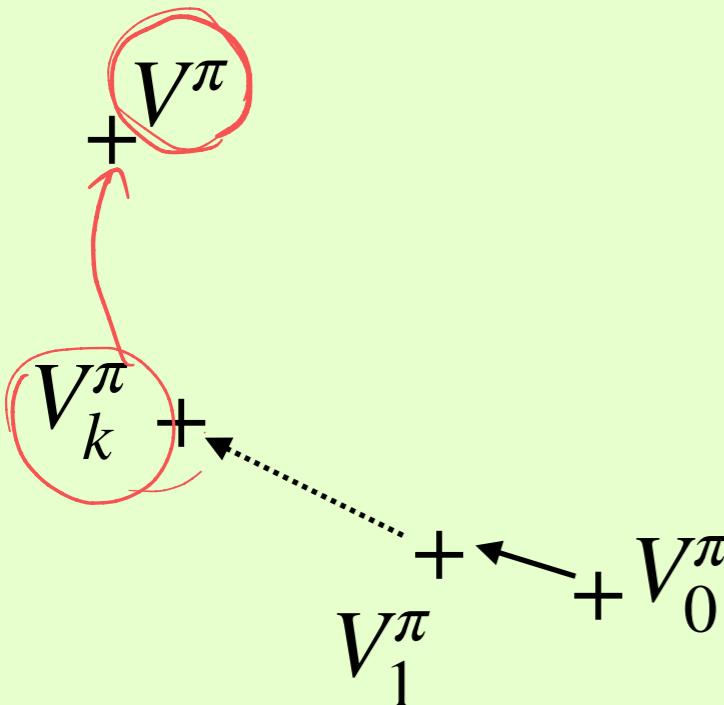
1. Initialize  $V_0^\pi(s) = 0$  for all  $s$

2. For  $k = 1, 2, \dots$

$$V_k^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a | s) \left( R_{s,a} + \gamma \sum_{s' \in \mathcal{S}} P_{ss'}^a V_{k-1}^\pi(s') \right) \text{ for all } s$$

- ▶ Question: What if we start from  $V_0^\pi(s) = V^\pi(s), \forall s$ ?
- ▶ Question: In general, does  $V_k^\pi(s)$  converge to the correct  $V^\pi(s)$ ?

# (Complete) Metric vector space $\mathbb{R}^{|\mathcal{S}|}$

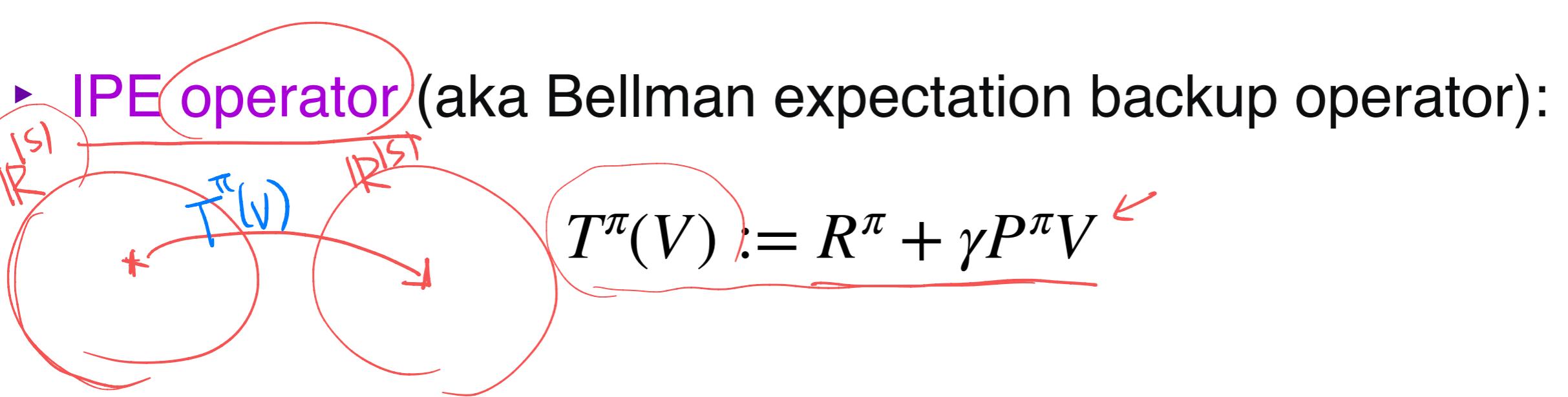


- ▶ **Question:** What does IPE do to points in this space?

Prove convergence in 2 steps:

- (A1): IPE brings points **closer** (formally, a contraction operator)
- (A2): Under any contraction operator, the points converges to a unique fixed point

# For (A1): IPE is a Contraction Map



- ▶ IPE operator (aka Bellman expectation backup operator):  
$$T^\pi(V) := R^\pi + \gamma P^\pi V$$
- ▶ Consider  $L_\infty$ -norm to measure distance between any two value functions  $V, V'$

$$\|V - V'\|_\infty := \max_{s \in \mathcal{S}} |V(s) - V'(s)|$$

# For (A1): IPE is a Contraction Map (Cont.)

- IPE operator:  $T^\pi(V) = R^\pi + \gamma P^\pi V$

For any two value functions  $V$  and  $V'$ ,

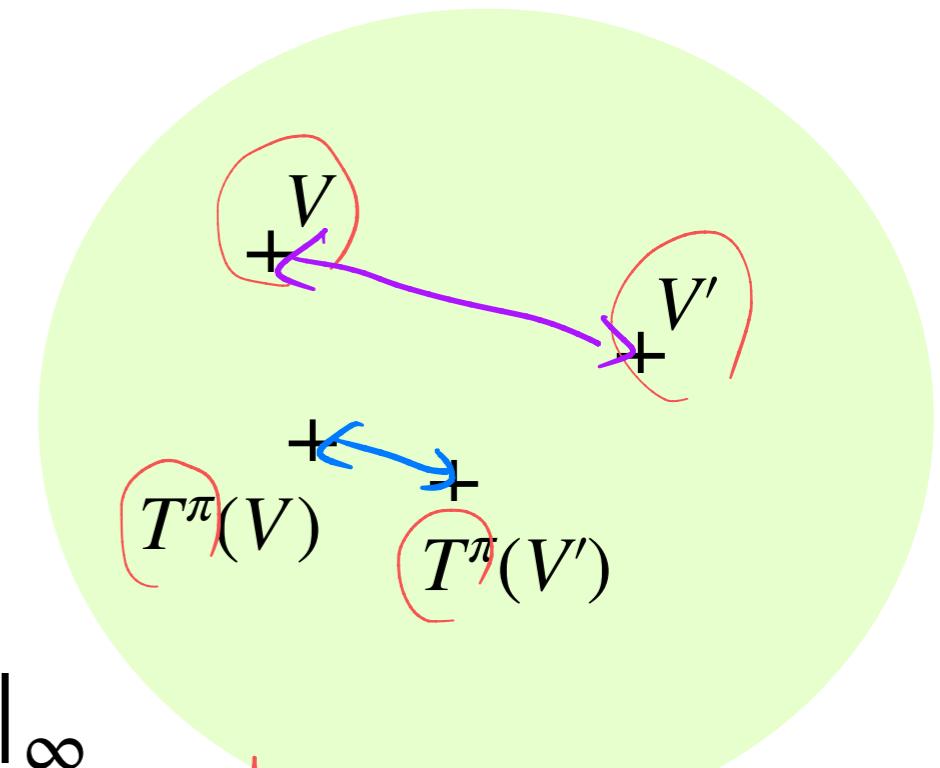
$$\|T^\pi(V) - T^\pi(V')\|_\infty$$

$$= \|(R^\pi + \gamma P^\pi V) - (R^\pi + \gamma P^\pi V')\|_\infty$$

$$= \gamma \|P^\pi(V - V')\|_\infty$$

$$\leq \gamma \|V - V'\|_\infty$$

We say  $T^\pi$  is a  $\gamma$ -contraction operator ( $\gamma < 1$ ).



$$\begin{aligned} & \|T^\pi(V) - T^\pi(V')\|_\infty \\ &= \|(R^\pi + \gamma P^\pi V) - (R^\pi + \gamma P^\pi V')\|_\infty \\ &= \gamma \|P^\pi(V - V')\|_\infty \\ &\leq \gamma \|V - V'\|_\infty \end{aligned}$$

$|S| = 5$

$$\begin{aligned} & \left[ \begin{array}{c} V(s^{(1)}) - V'(s^{(1)}) \\ \vdots \\ V(s^{(5)}) - V'(s^{(5)}) \end{array} \right] \\ & \quad \left[ \begin{array}{c} P_{s^{(1)}, \cdot}^\pi \\ \vdots \\ P_{s^{(5)}, \cdot}^\pi \end{array} \right]^\top \left[ \begin{array}{c} V - V' \\ \vdots \\ V - V' \end{array} \right] \\ & \quad \leq \max_{s \in S} |V(s) - V'(s)| \end{aligned}$$

# For (A2): Banach Fixed-Point Theorem

- ▶ **Banach Fixed-Point Theorem:** For any non-empty complete metric space, if  $T$  is a  $\gamma$ -contraction operator, then  $T$  has a unique fixed point.

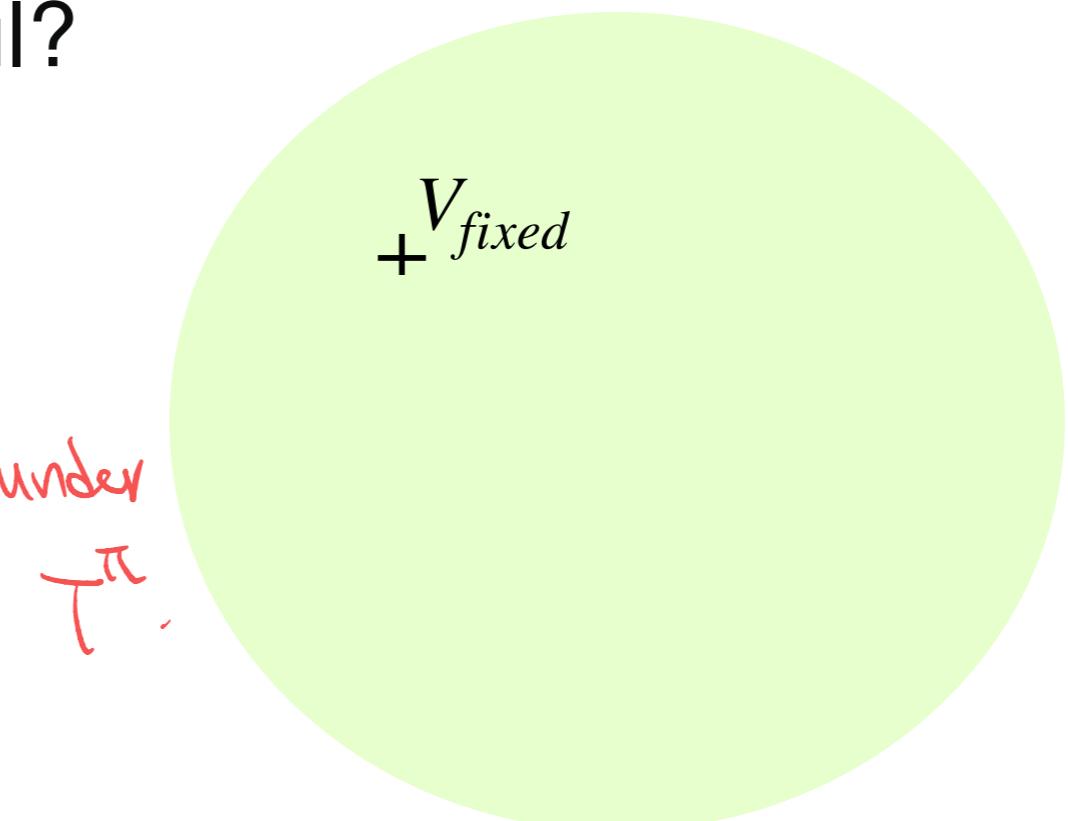
$$T(V) = V$$

- ▶ **Question:** Why is this useful?

By Banach Fixed-Point Theorem,

true  $V^\pi$  is our "unique" fixed point under  $T^\pi$ .

$+V_{fixed}$



$$(A_1): \quad \| T^\pi(v) - T^\pi(v') \|_\infty \leq \gamma \cdot \| v - v' \|_\infty$$

(A<sub>2</sub>):  $V^\pi$  is the unique fixed point (i.e.,  $T^\pi(V^\pi) = V^\pi$ )

Show:  $V_k^\pi \rightarrow V^\pi$

Set  $V = V_k^\pi, V' = V^\pi$

$$\| \underbrace{T^\pi(V_k^\pi)}_{V_{k+1}^\pi} - \underbrace{T^\pi(V^\pi)}_{V^\pi \text{ by (A}_2\text{)}} \|_\infty \leq \gamma \cdot \| V_k^\pi - V^\pi \|_\infty \cdots \text{ by (A}_1\text{)}$$

$$\Rightarrow \underbrace{\| V_{k+1}^\pi - V^\pi \|_\infty}_{\leq \gamma^{k+1} \cdot \| V_0^\pi - V^\pi \|_\infty}$$

# Quick Summary

- ▶ Under IPE, the value functions  $V_k^\pi$  converges to the correct  $V_k^\pi$ , for any initialization  $V_0^\pi$

# Contraction property is central to various RL algorithms:

## A Theory of Regularized Markov Decision Processes

Matthieu Geist<sup>1</sup> Bruno Scherrer<sup>2</sup> Olivier Pietquin<sup>1</sup>

### Abstract

Many recent successful (deep) reinforcement learning algorithms make use of regularization, generally based on entropy or Kullback-Leibler divergence. We propose a general theory of regularized Markov Decision Processes that generalizes these approaches in two directions: we consider a larger class of regularizers, and we consider the general modified policy iteration approach, encompassing both policy iteration and value iteration. The core building blocks of this theory are a notion of regularized Bellman operator and the Legendre-Fenchel transform, a classical tool of convex optimization. This approach allows for error propagation analyses of general algorithmic schemes of which (possibly variants of) classical algorithms such as Trust Region Policy Optimization, Soft Q-learning, Stochastic Actor Critic or Dynamic Policy Programming are special cases. This also draws connections to proximal convex optimization, especially to Mirror Descent.

Tsallis entropy (Lee et al., 2018) having a sparse regularized greedy are based on a notion of tempo somehow extending the notion of regularized case (Nachum et al. Nachum et al., 2018), or on policy Mnih et al., 2016).

This non-exhaustive set of algorithms regularization, but they are different principles, consider each iteration, and have ad-hoc analysis, a general theory of regularized M (MDPs). To do so, a key observation is dynamic programming, or (A)L from the core definition of the Bellman operator. The framework we propose is based on the Bellman operator, and on an associated transform. We study the theoretical properties of regularized MDPs and of the related reinforcement learning algorithms. This generalizes many existing theories and provides new ones. Notably, it allows for a unified analysis of various reinforcement learning algorithms.

(Geist et al., ICML 2019)

## A Generalized Algorithm for Multi-Objective Reinforcement Learning and Policy Adaptation

Runzhe Yang  
Department of Computer Science  
Princeton University  
runzhey@cs.princeton.edu

Xingyuan Sun  
Department of Computer Science  
Princeton University  
xs5@cs.princeton.edu

Karthik Narasimhan  
Department of Computer Science  
Princeton University  
karthikn@cs.princeton.edu

### Abstract

We introduce a new algorithm for multi-objective reinforcement learning (MORL) with linear preferences, with the goal of enabling few-shot adaptation to new tasks. In MORL, the aim is to learn policies over multiple competing objectives whose relative importance (*preferences*) is unknown to the agent. While this alleviates dependence on scalar reward design, the expected return of a policy can change significantly with varying preferences, making it challenging to learn a single model to produce optimal policies under different preference conditions. We propose a generalized version of the Bellman equation to learn a single parametric representation for optimal policies over the space of all possible preferences. After an initial learning phase, our agent can execute the optimal policy under any given preference, or automatically infer an underlying preference with very few samples. Experiments across four different domains demonstrate the effectiveness of our approach.<sup>[1]</sup>

(Yang et al., NeurIPS 2019)

Now we know how to evaluate a given policy  $\pi(a | s)$

How about an **optimal** policy?

# Optimal Value Functions

- ▶ **Optimal state-value function:**

$$V^*(s) = \max_{\pi \in \Pi} V^\pi(s)$$

$\Pi$  ↗ *Markov stationary policies*

- ▶ **Optimal action-value function:**

$$Q^*(s, a) = \max_{\pi \in \Pi} Q^\pi(s, a)$$

- ▶ **Question:** Suppose we know  $Q^*(s, a)$  for all pairs of  $(s, a)$ . What's your guess about the “optimal” policy (though we haven't defined it yet)?

$$\bar{\pi}^*(s) := \operatorname{argmax}_{a \in A} Q^*(s, a)$$

# Example: Shortest Path

(0,0)	(0,1)	(0,2)	(0,3)
			(3,3)

Problem

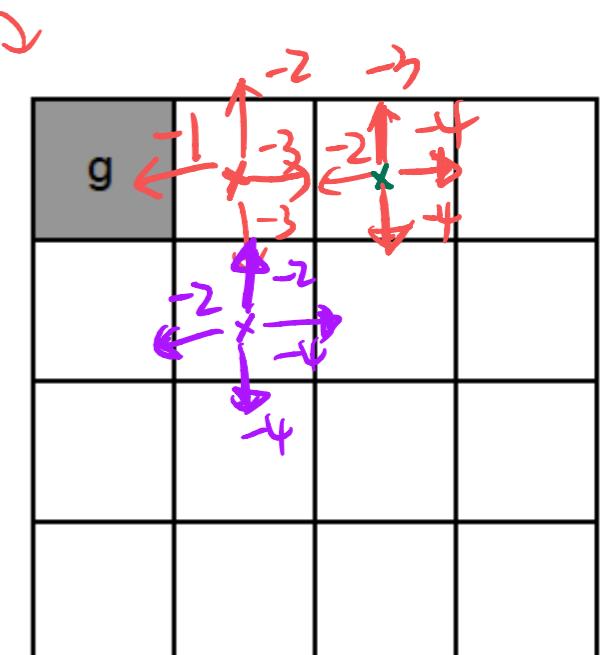
$$\pi \quad V^{\pi}(3,3) \quad V^*(3,3)$$

- Actions: Left / Right / Up / Down
- The learner stays at the same place after hitting a wall
- The episode ends after reaching the goal (denoted as  $g$ )
- Reward:  $R(s, a) = -1$ , for all  $(s, a)$
- In this example,  $V^*(s)$  captures the shortest distance from  $s$  to  $g$

0	-1	-2	-3
-1	-2	-3	-4
-2	-3	-4	-5
-3	-4	-5	-6

$V^*(s)$

Check:  $\pi^*(s) = \arg\max_{a \in A} Q^*(s, a)$



$Q^*(s, a)$

# Definition of an “Optimal” Policy (Formally)

- ▶ “Partial ordering” of policies:

$$\underline{\pi} \geq \underline{\pi'} \quad \text{if } \circled{V^\pi(s)} \geq V^{\pi'}(s), \forall s$$

- ▶ An optimal policy: A policy  $\pi^*$  is an optimal policy if it is better than or equal to all other policies, i.e.

$$\circled{\pi^* \geq \pi}, \text{ for all } \underline{\pi \in \Pi}$$

- ▶ **Question**: Given an MDP, does such  $\pi^*$  always exist? *Yes!*

# Existence of an Optimal Policy

- ▶ **Theorem (Existence of Optimal Policy):** For any finite MDP, there always exists an optimal policy  $\pi^*$  such that

$$\pi^* \geq \pi, \text{ for all } \pi \in \Pi$$

(This will be shown shortly)

[Puterman, 1994]

- ▶ A deterministic optimal policy can always be found by

$$\underline{\pi^*(a | s)} = \begin{cases} 1, & \text{if } a = \arg \max_{a' \in \mathcal{A}} Q^*(s, a') \\ 0, & \text{else} \end{cases}$$

(To be proved shortly)

- ▶ Therefore, knowing  $Q^*$  is equivalent to knowing  $\pi^*$
- ▶ **Question:** Number of deterministic policies?

How to find  $Q^*(s, a)$  (and thereafter  $\pi^*$ )?

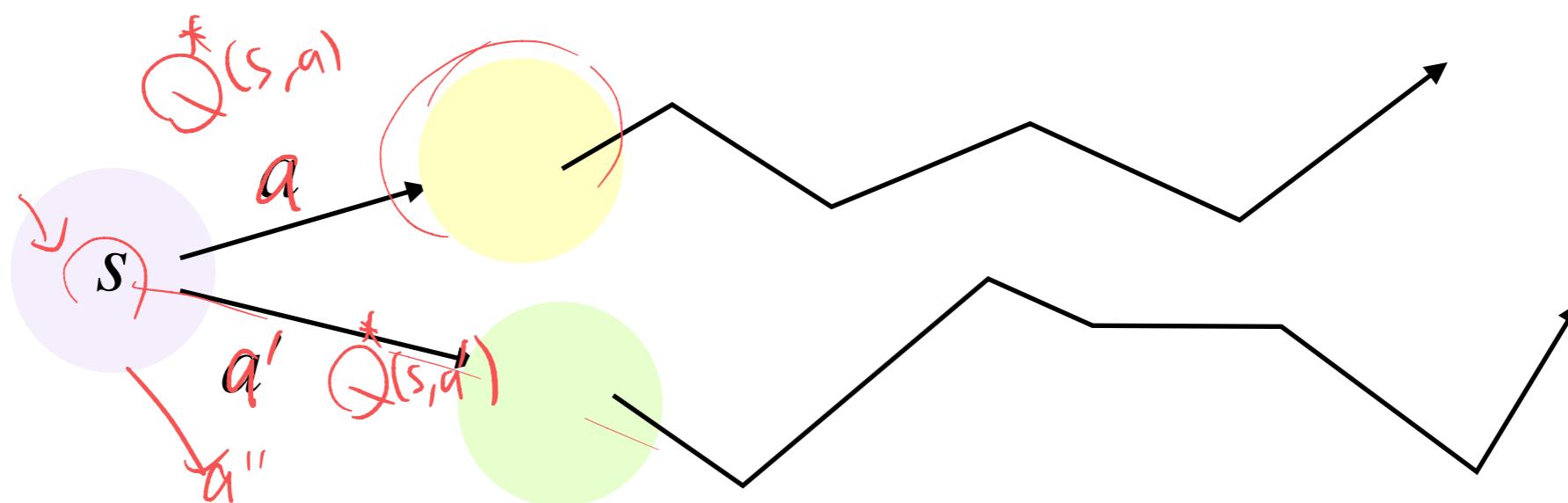
# Bellman "Optimality" Equations

(1)  $\underline{V^*}$  written in  $\underline{Q^*}$

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

(2)  $Q^*$  written in  $V^*$

$$Q^*(s, a) = R_s^a + \gamma \sum_{s'} P_{ss'}^a V^*(s')$$



(Proof: HW1 problems)

# Bellman Optimality Equations (Cont.)

(1)  $V^*$  written in  $Q^*$

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

(2)  $Q^*$  written in  $V^*$

$$Q^*(s, a) = R_s^a + \gamma \sum_{s'} P_{ss'}^a V^*(s')$$

(3)  $V^*$  written in  $V^*$

✓  $\underline{V^*(s)} = \max_a R_s^a + \gamma \sum_{s'} P_{ss'}^a \underline{V^*(s')}$

(4)  $Q^*$  written in  $Q^*$

$$\underline{Q^*(s, a)} = R_s^a + \gamma \sum_{s'} P_{ss'}^a \left( \max_{a'} \underline{Q^*(s', a')} \right)$$

# How to Solve the Bellman Optimality Equation?

$$V^*(s) = \max_a \left( R_s^a + \gamma \sum_{s'} P_{ss'}^a V^*(s') \right)$$

- ▶ **Question:** Solve this by linear algebra?
- ▶ The “max” operation makes it non-linear
- ▶ We need to resort to iterative methods:
  - ▶ Policy iteration
  - ▶ Value iteration

# Value Iteration (VI)

# From Bellman Optimality Equation to Bellman Optimality Backup Operator

- ▶ **Recall:** Bellman optimality equation

$$V^*(s) = \max_a \left( R_s^a + \gamma \sum_{s'} P_{ss'}^a V^*(s') \right)$$

- ▶ **Define:** Bellman optimality backup operator

$$T^*(V) := \max_{a \in \mathcal{A}} R^a + \gamma P^a V$$

(Comparison: IPE operator  $T^\pi(V) = R^\pi + \gamma P^\pi V$ )

# Value Iteration: Pseudo Code

1. Initialize  $k = 0$  and set  $V_0(s) = 0$  for all states
2. Repeat the following until convergence:

$$V_{k+1} \leftarrow T^*(V_k)$$

- ▶ Equivalently: for each state  $s$

$$V_{k+1}(s) = \max_a \left( R_s^a + \gamma \sum_{s'} P_{ss'}^a V_k(s') \right)$$

- ▶ **Remark:** Complexity per iteration is  $O(|\mathcal{S}|^2 |\mathcal{A}|)$
- ▶ **Remark:** Intermediate value functions  $V_k$  may not correspond to any policy

# Example: Shortest Path

$$V_{k+1}(s) = \max_a \left( R_s^a + \gamma \sum_{s'} P_{ss'}^a V_k(s') \right)$$

<b>g</b>			

Problem

0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

$V_1$

0	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	-1

$V_2$

0	-1	-2	-2
-1	-2	-2	-2
-2	-2	-2	-2
-2	-2	-2	-2

$V_3$

0	-1	-2	-3
-1	-2	-3	-3
-2	-3	-3	-3
-3	-3	-3	-3

$V_4$

0	-1	-2	-3
-1	-2	-3	-4
-2	-3	-4	-4
-3	-4	-4	-4

$V_5$

0	-1	-2	-3
-1	-2	-3	-4
-2	-3	-4	-5
-3	-4	-5	-5

$V_6$

0	-1	-2	-3
-1	-2	-3	-4
-2	-3	-4	-5
-3	-4	-5	-6

$V_7$

# Convergence of Value Iteration

- ▶ **Value iteration converges on  $V^*$ :** For any initial  $V_0$ , Value Iteration achieves that  $V_k \rightarrow V^*$  (in  $L_\infty$ -norm)
- ▶ **Question:** How to show this?

# Convergence of Value Iteration

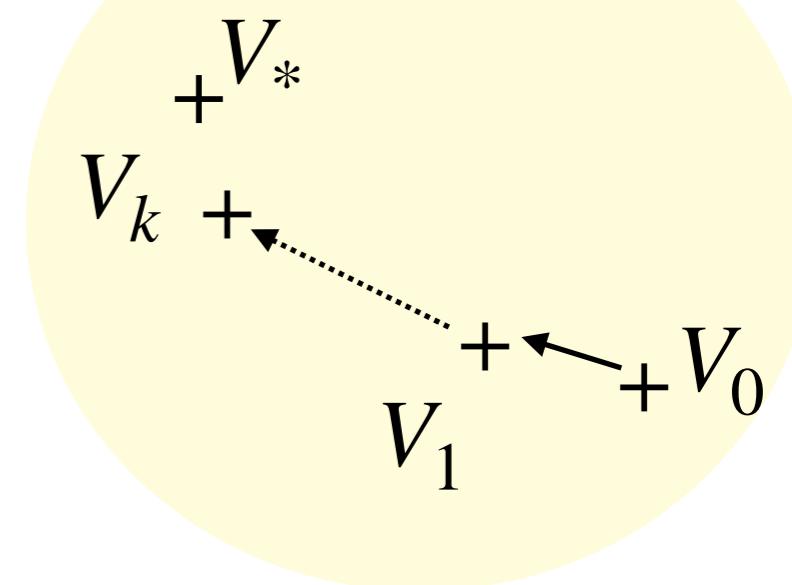
- ▶ **Value iteration converges on  $V^*$ :** For any initial  $V_0$ , Value Iteration achieves that  $V_k \rightarrow V^*$  (in  $L_\infty$ -norm)
- ▶ **Proof:** We prove convergence via the following steps.

(B1): Show that  $T^*$  is a contraction operator

(B2): Under a contraction operator,  $\{V_k\}$  shall converge to the unique fixed point (why?)

(B3): Since  $V^*$  is a fixed point, then  $V_k \rightarrow V^*$  due to uniqueness

Complete metric vector space  $\mathbb{R}^{|\mathcal{S}|}$



# For (B1): $T^*$ is a $\gamma$ -Contraction Operator on $V$

- Bellman optimality backup operator:  $T^*(V) = \max_{a \in \mathcal{A}} (R^a + \gamma P^a V)$

$$\|T^*(V) - T^*(\hat{V})\|_\infty$$

$$= \max_s |T^*(V)(s) - T^*(\hat{V})(s)|$$

$$= \max_s \left| \max_a \left( R_s^a + \gamma \sum_{s'} P_{ss'}^a V(s') \right) - \max_{a'} \left( R_s^{a'} + \gamma \sum_{s'} P_{ss'}^{a'} \hat{V}(s') \right) \right|$$

$$\leq \max_s \max_a \left| \left( R_s^a + \gamma \sum_{s'} P_{ss'}^a V(s') \right) - \left( R_s^a + \gamma \sum_{s'} P_{ss'}^a \hat{V}(s') \right) \right|$$

$$\leq \max_s \max_a \left| \gamma \sum_{s'} P_{ss'}^a (V(s') - \hat{V}(s')) \right|$$

$$\leq \gamma \|V - \hat{V}\|_\infty$$

Therefore,  $T^*$  is a  $\gamma$ -contraction operator ( $\gamma < 1$ )

## For (B2)

- ▶  $T^*$  is a  $\gamma$ -contraction operator in a complete metric space
- ▶ By Banach Fixed Point Theorem,  $T^*$  converges to the unique fixed point

# Discussion: Issues With Value Iteration

- ▶ **Question 1:** In how many iterations will VI converge?
- ▶ **Question 2:** By applying VI, could we find an optimal policy?
- ▶ **Question 3:** By using VI, could we directly prove the **existence** of an optimal policy?

# Some Historical Accounts

A seminal paper on “Value Iteration”  
(published in 1957)

## *A Markovian Decision Process*

RICHARD BELLMAN

**1. Introduction.** The purpose of this paper is to discuss the asymptotic behavior of the sequence  $\{f_N(i)\}$ ,  $i = 1, 2, \dots, M$ ,  $N = 1, 2, \dots$ , generated by the non-linear recurrence relations

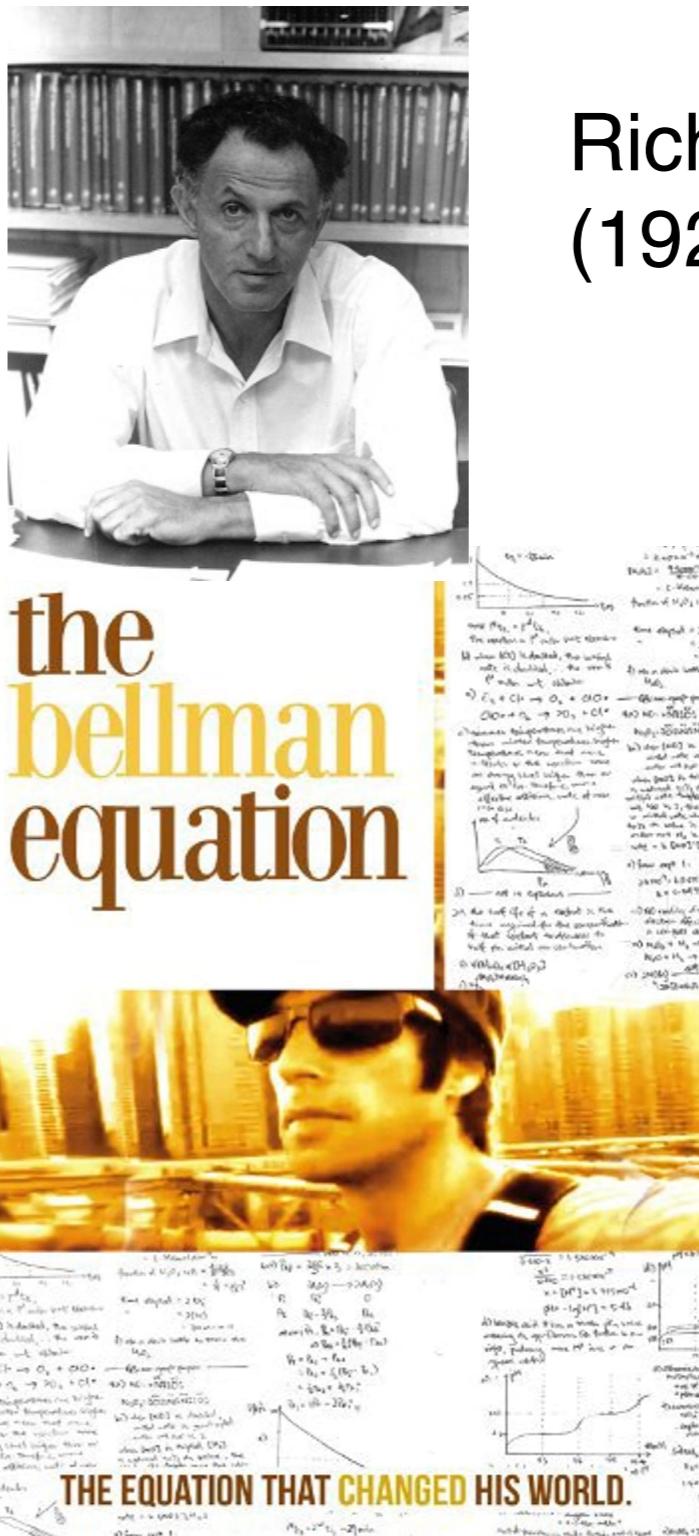
$$(1) \quad f_N(i) = \max_q \left[ b_i(q) + \sum_{j=1}^M a_{ij}(q) f_{N-1}(j) \right], \quad N = 1, 2, \dots,$$
$$f_0(i) = c_i, \quad i = 1, 2, \dots, M.$$

Although these equations are non-linear, they possess certain quasi-linear properties which permit a more thorough discussion than might be imagined upon first glance.

As we shall discuss below, this question arises from the consideration of a dynamic programming process. A related process gave rise to an equation of the above form which was discussed in [1], under a particular set of assumptions concerning the functions  $b_i(q)$  and the matrices  $A(q) = (a_{ij}(q))$ . Here we shall impose restrictions of a quite different type. Any complete discussion of relations of the foregoing type is at least as detailed as a corresponding discussion of the linear case, and, as in the linear case, the assumptions made determine the techniques employed to a considerable extent.

We shall discuss elsewhere some interesting quadratically non-linear recurrence relations which arise from specializing the forms of  $b_i(q)$  and  $a_{ij}(q)$ . These are related to the types of differential equations discussed in [3], being essentially a particular type of discrete version.

It will be clear from what follows that similar techniques can be utilized to treat the determination of the asymptotic behavior of the sequences defined by

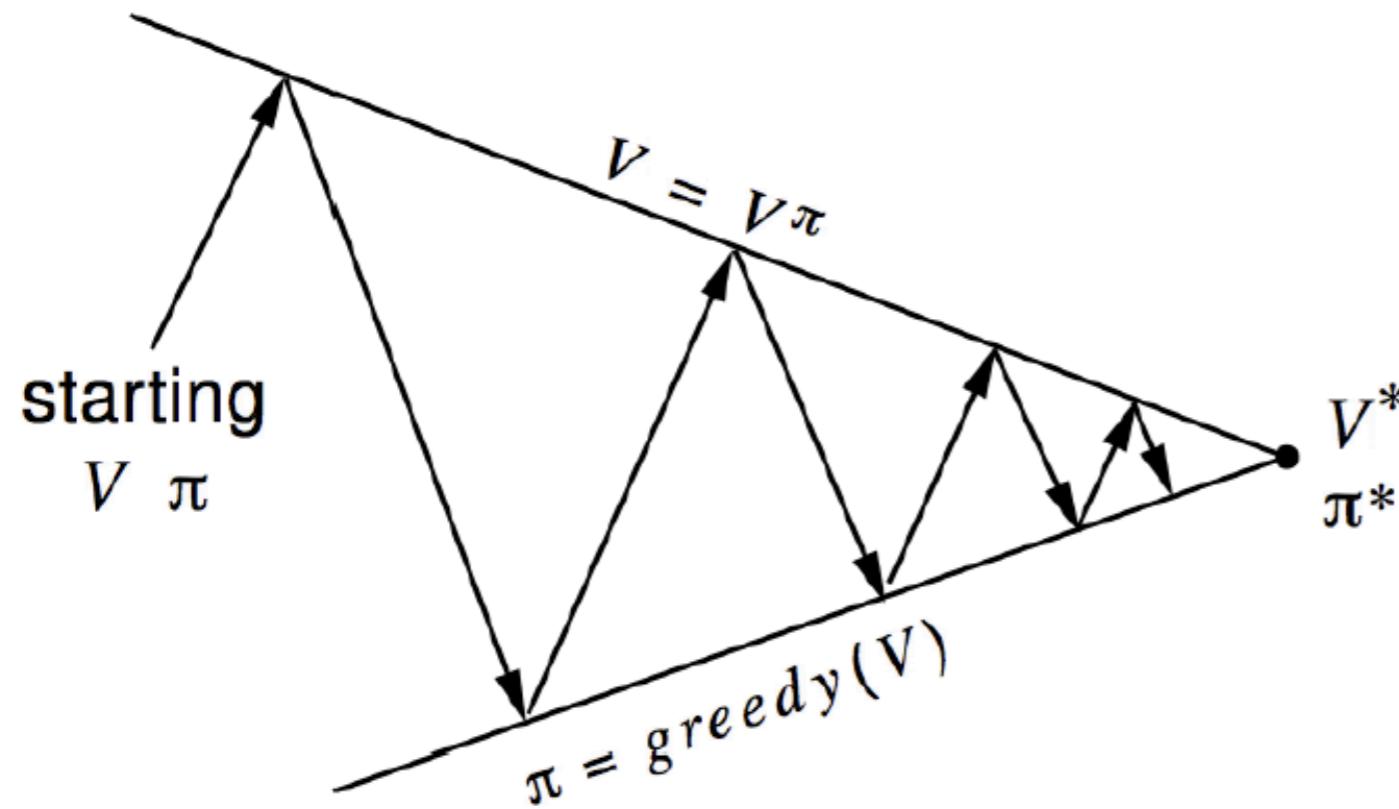


Richard Bellman  
(1920-1984)

A documentary released in 2011

# Policy Iteration (PI)

# Policy Iteration: Generic Procedure

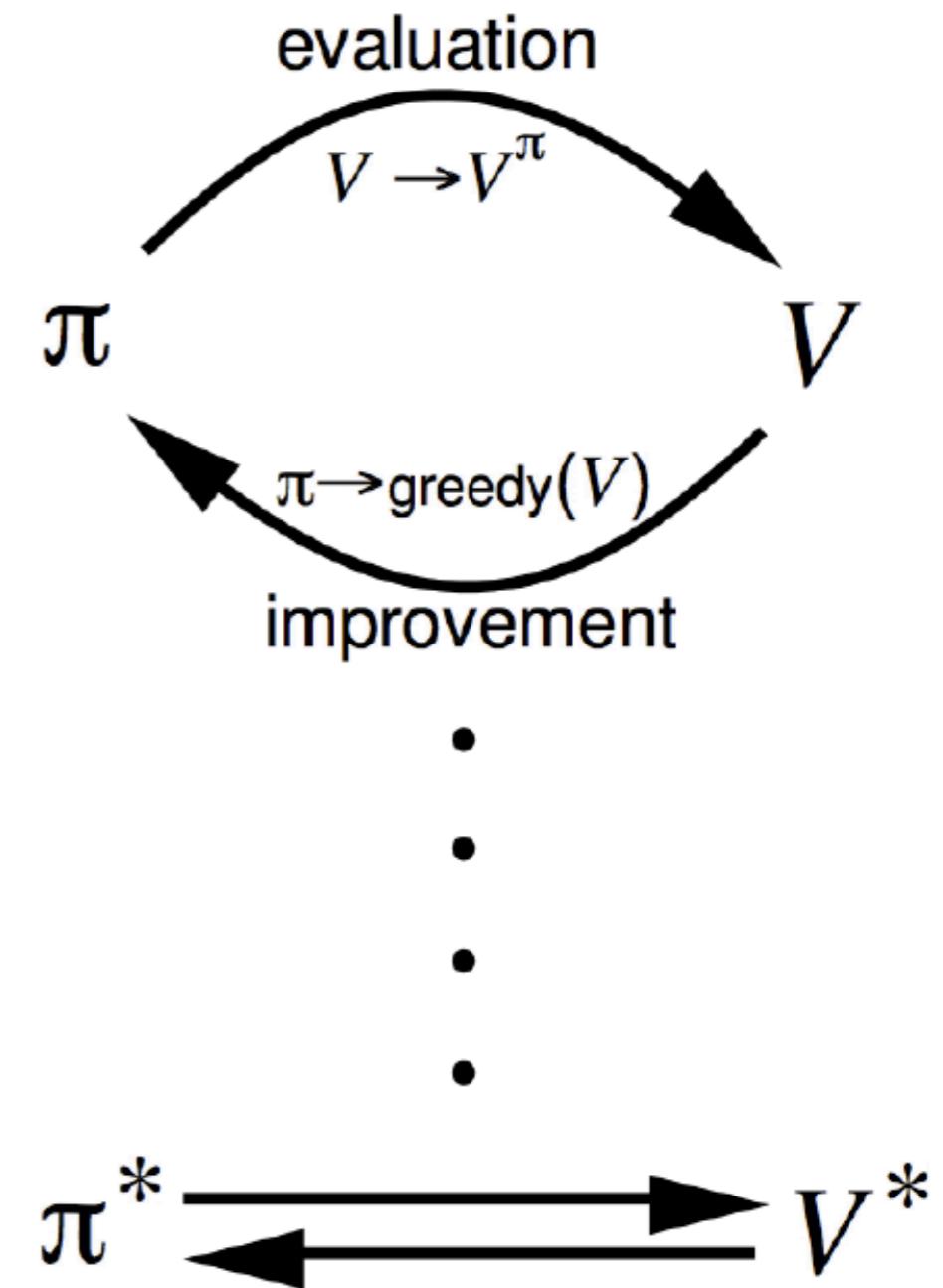


Policy evaluation Estimate  $v_\pi$

Iterative policy evaluation

Policy improvement Generate  $\pi' \geq \pi$

Greedy policy improvement



# Policy Iteration: Pseudo Code

(We focus on deterministic policies)

1. Initialize  $k = 0$  and set  $\pi_0(s)$  arbitrarily for all states
2. While  $k$  is zero or  $\pi_k \neq \pi_{k-1}$ :
  - ▶ Derive  $V^{\pi_k}$  via **policy evaluation** for  $\pi_k$  (iterative/non-iterative)
  - ▶ Derive  $\pi_{k+1}$  by greedy **one-step policy improvement**

# One-Step Policy Improvement

- Given  $V^{\pi_k}$ , compute  $Q^{\pi_k}(s, a)$ :

$$Q^{\pi_k}(s, a) = R(s, a) + \gamma \sum_s P(s' | s, a) V^{\pi_k}(s')$$

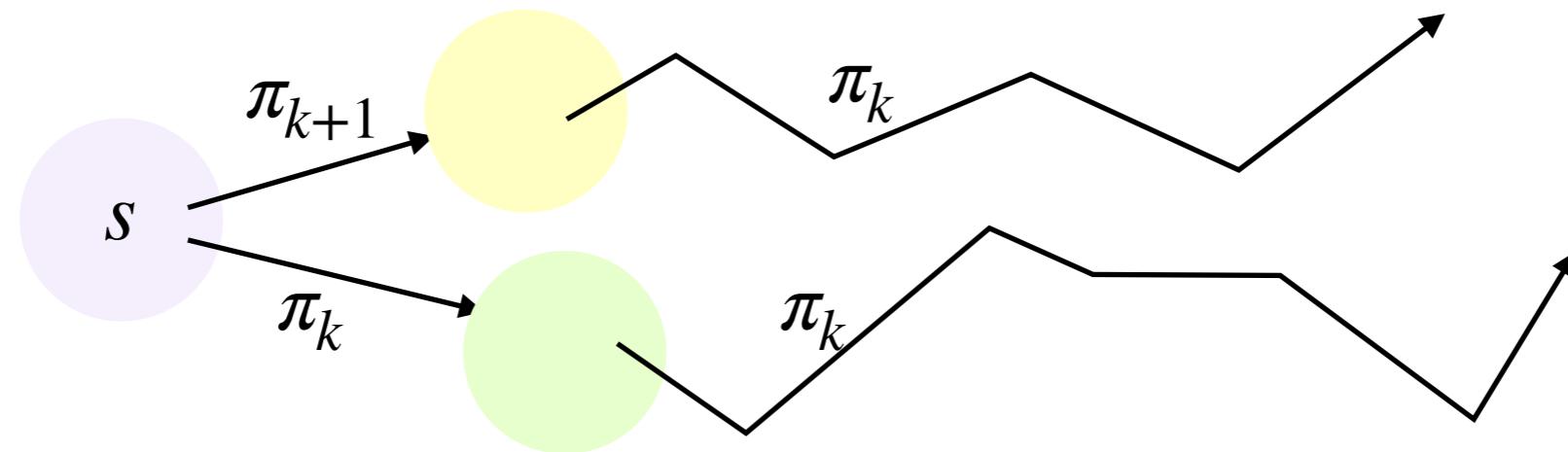
- Derive the new policy  $\pi_{k+1}$ : For all states  $s$ ,

$$\pi_{k+1}(s) = \arg \max_a Q^{\pi_k}(s, a)$$

- Note: We will use  $R(s, a)$  and  $R_s^a$  interchangeably

# Why is One-Step Policy Improvement Reasonable?

- ▶ **Question:** Suppose we take  $\pi_{k+1}(s)$  for one step and then follow  $\pi_k$  subsequently. Is this better than just following  $\pi_k$ ?



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

$$\max_a Q^{\pi_k}(s, a) \geq R(s, \pi_k(s)) + \gamma \sum_s P(s' | s, a) V^{\pi_k}(s') = V^{\pi_k}(s)$$

$$\pi_{k+1}(s) = \arg \max_a Q^{\pi_k}(s, a)$$

- ▶ **Question:** But how about following  $\pi_{k+1}$  all the way?

# Monotonic Improvement in Policy

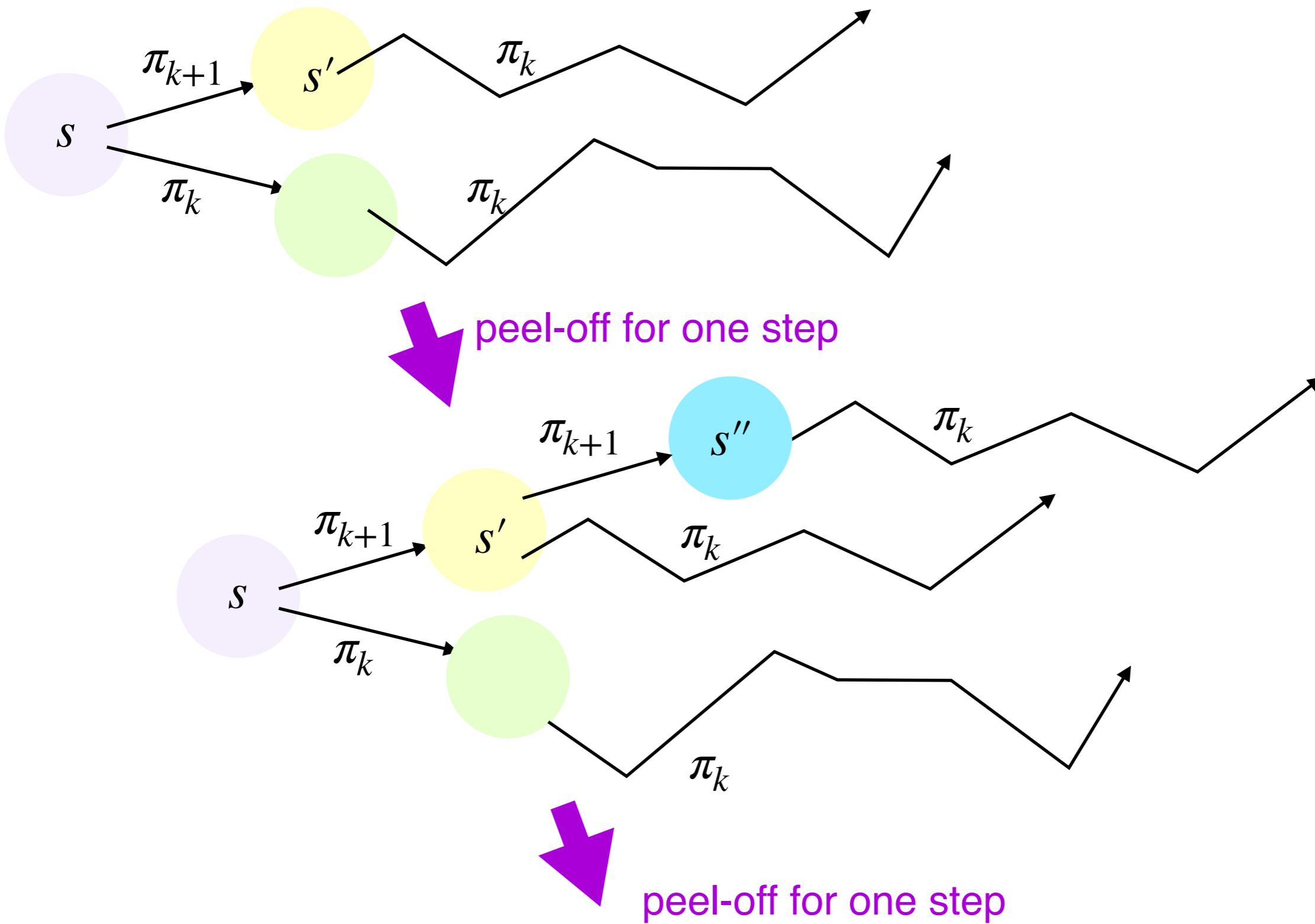
- ▶ Recall: Partial ordering of policies

$$\pi \geq \pi' \quad \text{if} \quad V^\pi(s) \geq V^{\pi'}(s), \forall s$$

- ▶ Question: Do we have  $\pi_{k+1} \geq \pi_k$  ?

- ▶ **Theorem (Monotonic Policy Improvement):** Under the one-step policy improvement step, we have  $V^{\pi_{k+1}}(s) \geq V^{\pi_k}(s), \forall s$  and hence  $\pi_{k+1} \geq \pi_k$

## Proof Idea: “Peeling off”



# Proof: Monotonic Policy Improvement

$$\begin{aligned} V^{\pi_k}(s) &\leq \max_a Q^{\pi_k}(s, a) \\ &= \max_a R(s, a) + \gamma \sum_{s'} P(s' | s, a) V^{\pi_k}(s') \\ &= R(s, \pi_{k+1}(s)) + \gamma \sum_{s'} P(s' | s, \pi_{k+1}(s)) V^{\pi_k}(s') \\ &\leq R(s, \pi_{k+1}(s)) + \gamma \sum_{s'} P(s' | s, \pi_{k+1}(s)) \max_{a'} Q^{\pi_k}(s', a') \\ &= R(s, \pi_{k+1}(s)) + \gamma \sum_{s'} P(s' | s, \pi_{k+1}(s)) \\ &\quad \times \left( R(s', \pi_{k+1}(s')) + \gamma \sum_{s''} P(s'' | s', \pi_{k+1}(s')) V^{\pi_k}(s'') \right) \\ &\quad \dots \\ &= V^{\pi_{k+1}}(s) \end{aligned}$$

# Discussions: Policy Iteration

- ▶ **Question:** If we have  $\pi_{k+1} = \pi_k$ , what shall we expect about  $\pi_{k+2}$ ?

- ▶ **Recall:**

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

$$\pi_{k+1}(s) = \arg \max_a Q^{\pi_k}(s, a)$$

$$\pi_{k+2}(s) = \arg \max_a Q^{\pi_{k+1}}(s, a) = \arg \max_a Q^{\pi_k}(s, a)$$

Therefore, policy iteration can terminate when  $\pi_{k+1} = \pi_k$

Moreover,  $\pi_{k+1} = \pi_k$  implies **Bellman optimality equation** is satisfied by  $\pi_k$ :

$$V^{\pi_k}(s) = \max_a Q^{\pi_k}(s, a)$$

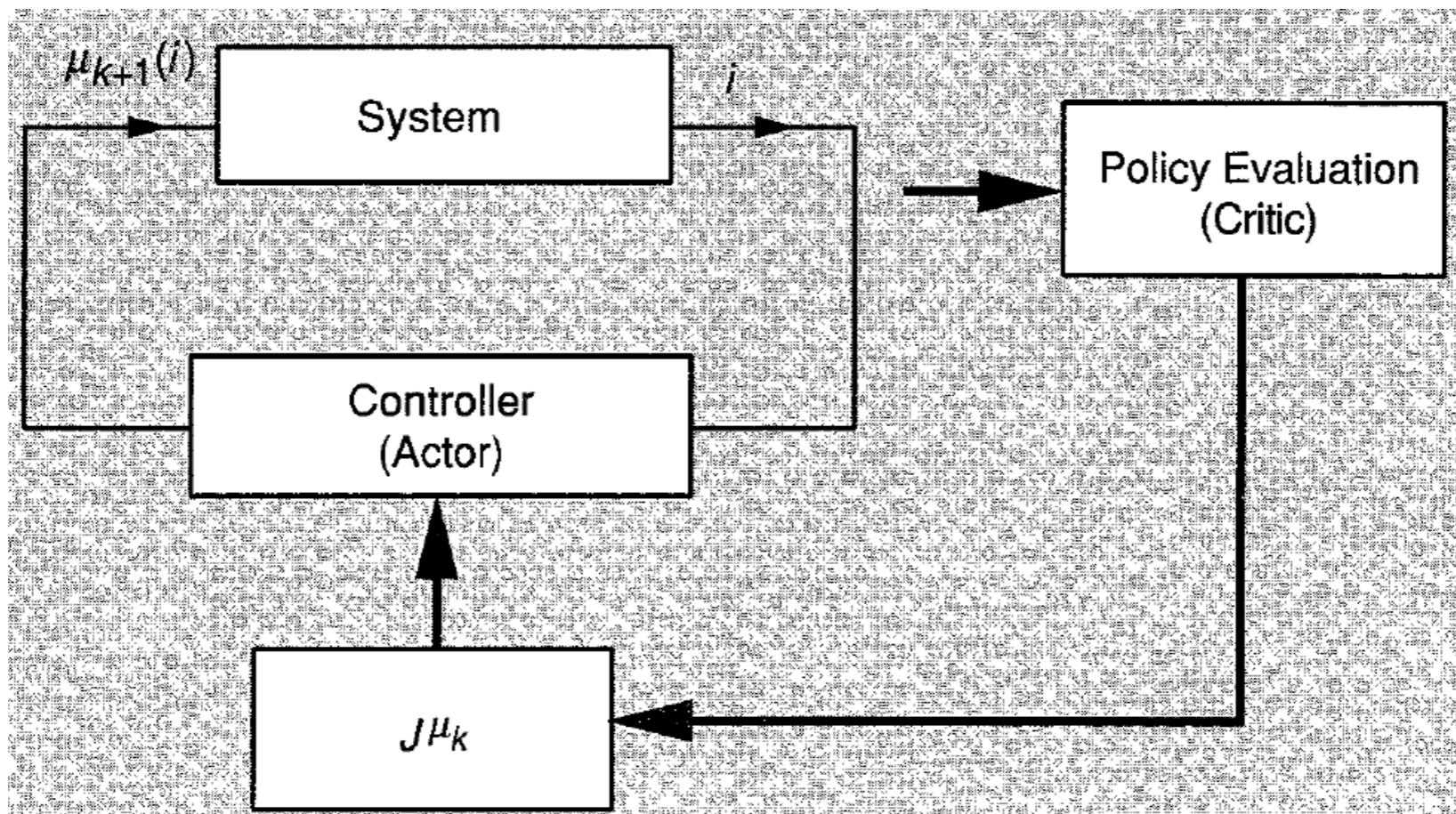
Therefore,  $\pi_k$  must be a (deterministic) optimal policy (**Why?**)

Hence, PI+VI prove the existence of an optimal policy

# Discussions: Policy Iteration (Cont.)

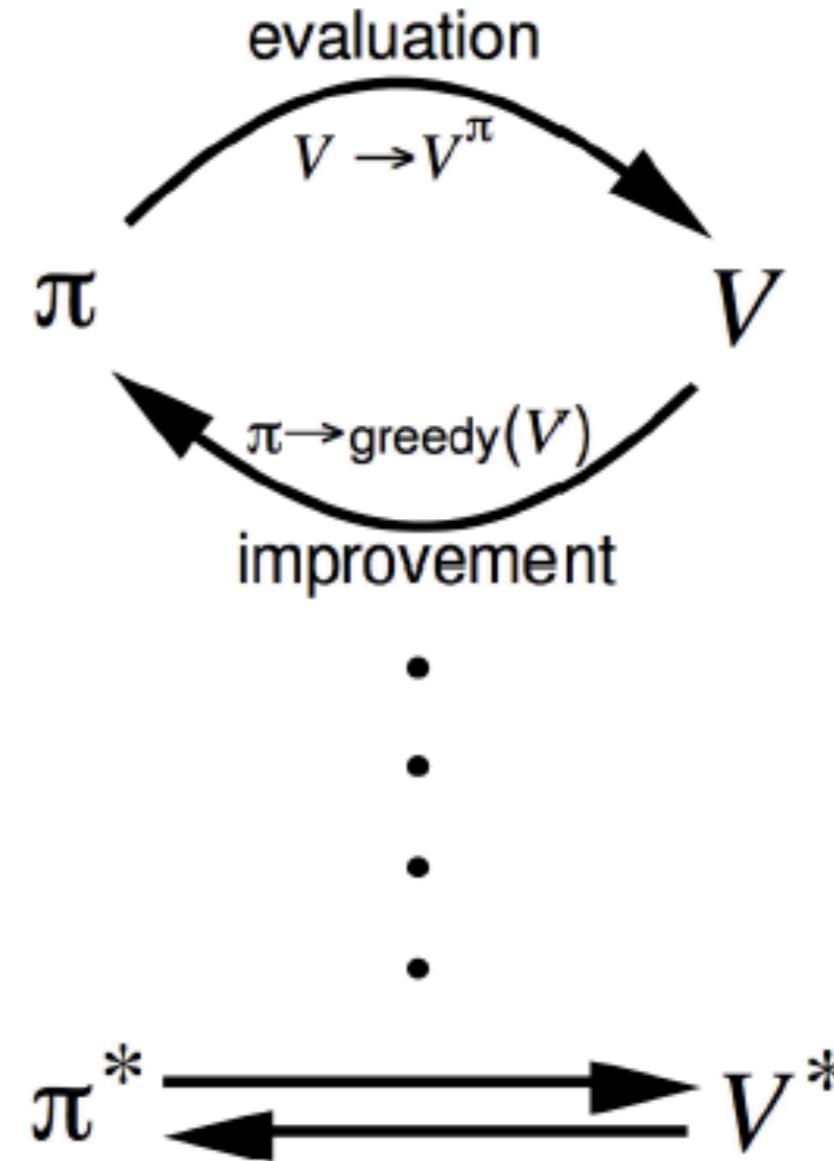
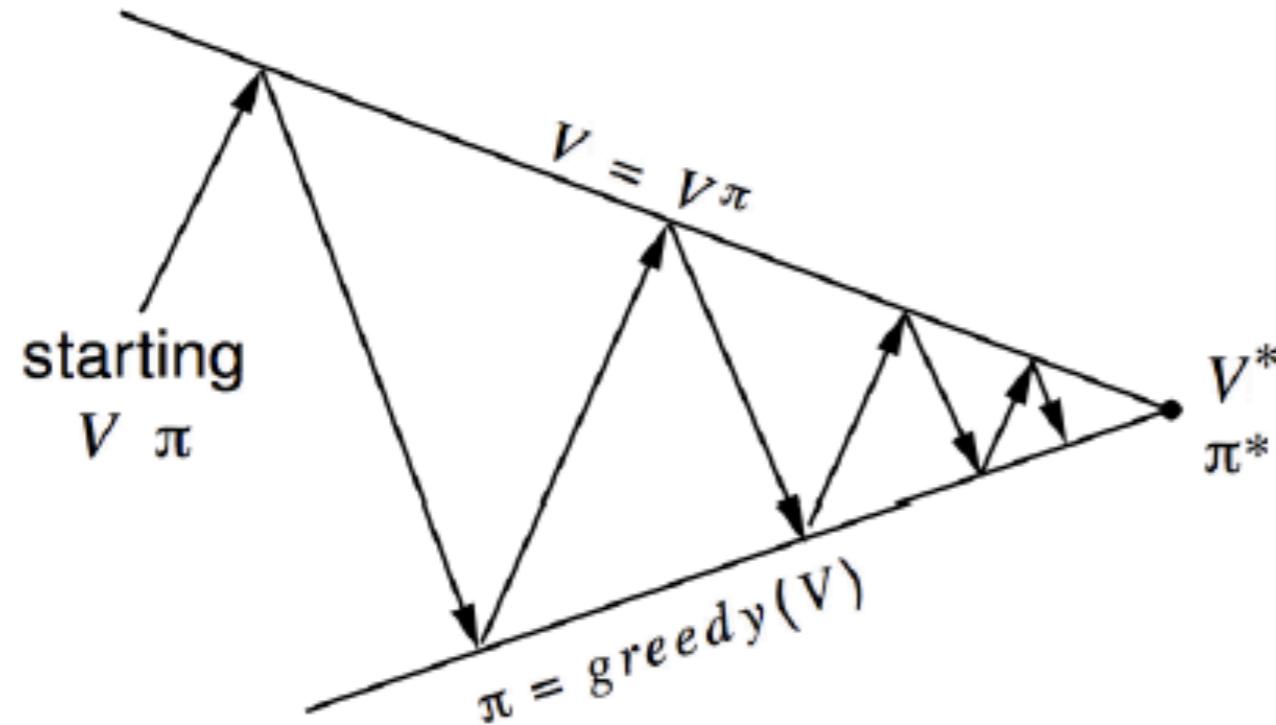
- ▶ **Question:** Will policy iteration terminate in finite number of iterations?
- ▶ Yes, in at most  $|\mathcal{A}|^{|\mathcal{S}|}$  iterations (assume  $|\mathcal{A}|, |\mathcal{S}|$  are finite)
  - (There are  $|\mathcal{A}|^{|\mathcal{S}|}$  deterministic policies)
  - (If  $\pi_{k+1} \neq \pi_k$ , then they must differ by at least 1 entry)
  - (Monotonic policy improvement:  $\pi_{k+1} \geq \pi_k$ )

# Policy Iteration: Actor-Critic Viewpoint



- ▶ **Critic:** evaluate the performance of the current policy
- ▶ **Actor:** take into account the latest evaluation and output an improved policy

# Generalized Policy Iteration



Policy evaluation Estimate  $v_\pi$

Any policy evaluation algorithm

Policy improvement Generate  $\pi' \geq \pi$

Any policy improvement algorithm

- ▶ Remark: Policy gradient methods can be interpreted as an instance of generalized policy iteration (discussed in Lectures 5-7)

# Summary

Problem	Bellman Equation	Algorithm
Prediction	Bellman Expectation Equation	Iterative Policy Evaluation
Control	Bellman Expectation Equation + Greedy Policy Improvement	Policy Iteration
Control	Bellman Optimality Equation	Value Iteration