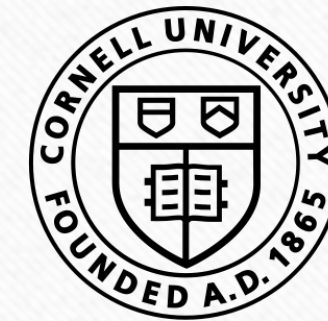


Value Iteration

**The Learning
Strategies
Center**

Find Study Partners!



Studying with peers is a great way to connect with other Cornell students and is a powerful tool for learning.

Cornell's Learning Strategies Center ([LSC](#)) helps match you with study partners.

To learn more, visit the LSC's [Studying Together webpage](#) or scan the code →



Scan the QR code to find out more about Study Partners

or visit
<http://lsc.cornell.edu/studying-together/>

lsc.cornell.edu

Recap: Bellman Optimality

$$\mathcal{M} = \{S, A, P, r, \gamma\}$$

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \rightarrow [0,1], \quad \gamma \in [0,1)$$

$$\text{Policy } \pi : S \mapsto A$$

Recap: Bellman Optimality

$$\mathcal{M} = \{S, A, P, r, \gamma\}$$

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \rightarrow [0,1], \quad \gamma \in [0,1)$$

$$\text{Policy } \pi : S \mapsto A$$

Bellman Optimality—-the Q version (HW0 problem)

$$Q^*(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[\max_{a' \in A} Q^*(s', a') \right]$$

Recap: Bellman Optimality

$$\mathcal{M} = \{S, A, P, r, \gamma\}$$

$$P : S \times A \mapsto \Delta(S), \quad r : S \times A \rightarrow [0,1], \quad \gamma \in [0,1)$$

$$\text{Policy } \pi : S \mapsto A$$

Bellman Optimality—the Q version (HW0 problem)

$$Q^*(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[\max_{a' \in A} Q^*(s', a') \right]$$

For any $Q : S \times A \rightarrow \mathbb{R}$, if $Q(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[\max_{a'} Q(s', a') \right]$ for all s, a , then $Q(s, a) = Q^*(s, a), \forall s, a$

Recap: Fixed-point solution

Find the fixed point solution of $x^\star = f(x^\star)$, $x \in \mathbb{R}$

Recap: Fixed-point solution

Find the fixed point solution of $x^\star = f(x^\star)$, $x \in \mathbb{R}$

Start with some x_0 , set $x_{t+1} \leftarrow f(x_t)$

Recap: Fixed-point solution

Find the fixed point solution of $x^\star = f(x^\star)$, $x \in \mathbb{R}$

Start with some x_0 , set $x_{t+1} \leftarrow f(x_t)$

Suppose f is contraction, i.e., $\forall x, x', |f(x') - f(x)| \leq \gamma |x' - x|$, $\gamma \in [0, 1)$,
then $x_t \rightarrow x^\star$

Recap: Fixed-point solution

Find the fixed point solution of $x^\star = f(x^\star)$, $x \in \mathbb{R}$

Start with some x_0 , set $x_{t+1} \leftarrow f(x_t)$

Suppose f is contraction, i.e., $\forall x, x', |f(x') - f(x)| \leq \gamma |x' - x|$, $\gamma \in [0, 1)$,
then $x_t \rightarrow x^\star$

For Policy Evaluation (i.e., given \mathcal{M} and π , compute V^π)

Recap: Fixed-point solution

Find the fixed point solution of $x^\star = f(x^\star)$, $x \in \mathbb{R}$

Start with some x_0 , set $x_{t+1} \leftarrow f(x_t)$

Suppose f is contraction, i.e., $\forall x, x', |f(x') - f(x)| \leq \gamma |x' - x|$, $\gamma \in [0, 1)$,
then $x_t \rightarrow x^\star$

For Policy Evaluation (i.e., given \mathcal{M} and π , compute V^π)

$$V^\pi = R + \underbrace{\gamma P V^\pi}_{:= \mathcal{T}^\pi V^\pi}$$

Recap: Fixed-point solution

Find the fixed point solution of $x^\star = f(x^\star)$, $x \in \mathbb{R}$

Start with some x_0 , set $x_{t+1} \leftarrow f(x_t)$

Suppose f is contraction, i.e., $\forall x, x', |f(x') - f(x)| \leq \gamma |x' - x|$, $\gamma \in [0, 1)$,
then $x_t \rightarrow x^\star$

For Policy Evaluation (i.e., given \mathcal{M} and π , compute V^π)

$$V^\pi = R + \underbrace{\gamma P V^\pi}_{:= \mathcal{T}^\pi V^\pi} \quad V^{t+1} \leftarrow \mathcal{T}^\pi V^t$$

Question for Today:

Given an MDP $\mathcal{M} = (S, A, P, r, \gamma)$, How to find $\pi^\star : S \mapsto A$ (approximately)

Motivation for Finding the Optimal Policy

Motivation for Finding the Optimal Policy

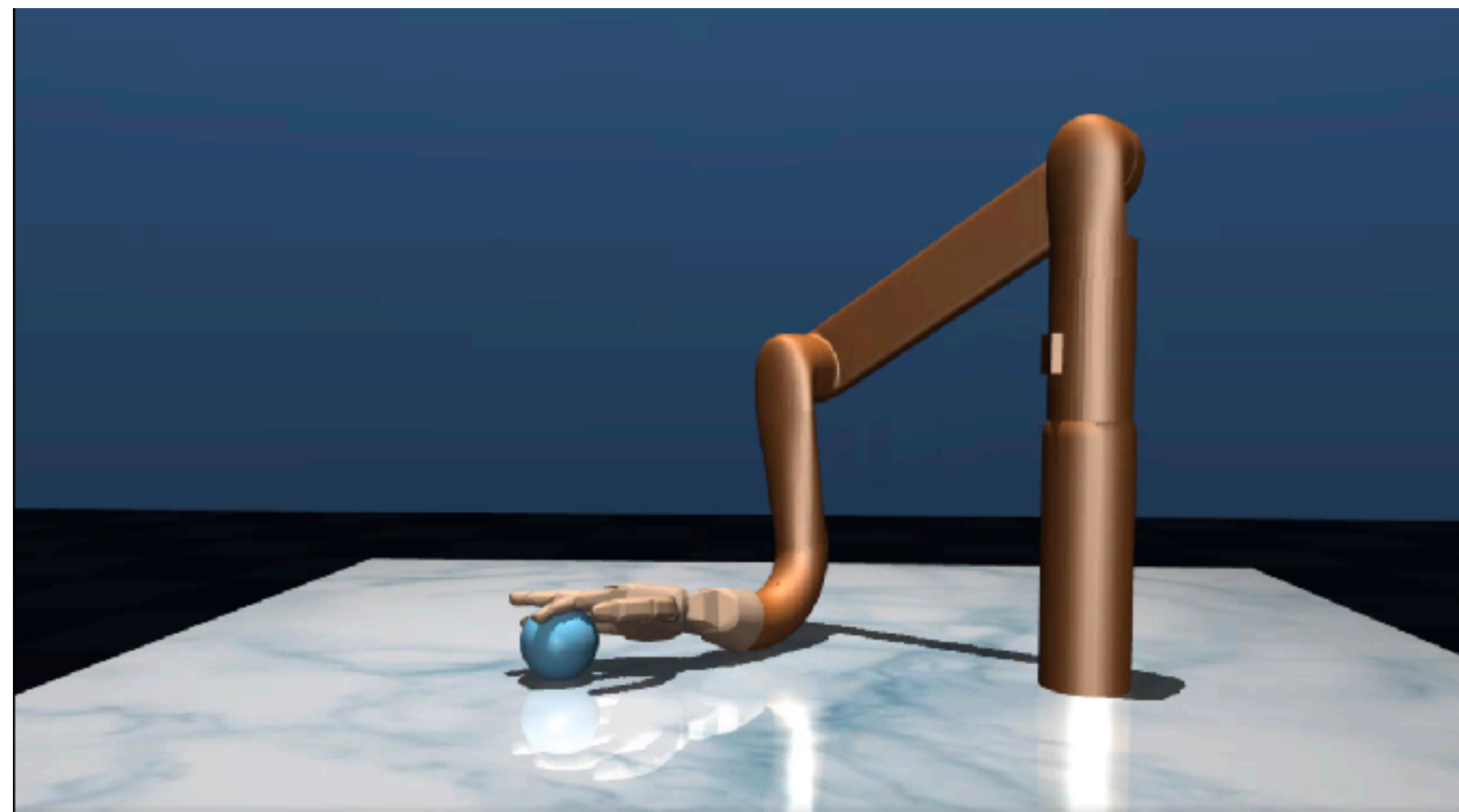


Find the strategy w/ the highest
prob of winning
(i.e., a policy that maps the board
position to the next move)

Motivation for Finding the Optimal Policy



Find the strategy w/ the highest prob of winning
(i.e., a policy that maps the board position to the next move)



Find the strategy (i.e., a mapping from robot & ball configuration to torques) that picks the ball and moves it to a goal position ASAP

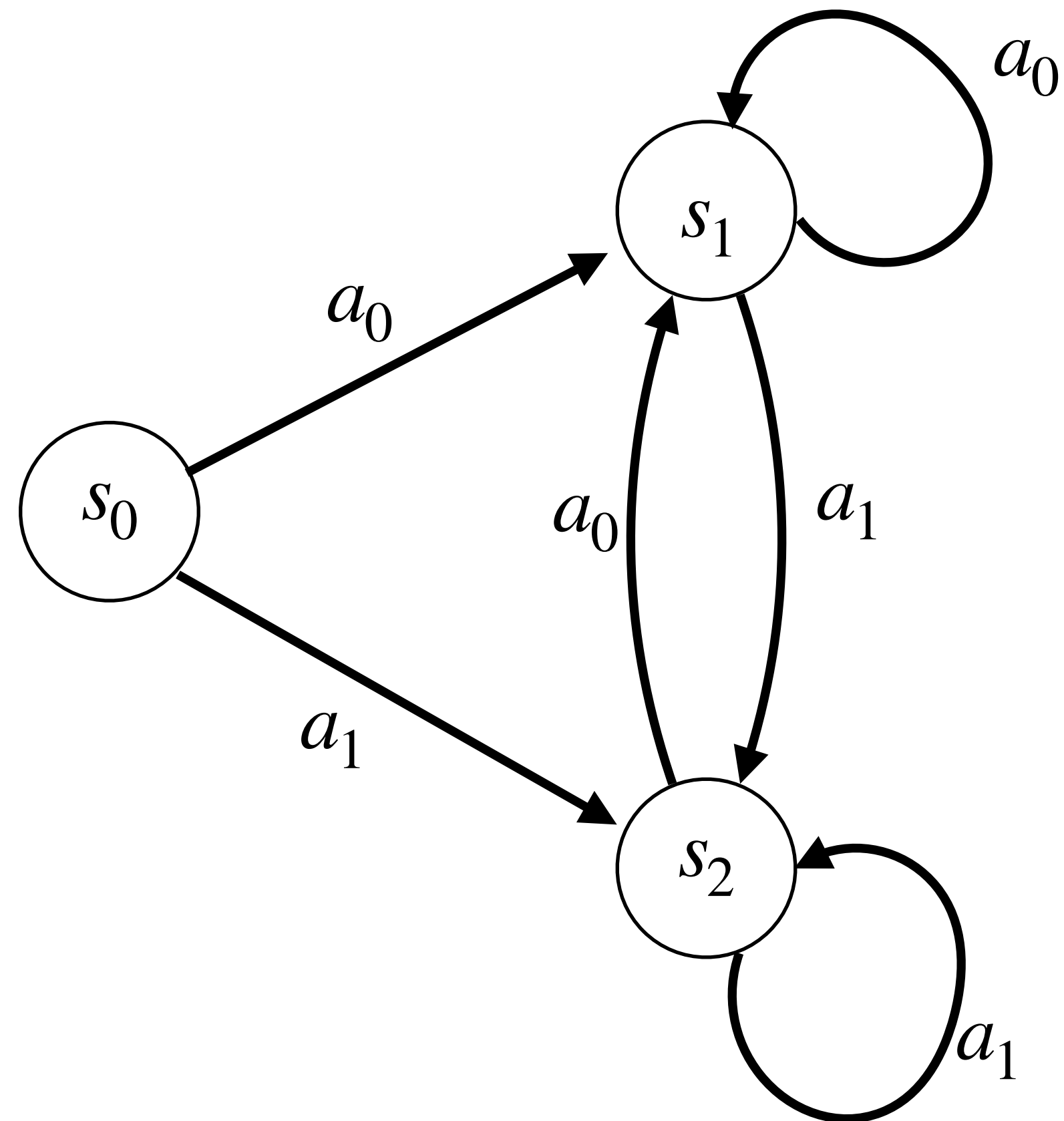
Outline:

1: An Iterative Algorithm: Value Iteration
(a fix-point iteration algorithm again!)

2: Convergence? How fast?
(Via the contraction argument again!)

Example of Optimal Policy π^\star

Consider the following **deterministic** MDP w/ 3 states & 2 actions

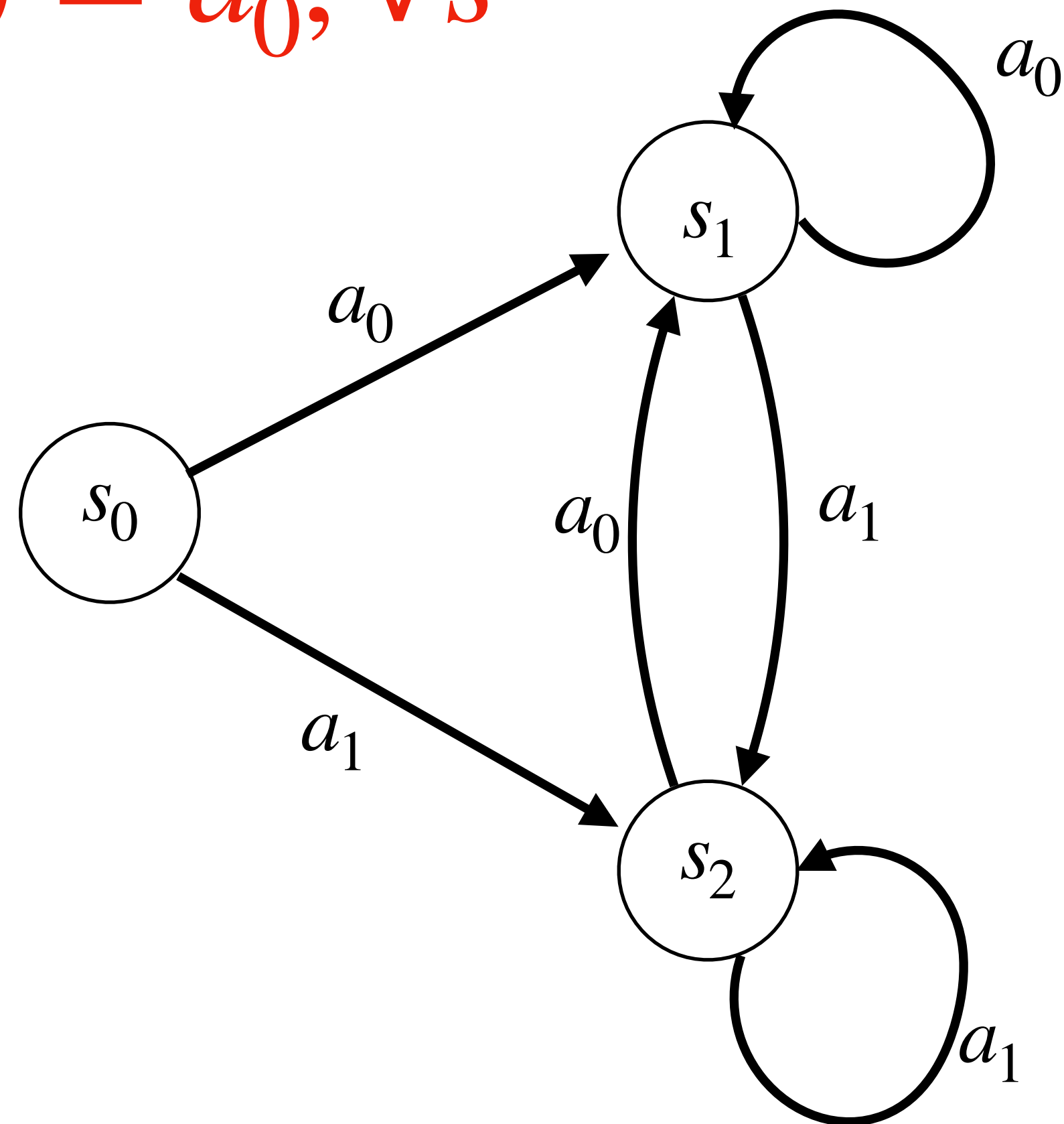


Reward: $r(s_1, a_0) = 1$, 0 everywhere else

Example of Optimal Policy π^\star

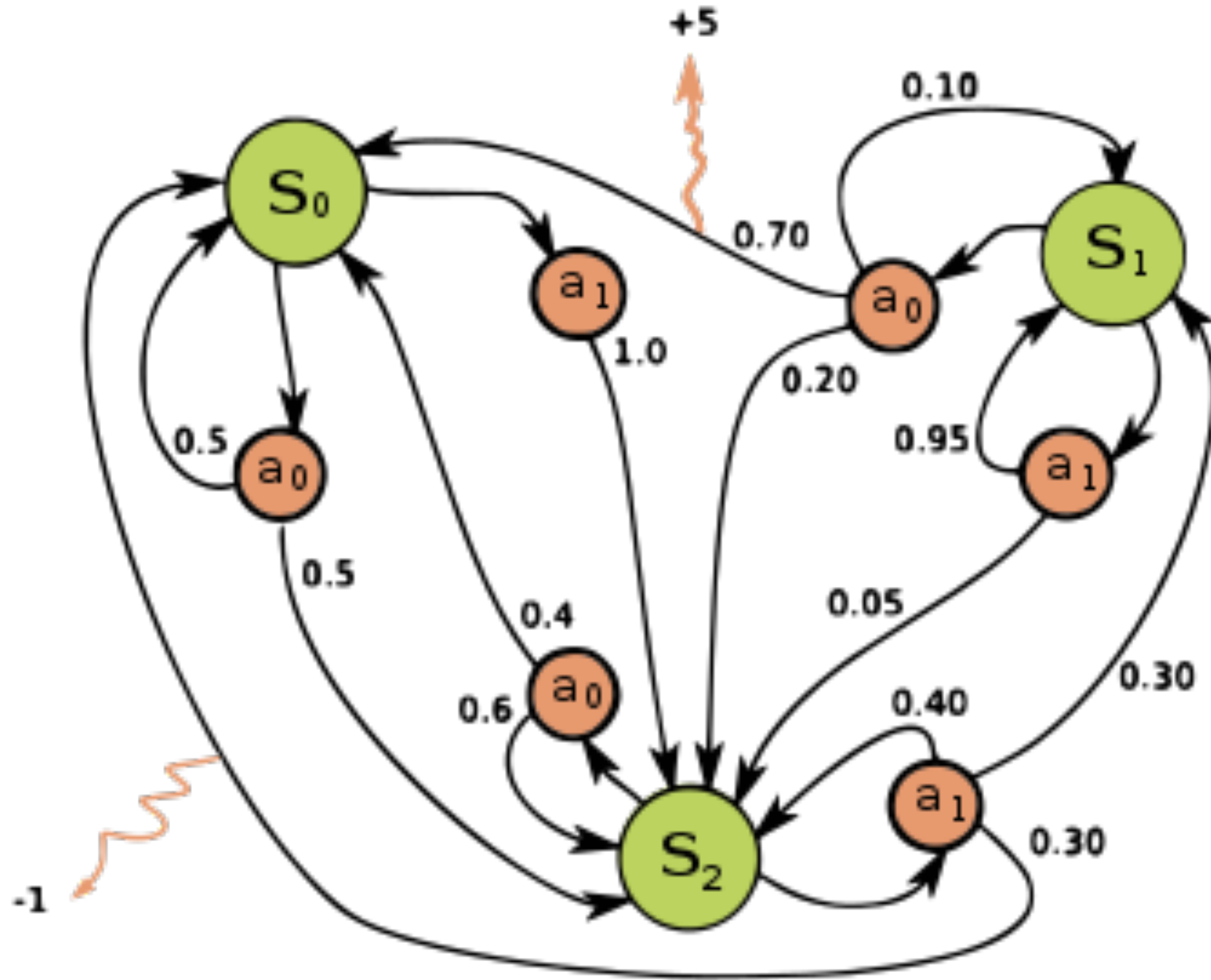
Consider the following **deterministic** MDP w/ 3 states & 2 actions

$$\pi^\star(s) = a_0, \forall s$$

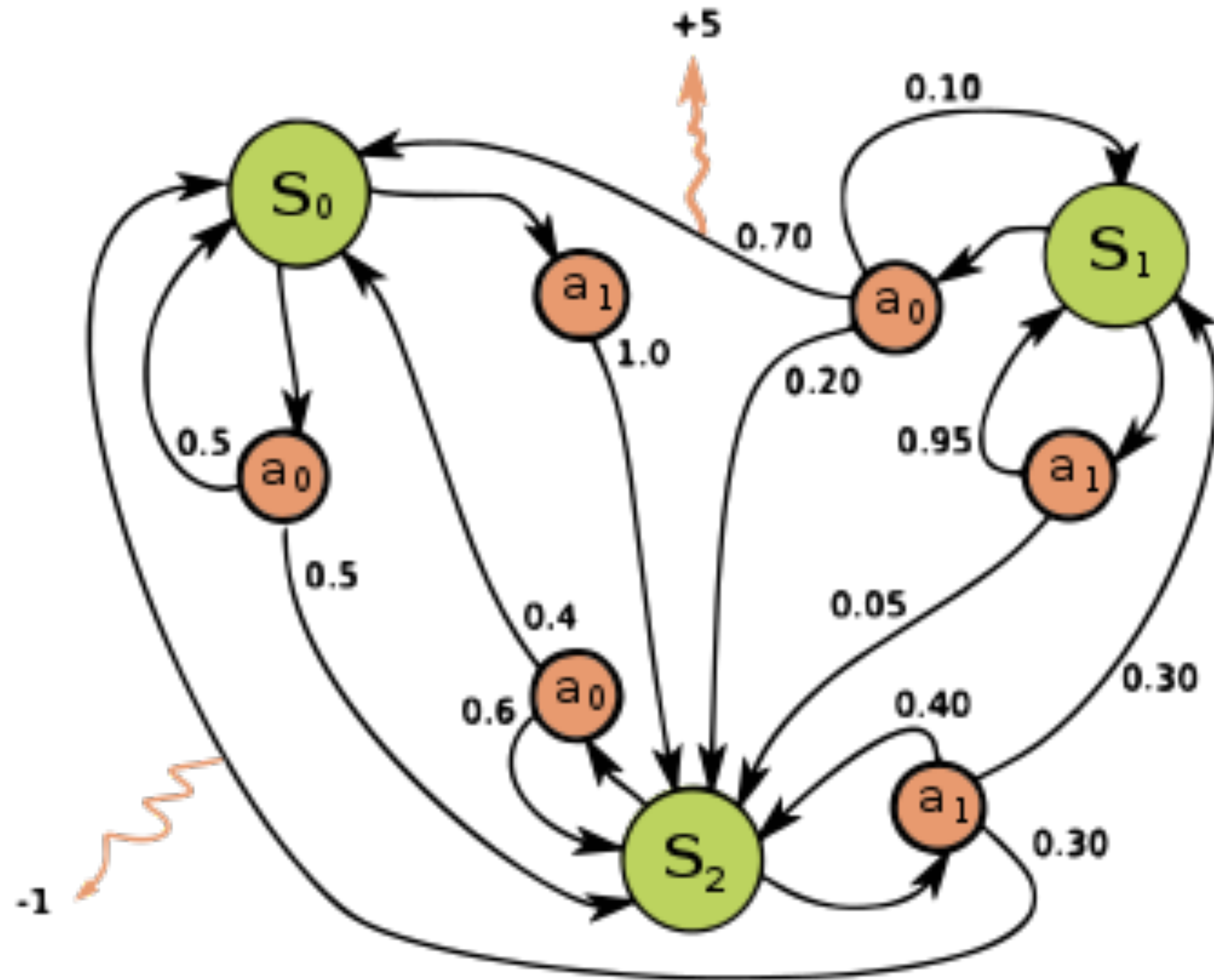


Reward: $r(s_1, a_0) = 1$, 0 everywhere else

What about this one...



What about this one...



Let's design an algorithm that computes V^*/Q^* for any given $r \in \mathbb{R}^{|S| \times |A|}$ & $P \in \mathbb{R}^{|S| \times (|S| |A|)}$

A Naive Approach (not computationally efficient)

Well, we know how to do policy evaluation for any given $\pi : S \mapsto A$, so...

A Naive Approach (not computationally efficient)

Well, we know how to do policy evaluation for any given $\pi : S \mapsto A$, so...

Enumeration:

$\forall \pi \in S \mapsto A$, do PE, i.e., $V^\pi = \text{Exact-PE}(\pi)$,

then pick the policy π' , such that:

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

A Naive Approach (not computationally efficient)

Well, we know how to do policy evaluation for any given $\pi : S \mapsto A$, so...

Enumeration:

$\forall \pi \in S \mapsto A$, do PE, i.e., $V^\pi = \text{Exact-PE}(\pi)$,

then pick the policy π' , such that:

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

Computation time: $O(A^S \cdot S^3)$

A Naive Approach (not computationally efficient)

Well, we know how to do policy evaluation for any given $\pi : S \mapsto A$, so...

Enumeration:

$\forall \pi \in S \mapsto A$, do PE, i.e., $V^\pi = \text{Exact-PE}(\pi)$,

then pick the policy π' , such that:

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

Computation time: $O(A^S \cdot S^3)$

Can we do better? We definitely want to avoid A^S ...

Define Bellman Operator \mathcal{T} :

Given a function $Q : S \times A \mapsto \mathbb{R}$,

$$\mathcal{T}Q : S \times A \mapsto \mathbb{R},$$

$$(\mathcal{T}Q)(s, a) := r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a' \in A} Q(s', a'), \forall s, a \in S \times A$$

Define Bellman Operator \mathcal{T} :

Given a function $Q : S \times A \mapsto \mathbb{R}$,

$$\mathcal{T}Q : S \times A \mapsto \mathbb{R},$$

$$(\mathcal{T}Q)(s, a) := r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a' \in A} Q(s', a'), \forall s, a \in S \times A$$

We can express $Q \in \mathbb{R}^{|S||A|}$, so $\mathcal{T}Q \in \mathbb{R}^{|S||A|}$

Define Bellman Operator \mathcal{T} :

Given a function $Q : S \times A \mapsto \mathbb{R}$,

$$\mathcal{T}Q : S \times A \mapsto \mathbb{R},$$

$$(\mathcal{T}Q)(s, a) := r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a' \in A} Q(s', a'), \forall s, a \in S \times A$$

We can express $Q \in \mathbb{R}^{|S||A|}$, so $\mathcal{T}Q \in \mathbb{R}^{|S||A|}$

i.e., think about \mathcal{T} as a (non-linear) mapping that maps from $\mathbb{R}^{|S||A|}$ to $\mathbb{R}^{|S||A|}$

High Level idea for Algorithm Design

Fix-point iteration again!

High Level idea for Algorithm Design

Fix-point iteration again!

Recall Bellman Optimality for Q^* :

$$Q^*(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a'} Q^*(s', a')$$

High Level idea for Algorithm Design

Fix-point iteration again!

Recall Bellman Optimality for Q^* :

$$Q^*(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a'} Q^*(s', a')$$

We have $Q^* = \mathcal{T} Q^*$,
i.e., Q^* is a fix-point solution of $Q = \mathcal{T} Q$

Value Iteration Algorithm:

1. Initialization: $Q^0 : \|Q^0\|_\infty \in \left[0, \frac{1}{1-\gamma}\right]$
2. Iterate until convergence: $Q^{t+1} \Leftarrow \mathcal{T} Q^t$

Value Iteration Algorithm:

1. Initialization: $Q^0 : \|Q^0\|_\infty \in \left[0, \frac{1}{1-\gamma}\right]$
2. Iterate until convergence: $Q^{t+1} \Leftarrow \mathcal{T} Q^t$

Guarantee of VI:

The fix-point iteration converges, i.e., $Q^t \rightarrow Q^*$, as $t \rightarrow \infty$

Summary so far:

Zooming in $Q^{t+1} \Leftarrow \mathcal{T} Q^t$:

Summary so far:

Zooming in $Q^{t+1} \Leftarrow \mathcal{T} Q^t$:

Given Q^t , we set:

$$\forall s, a : Q^{t+1}(s, a) \Leftarrow r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a'} Q^t(s', a')$$

Outline:

✓ 1: An Iterative Algorithm: Value Iteration
(a fix-point iteration algorithm again!)

2: Convergence? How fast?
(Via the contraction argument again!)

Convergence of Value Iteration:

Lemma [contraction]: Given any Q, Q' , we have:

$$\|\mathcal{T}Q - \mathcal{T}Q'\|_{\infty} \leq \gamma \|Q - Q'\|_{\infty}$$

Proof:

Convergence of Value Iteration:

Lemma [contraction]: Given any Q, Q' , we have:

$$\|\mathcal{T}Q - \mathcal{T}Q'\|_\infty \leq \gamma \|Q - Q'\|_\infty$$

Proof:

$$|(\mathcal{T}Q)(s, a) - (\mathcal{T}Q')(s, a)| = \left| r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} \max_{a'} Q(s', a') - \left(r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} \max_{a'} Q'(s', a') \right) \right|$$

Convergence of Value Iteration:

Lemma [contraction]: Given any Q, Q' , we have:

$$\|\mathcal{T}Q - \mathcal{T}Q'\|_\infty \leq \gamma \|Q - Q'\|_\infty$$

Proof:

$$\begin{aligned} |(\mathcal{T}Q)(s, a) - (\mathcal{T}Q')(s, a)| &= \left| r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} \max_{a'} Q(s', a') - \left(r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} \max_{a'} Q'(s', a') \right) \right| \\ &\leq \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left| \left(\max_{a'} Q(s', a') - \max_{a'} Q'(s', a') \right) \right| \end{aligned}$$

Convergence of Value Iteration:

Lemma [contraction]: Given any Q, Q' , we have:

$$\|\mathcal{T}Q - \mathcal{T}Q'\|_\infty \leq \gamma \|Q - Q'\|_\infty$$

Proof:

$$\begin{aligned} |(\mathcal{T}Q)(s, a) - (\mathcal{T}Q')(s, a)| &= \left| r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} \max_{a'} Q(s', a') - \left(r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} \max_{a'} Q'(s', a') \right) \right| \\ &\leq \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left| \left(\max_{a'} Q(s', a') - \max_{a'} Q'(s', a') \right) \right| \\ &\leq \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a'} \left| (Q(s', a') - Q'(s', a')) \right| \end{aligned}$$

Convergence of Value Iteration:

Lemma [contraction]: Given any Q, Q' , we have:

$$\|\mathcal{T}Q - \mathcal{T}Q'\|_\infty \leq \gamma \|Q - Q'\|_\infty$$

Proof:

$$\begin{aligned} |(\mathcal{T}Q)(s, a) - (\mathcal{T}Q')(s, a)| &= \left| r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} \max_{a'} Q(s', a') - \left(r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} \max_{a'} Q'(s', a') \right) \right| \\ &\leq \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left| \left(\max_{a'} Q(s', a') - \max_{a'} Q'(s', a') \right) \right| \\ &\leq \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a'} \left| (Q(s', a') - Q'(s', a')) \right| \\ &\leq \gamma \max_{s'} \max_{a'} \left| (Q(s', a') - Q'(s', a')) \right| \end{aligned}$$

Convergence of Value Iteration:

Lemma [contraction]: Given any Q, Q' , we have:

$$\|\mathcal{T}Q - \mathcal{T}Q'\|_\infty \leq \gamma \|Q - Q'\|_\infty$$

Proof:

$$\begin{aligned} |(\mathcal{T}Q)(s, a) - (\mathcal{T}Q')(s, a)| &= \left| r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} \max_{a'} Q(s', a') - \left(r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} \max_{a'} Q'(s', a') \right) \right| \\ &\leq \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left| \left(\max_{a'} Q(s', a') - \max_{a'} Q'(s', a') \right) \right| \\ &\leq \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a'} \left| (Q(s', a') - Q'(s', a')) \right| \\ &\leq \gamma \max_{s'} \max_{a'} \left| (Q(s', a') - Q'(s', a')) \right| = \gamma \|Q - Q'\|_\infty \end{aligned}$$

Convergence of Value Iteration:

Lemma [Convergence]: Given Q^0 , we have:

$$\|Q^t - Q^*\|_\infty \leq \gamma^t \|Q^0 - Q^*\|_\infty$$

Proof ??

Convergence of Value Iteration:

Lemma [Convergence]: Given Q^0 , we have:

$$\|Q^t - Q^*\|_\infty \leq \gamma^t \|Q^0 - Q^*\|_\infty$$

Proof ??

$$\|Q^{t+1} - Q^*\|_\infty = \|\mathcal{T}Q^t - \mathcal{T}Q^*\|_\infty \leq \gamma \|Q^t - Q^*\|_\infty$$

Convergence of Value Iteration:

Lemma [Convergence]: Given Q^0 , we have:

$$\|Q^t - Q^*\|_\infty \leq \gamma^t \|Q^0 - Q^*\|_\infty$$

Proof ??

$$\|Q^{t+1} - Q^*\|_\infty = \|\mathcal{T}Q^t - \mathcal{T}Q^*\|_\infty \leq \gamma \|Q^t - Q^*\|_\infty$$

$$\dots \leq \gamma^{t+1} \|\widehat{Q}^0 - Q^*\|_\infty$$

Summary so far:

VI (a fix point iteration alg):

$$Q^{t+1} \Leftarrow \mathcal{T} Q^t$$

VI convergence (via contraction)

$$\text{i.e., } \|Q^t - Q^*\|_\infty \leq \gamma^t \|Q^0 - Q^*\|_\infty$$

Summary so far:

VI (a fix point iteration alg):

$$Q^{t+1} \leftarrow \mathcal{T} Q^t$$

VI convergence (via contraction)

$$\text{i.e., } \|Q^t - Q^*\|_\infty \leq \gamma^t \|Q^0 - Q^*\|_\infty$$

Next: what about the policy? Ultimately, we do want π^* ...



From Q functions to policies...

We know that $\pi^\star(s) = \arg \max_a Q^\star(s, a)$

Recall that VI ensures that $Q^t(s, a) \approx Q^\star(s, a), \forall s, a, \dots$

From Q functions to policies...

We know that $\pi^\star(s) = \arg \max_a Q^\star(s, a)$

Recall that VI ensures that $Q^t(s, a) \approx Q^\star(s, a), \forall s, a, \dots$

then maybe $\pi(s) := \arg \max_a Q^t(s, a)$ is a good choice?

The Quality of Policy:

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

Theorem: $V^{\pi^t}(s) \geq V^*(s) - \frac{2\gamma^t}{1-\gamma} \|Q^0 - Q^*\|_\infty \forall s \in S$

The Quality of Policy:

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

Theorem: $V^{\pi^t}(s) \geq V^*(s) - \frac{2\gamma^t}{1-\gamma} \|Q^0 - Q^*\|_\infty \forall s \in S$

Proof:

The Quality of Policy:

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

$$\textbf{Theorem: } V^{\pi^t}(s) \geq V^*(s) - \frac{2\gamma^t}{1-\gamma} \|Q^0 - Q^*\|_\infty \forall s \in S$$

Proof:

$$V^{\pi^t}(s) - V^*(s) = Q^{\pi^t}(s, \pi^t(s)) - Q^*(s, \pi^*(s))$$

The Quality of Policy:

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

$$\textbf{Theorem: } V^{\pi^t}(s) \geq V^*(s) - \frac{2\gamma^t}{1-\gamma} \|Q^0 - Q^*\|_\infty \forall s \in S$$

Proof:

$$\begin{aligned} V^{\pi^t}(s) - V^*(s) &= Q^{\pi^t}(s, \pi^t(s)) - Q^*(s, \pi^*(s)) \\ &= Q^{\pi^t}(s, \pi^t(s)) - Q^*(s, \pi^t(s)) + Q^*(s, \pi^t(s)) - Q^*(s, \pi^*(s)) \end{aligned}$$

The Quality of Policy:

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

$$\textbf{Theorem: } V^{\pi^t}(s) \geq V^*(s) - \frac{2\gamma^t}{1-\gamma} \|Q^0 - Q^*\|_\infty \forall s \in S$$

Proof:

$$\begin{aligned} V^{\pi^t}(s) - V^*(s) &= Q^{\pi^t}(s, \pi^t(s)) - Q^*(s, \pi^*(s)) \\ &= Q^{\pi^t}(s, \pi^t(s)) - Q^*(s, \pi^t(s)) + Q^*(s, \pi^t(s)) - Q^*(s, \pi^*(s)) \\ &= \gamma \mathbb{E}_{s' \sim P(s, \pi^t(s))} \left(V^{\pi^t}(s') - V^*(s') \right) + Q^*(s, \pi^t(s)) - Q^*(s, \pi^*(s)) \end{aligned}$$

The Quality of Policy:

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

$$\textbf{Theorem: } V^{\pi^t}(s) \geq V^*(s) - \frac{2\gamma^t}{1-\gamma} \|Q^0 - Q^*\|_\infty \forall s \in S$$

Proof:

$$\begin{aligned} V^{\pi^t}(s) - V^*(s) &= Q^{\pi^t}(s, \pi^t(s)) - Q^*(s, \pi^*(s)) \\ &= Q^{\pi^t}(s, \pi^t(s)) - Q^*(s, \pi^t(s)) + Q^*(s, \pi^t(s)) - Q^*(s, \pi^*(s)) \\ &= \gamma \mathbb{E}_{s' \sim P(s, \pi^t(s))} \left(V^{\pi^t}(s') - V^*(s') \right) + Q^*(s, \pi^t(s)) - Q^*(s, \pi^*(s)) \\ &\geq \gamma \mathbb{E}_{s' \sim P(s, \pi^t(s))} \left(V^{\pi^t}(s') - V^*(s') \right) + Q^*(s, \pi^t(s)) - Q^t(s, \pi^t(s)) + Q^t(s, \pi^*(s)) - Q^*(s, \pi^*(s)) \end{aligned}$$

The Quality of Policy:

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

$$\textbf{Theorem: } V^{\pi^t}(s) \geq V^*(s) - \frac{2\gamma^t}{1-\gamma} \|Q^0 - Q^*\|_\infty \forall s \in S$$

Proof:

$$\begin{aligned} V^{\pi^t}(s) - V^*(s) &= Q^{\pi^t}(s, \pi^t(s)) - Q^*(s, \pi^*(s)) \\ &= Q^{\pi^t}(s, \pi^t(s)) - Q^*(s, \pi^t(s)) + Q^*(s, \pi^t(s)) - Q^*(s, \pi^*(s)) \\ &= \gamma \mathbb{E}_{s' \sim P(s, \pi^t(s))} \left(V^{\pi^t}(s') - V^*(s') \right) + Q^*(s, \pi^t(s)) - Q^*(s, \pi^*(s)) \\ &\geq \gamma \mathbb{E}_{s' \sim P(s, \pi^t(s))} \left(V^{\pi^t}(s') - V^*(s') \right) + Q^*(s, \pi^t(s)) - Q^t(s, \pi^t(s)) + Q^t(s, \pi^*(s)) - Q^*(s, \pi^*(s)) \\ &\geq \gamma \mathbb{E}_{s' \sim P(s, \pi^t(s))} \left(V^{\pi^t}(s') - V^*(s') \right) - 2\gamma^t \|Q^0 - Q^*\|_\infty \end{aligned}$$

The Quality of Policy:

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

$$\textbf{Theorem: } V^{\pi^t}(s) \geq V^*(s) - \frac{2\gamma^t}{1-\gamma} \|Q^0 - Q^*\|_\infty \forall s \in S$$

Proof:

$$\begin{aligned} V^{\pi^t}(s) - V^*(s) &= Q^{\pi^t}(s, \pi^t(s)) - Q^*(s, \pi^*(s)) \\ &= Q^{\pi^t}(s, \pi^t(s)) - Q^*(s, \pi^t(s)) + Q^*(s, \pi^t(s)) - Q^*(s, \pi^*(s)) \\ &= \gamma \mathbb{E}_{s' \sim P(s, \pi^t(s))} \left(V^{\pi^t}(s') - V^*(s') \right) + Q^*(s, \pi^t(s)) - Q^*(s, \pi^*(s)) \\ &\geq \gamma \mathbb{E}_{s' \sim P(s, \pi^t(s))} \left(V^{\pi^t}(s') - V^*(s') \right) + Q^*(s, \pi^t(s)) - Q^t(s, \pi^t(s)) + Q^t(s, \pi^*(s)) - Q^*(s, \pi^*(s)) \\ &\geq \gamma \mathbb{E}_{s' \sim P(s, \pi^t(s))} \left(V^{\pi^t}(s') - V^*(s') \right) - 2\gamma^t \|Q^0 - Q^*\|_\infty \dots \text{Recursion} \end{aligned}$$

Summary for VI:

1. VI

(a fix point iteration alg):

$$Q^{t+1} \leftarrow \mathcal{T} Q^t$$

Summary for VI:

1. VI

(a fix point iteration alg):

$$Q^{t+1} \leftarrow \mathcal{T} Q^t$$

Contraction



2. VI convergence: exponentially fast,
i.e., $\|Q^t - Q^*\|_\infty \leq \gamma^t \|Q^0 - Q^*\|_\infty$

Summary for VI:

1. VI

(a fix point iteration alg):

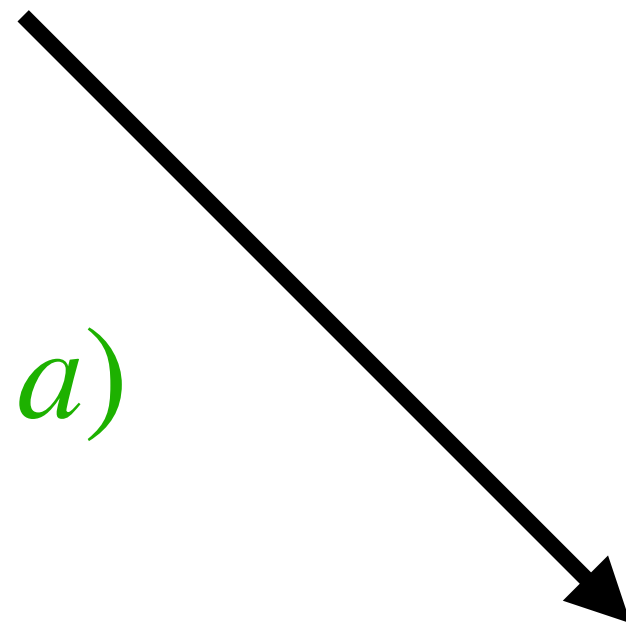
$$Q^{t+1} \leftarrow \mathcal{T} Q^t$$

Contraction



2. VI convergence: exponentially fast,
i.e., $\|Q^t - Q^*\|_\infty \leq \gamma^t \|Q^0 - Q^*\|_\infty$

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



Summary for VI:

1. VI

Contraction

(a fix point iteration alg):

$$Q^{t+1} \leftarrow \mathcal{T} Q^t$$

2. VI convergence: exponentially fast,
i.e., $\|Q^t - Q^*\|_\infty \leq \gamma^t \|Q^0 - Q^*\|_\infty$

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

3. Policy Performance: $V^{\pi^t}(s) \geq V^*(s) - \frac{2\gamma^t}{1-\gamma} \|Q^0 - Q^*\|_\infty \forall s \in S$

Summary for this week

Bellman Equation:

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

Bellman Optimality
(Q-version):

$$Q^\star = \mathcal{T} Q^\star$$

Summary for this week

Bellman Equation:

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

Fix-point Iteration
framework

Bellman Optimality
(Q-version):

$$Q^\star = \mathcal{T} Q^\star$$

Summary for this week

Bellman Equation:

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

Fix-point Iteration
framework

Bellman Optimality
(Q-version):

$$Q^\star = \mathcal{T} Q^\star$$

Iterative PE:

$$V^{t+1} \leftarrow R + P V^t$$

Summary for this week

Bellman Equation:

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

Fix-point Iteration
framework

Bellman Optimality
(Q-version):

$$Q^\star = \mathcal{T} Q^\star$$

Iterative PE:

$$V^{t+1} \Leftarrow R + P V^t$$

VI:

$$Q^{t+1} \Leftarrow \mathcal{T} Q^t$$

Summary for this week

Contraction

Bellman Equation:

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

Fix-point Iteration
framework

Bellman Optimality
(Q-version):

$$Q^\star = \mathcal{T} Q^\star$$

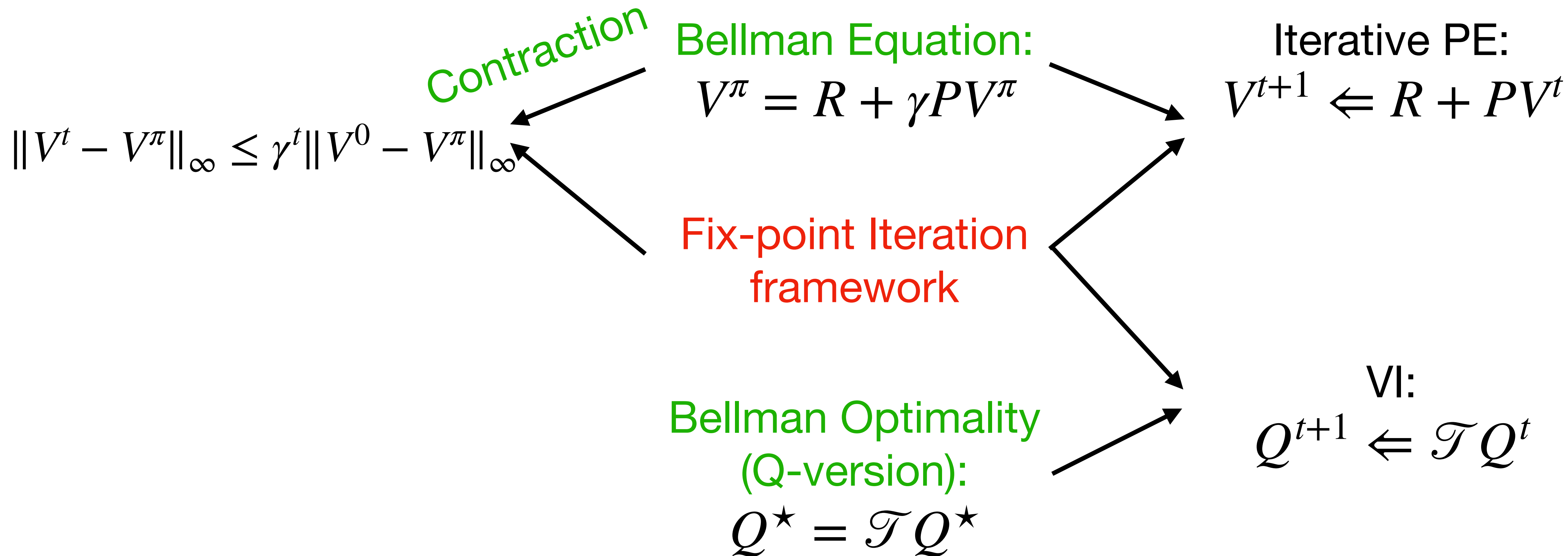
Iterative PE:

$$V^{t+1} \leftarrow R + P V^t$$

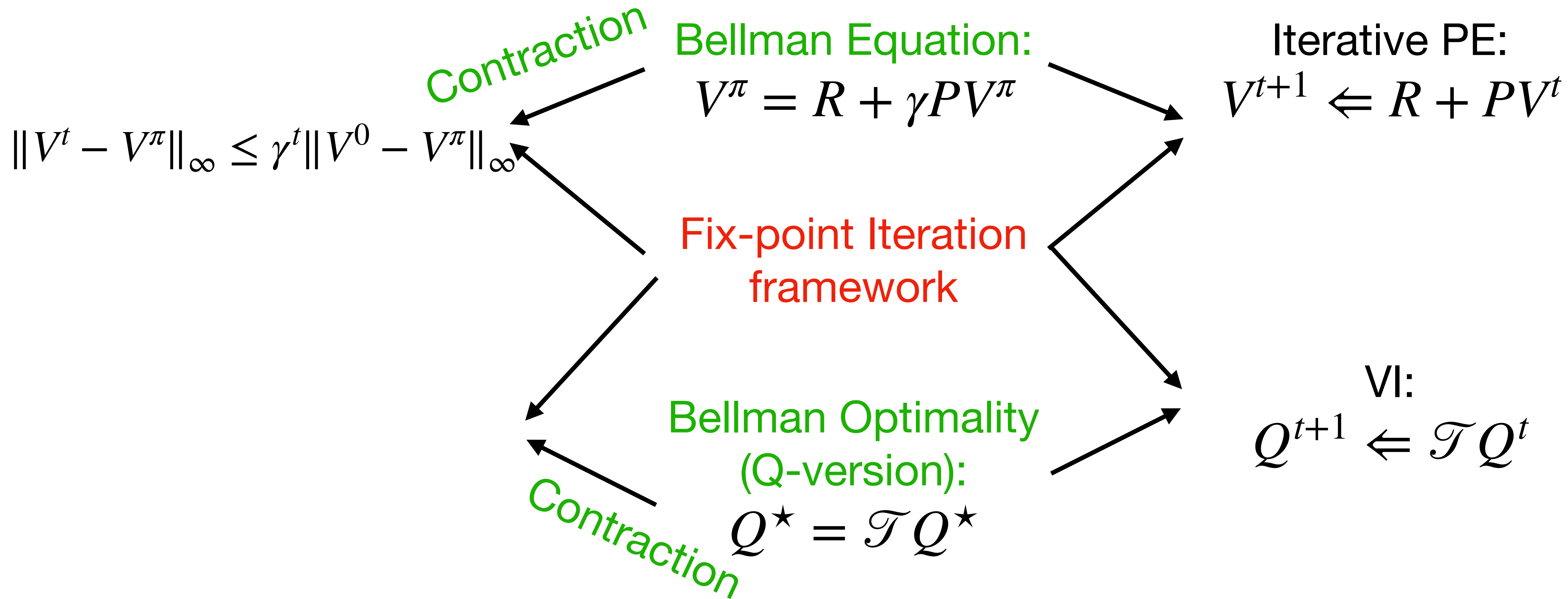
VI:

$$Q^{t+1} \leftarrow \mathcal{T} Q^t$$

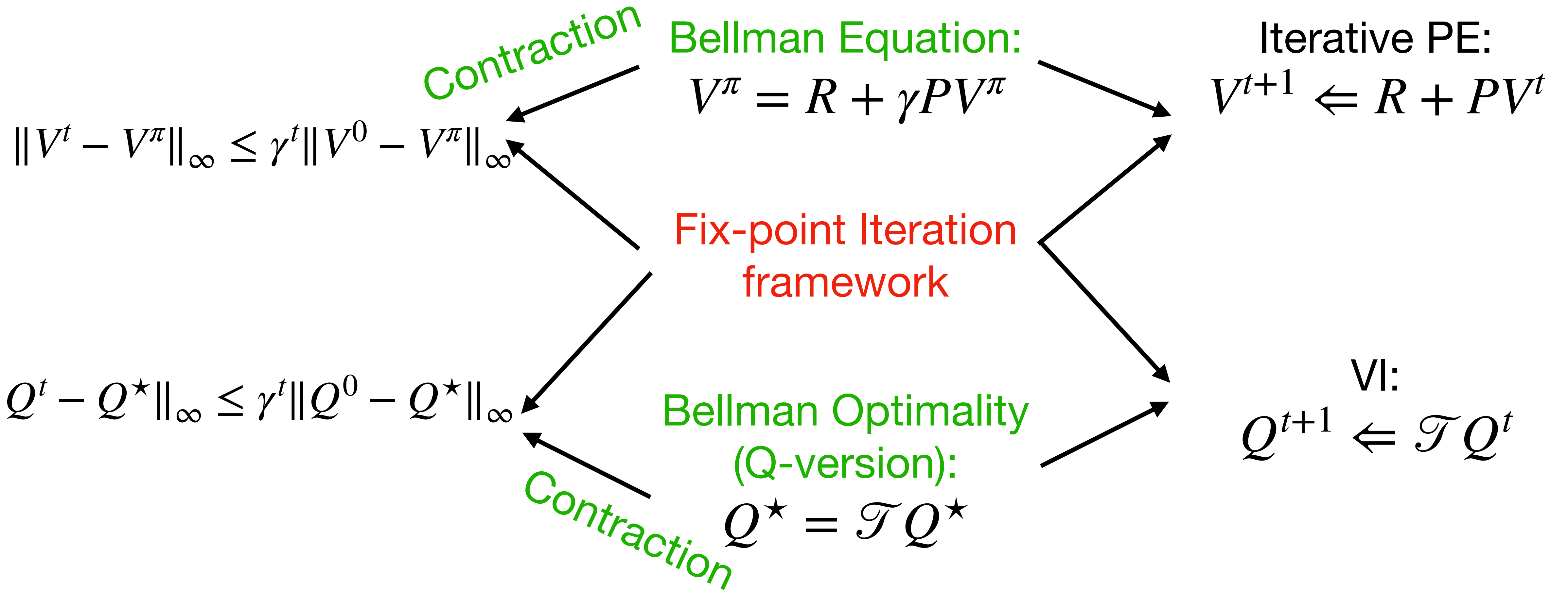
Summary for this week



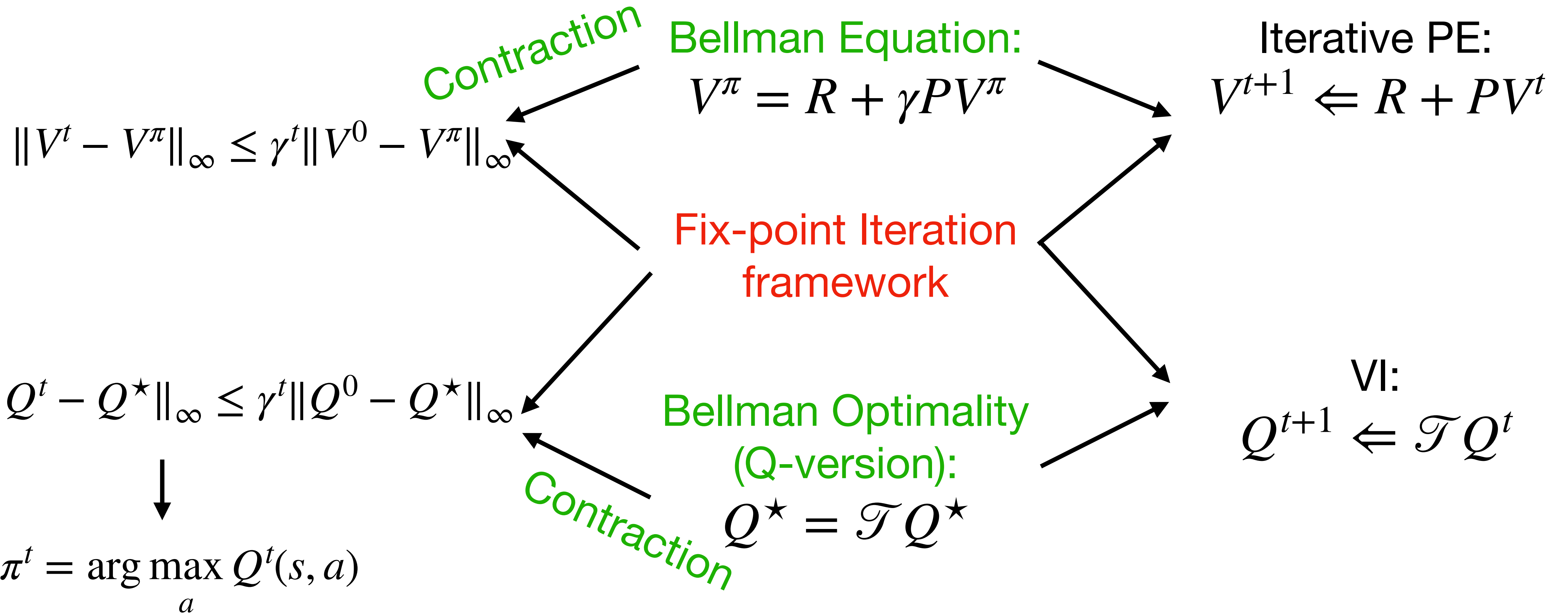
Summary for this week



Summary for this week



Summary for this week



Next week:

1. One more algorithm (Policy Iteration) for computing π^\star
2. A continuous control model: Linear Quadratic Regulator