

Value Iteration

Recap: Optimal Policy π^\star

rc1

For infinite horizon discounted MDP, there always exists a deterministic policy

$$\pi^\star : S \mapsto A, \text{ s.t.}, V^{\pi^\star}(s) \geq V^\pi(s), \forall s, \pi$$

[Puterman 94 chapter 6, also see theorem 1.7 in the RL monograph—no need to understand the proof]

Recap: Optimal Policy π^\star

For infinite horizon discounted MDP, there always exists a deterministic policy

$$\pi^\star : S \mapsto A, \text{ s.t.}, V^{\pi^\star}(s) \geq V^\pi(s), \forall s, \pi$$

[Puterman 94 chapter 6, also see theorem 1.7 in the RL monograph—no need to understand the proof]

i.e., π^\star dominates any other policy π , everywhere!

Recap: Optimal Policy π^\star

For infinite horizon discounted MDP, there always exists a deterministic policy

$$\pi^\star : S \mapsto A, \text{ s.t.}, V^{\pi^\star}(s) \geq V^\pi(s), \forall s, \pi$$

[Puterman 94 chapter 6, also see theorem 1.7 in the RL monograph—no need to understand the proof]

i.e., π^\star dominates any other policy π , everywhere!

We often denote V^\star, Q^\star in short for $V^{\pi^\star}, Q^{\pi^\star}$

Recap: Optimal Policy π^\star

For infinite horizon discounted MDP, there always exists a deterministic policy

$$\pi^\star : S \mapsto A, \text{ s.t., } V^{\pi^\star}(s) \geq V^\pi(s), \forall s, \pi$$

[Puterman 94 chapter 6, also see theorem 1.7 in the RL monograph—no need to understand the proof]

i.e., π^\star dominates any other policy π , everywhere!

We often denote V^\star, Q^\star in short for $V^{\pi^\star}, Q^{\pi^\star}$

$V^\star(s)$: the maximum value we could possibly achieve

Question for Today and Wed:

Given an MDP $\mathcal{M} = (S, A, P, r, \gamma)$, How to find Q^* and π^* (approximately)

π or V^π or Q^π

Motivation for Finding the Optimal Policy

Motivation for Finding the Optimal Policy

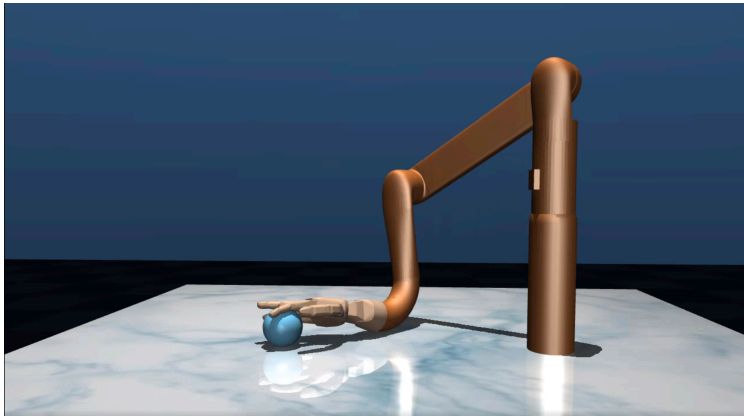


Find the strategy w/ the highest
prob of winning
(i.e., a policy that maps the board
position to the next move)

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:

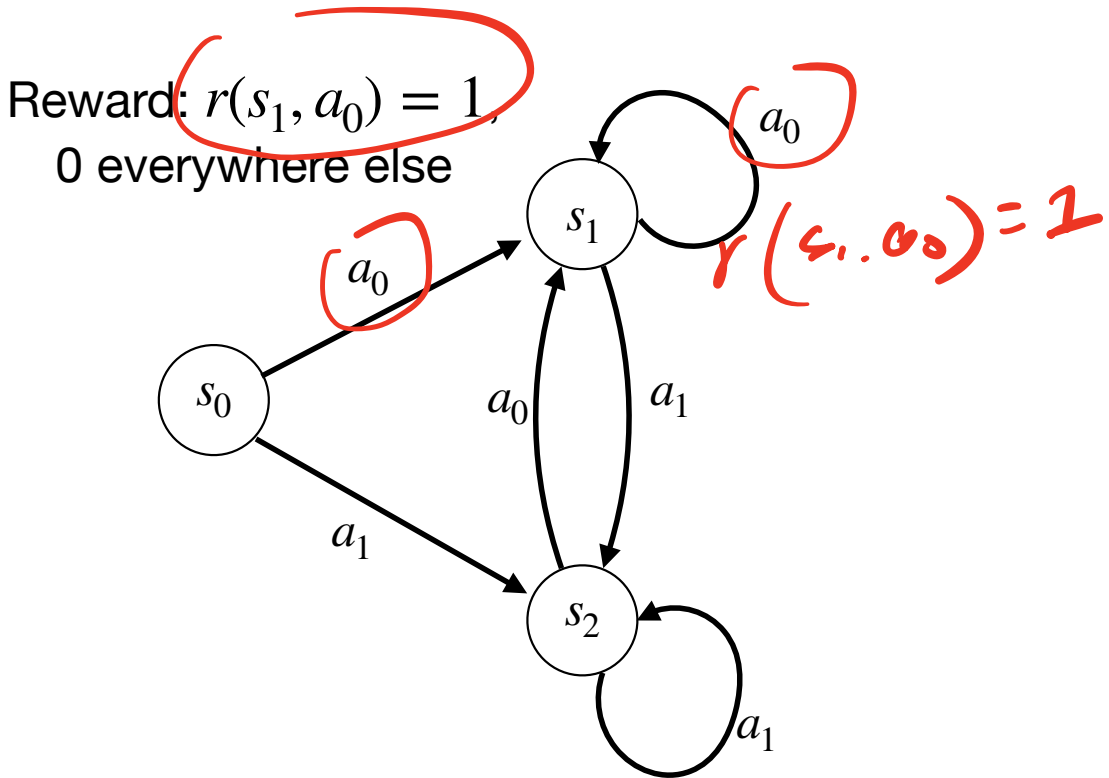
1: Bellman optimality

 Bellman
Equation

2: An Iterative Algorithm: Value Iteration

Example of Optimal Policy π^\star

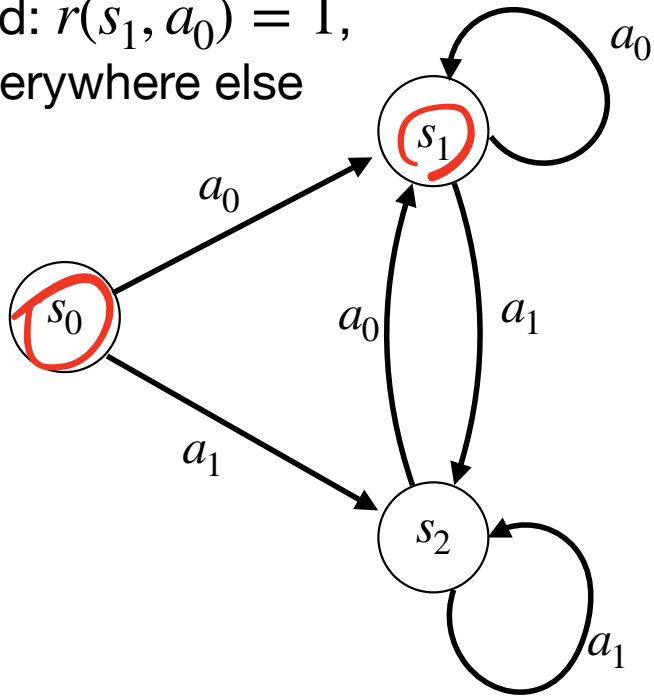
Consider the following **deterministic** MDP w/ 3 states & 2 actions



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



Let's say $\gamma \in (0,1)$
What's the optimal policy?

$$\pi^\star(s_0) = a_0$$

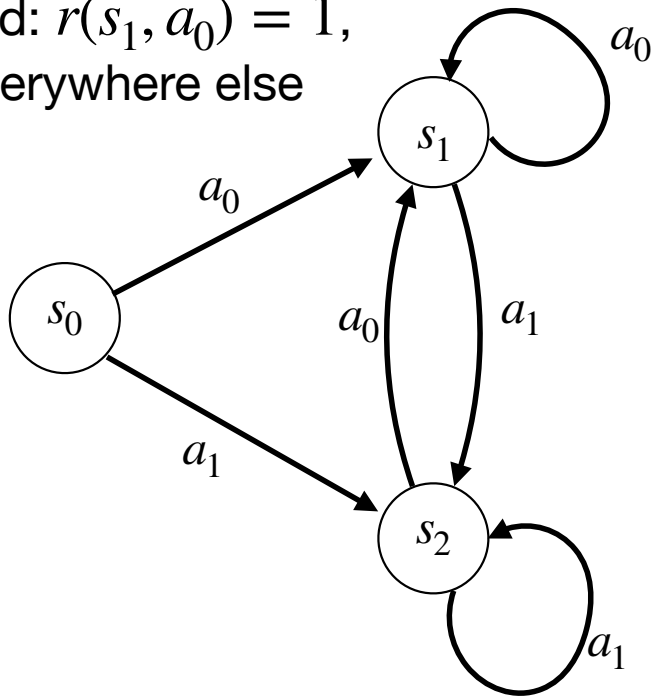
$$\pi^\star(s_1) = a_0$$

$$\pi^\star(s_2) = a_0$$

Example of Optimal Policy π^\star

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

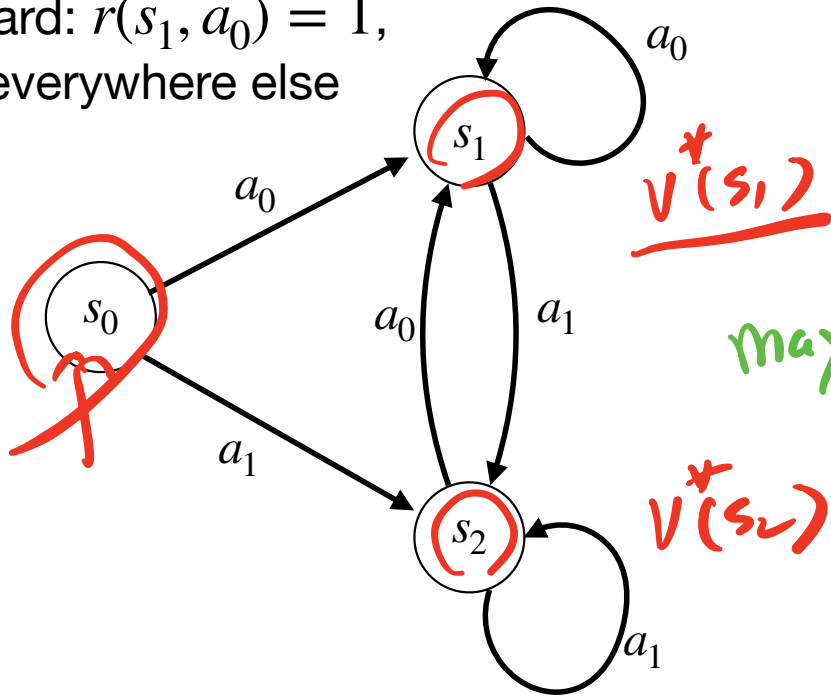
Reward: $r(s_1, a_0) = 1$,
0 everywhere else



Example of Optimal Policy π^*

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

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



If we were told $V^*(s_1)$ & $V^*(s_2)$,
how to compute $V^*(s_0)$

$$\max \left\{ \begin{array}{l} r(s_0, a_0) + \gamma \cdot \underbrace{V^*(s_1)} \\ \underbrace{=0} \end{array} \right., \left. \begin{array}{l} r(s_0, a_1) + \gamma \cdot \underbrace{V^*(s_2)} \\ \underbrace{=0} \end{array} \right\}$$

Bellman Optimality

Bellman Optimality

$$V^*(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^*(s') \right], \forall s$$

$V^*(s)$

Bellman Optimality

Bellman Optimality

Value function of π^*

$$V^*(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^*(s') \right], \forall s$$

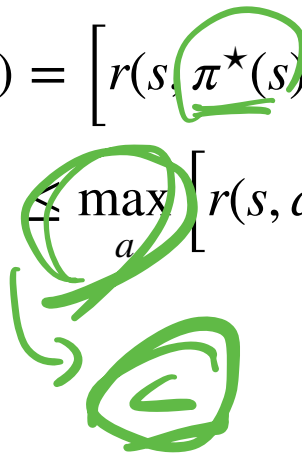
$$V^*(s) = \left[r(s, \pi^*(s)) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi^*(s))} V^*(s') \right] \quad (\text{By BE of } \pi^*)$$

Bellman Optimality

Bellman Optimality

$$V^*(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^*(s') \right], \forall s$$

$$V^*(s) = \left[r(s, \pi^*(s)) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi^*(s))} V^*(s') \right] \quad (\text{By BE of } \pi^*)$$

$$\leq \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^*(s') \right]$$


Bellman Optimality

Bellman Optimality

$$V^*(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^*(s') \right], \forall s$$

$$V^*(s) = \left[r(s, \pi^*(s)) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi^*(s))} V^*(s') \right] \quad (\text{By BE of } \pi^*)$$

$$\leq \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^*(s') \right]$$

$> \underline{V^*(s)}$

If we took this $\arg \max_a$ at s , then follow π^* , we **would have higher value**

Bellman Optimality

Bellman Optimality

$$V^*(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^*(s') \right], \forall s$$

$$V^*(s) = \left[r(s, \pi^*(s)) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi^*(s))} V^*(s') \right] \quad (\text{By BE of } \pi^*)$$

$$\leq \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^*(s') \right]$$

If we took this $\arg \max_a$ at s , then follow π^* , we **would have higher value**

Contradicts to the fact that $V^*(s)$ is the maximum value at s one could possibly achieve



Bellman Optimality

Bellman Optimality

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

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

Q^π
 Q^{π^*}
 (Q^*)

$Q^*(s, a)$ s a π^*

Bellman Optimality

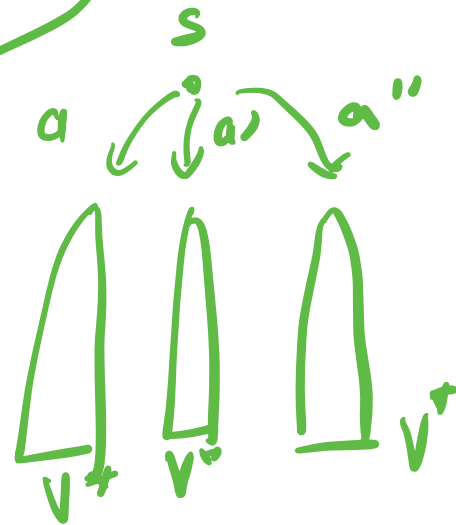
Bellman Optimality

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

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

This implies that $\arg \max_a Q^*(s, a)$ is an optimal policy

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



Bellman Optimality

Bellman Optimality

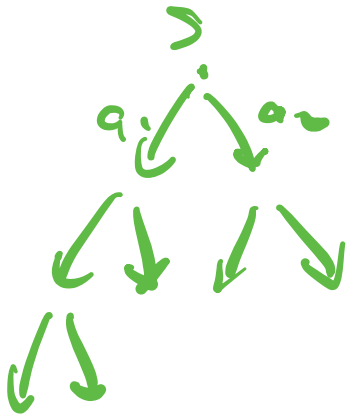
$$V^*(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^*(s') \right]$$

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

This implies that $\arg \max_a Q^*(s, a)$ is an optimal policy

a

An optimal policy should pick
this action at s



Outline:

1: Bellman optimality

2: An Iterative Algorithm: Value Iteration

$$\Rightarrow \hat{Q} \approx Q^*$$

Define the Bellman optimality for Q^*

We now know:

$$\underline{V^*(s)} = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \underline{V^*(s')} \right]$$

What's the Q version?

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

Define the Bellman optimality for Q^*

We now know:

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

What's the Q version?

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

$= V^*(s')$

Define the Bellman optimality for Q^*

We now know:

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

What's the Q version?

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

To estimate Q^* , we will use the fix-point iterative approach again

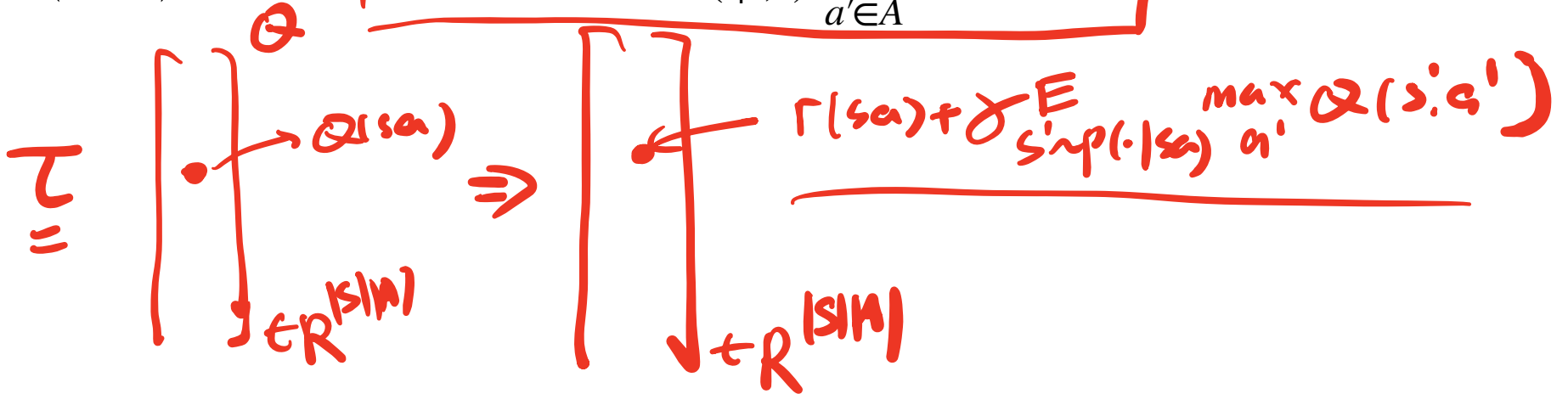
Define the Bellman operator

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

\mathcal{T}

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

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



Define the Bellman operator

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 the Bellman operator

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 $\mathbb{R}^{|S||A|}$

$$\mathcal{T}(Q + Q') \neq \mathcal{T}Q + \mathcal{T}Q'$$

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^* .

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

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

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:

$$T(Q_1 + Q_2) \\ \neq TQ_1 + TQ_2$$

1. Initialization: Q^0

2. Iterate until convergence: $Q^{t+1} \leftarrow TQ^t$

For t in $t=0$ to T :

For All $(s,a) \in S \times A$:

$$\text{set } Q^{t+1}(s,a) \leftarrow r(s,a) + \gamma \max_{a'} Q^t(s',a')$$

Value Iteration Algorithm:

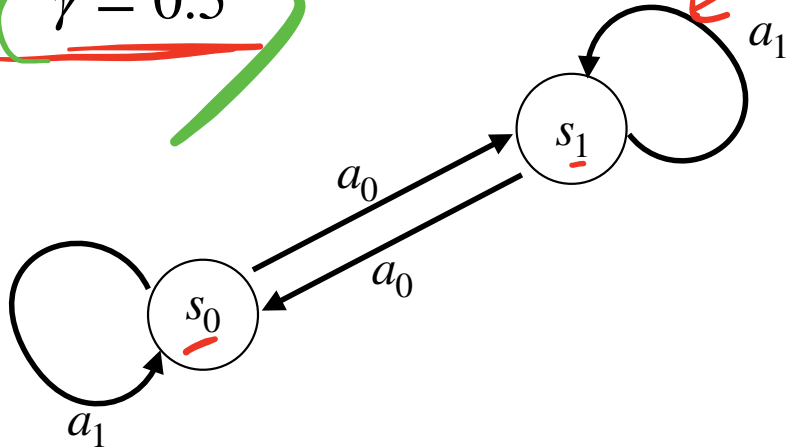
1. Initialization: Q^0

2. Iterate until convergence: $Q^{t+1} \leftarrow \mathcal{T} Q^t$

We hope $Q^t \rightarrow Q^*$, as $t \rightarrow \infty$

Exercise

Reward: $r(s_1, a_1) = 1$,
0 everywhere else,
 $\gamma = 0.5$



$$Q^0(s, a) = 0, \forall s, a$$

1. Compute $Q^1(s, a), Q^2(s, a), \forall s, a$

2. Compute $\|Q^i - Q^*\|_\infty$ for $i \in \{0, 1, 2\}$

3. How does $\|Q^i - Q^*\|_\infty$ behave as i increases

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

$$\|Q^0 - Q^*\|_\infty = 2$$

$$\|Q^1 - Q^*\|_\infty = 1$$

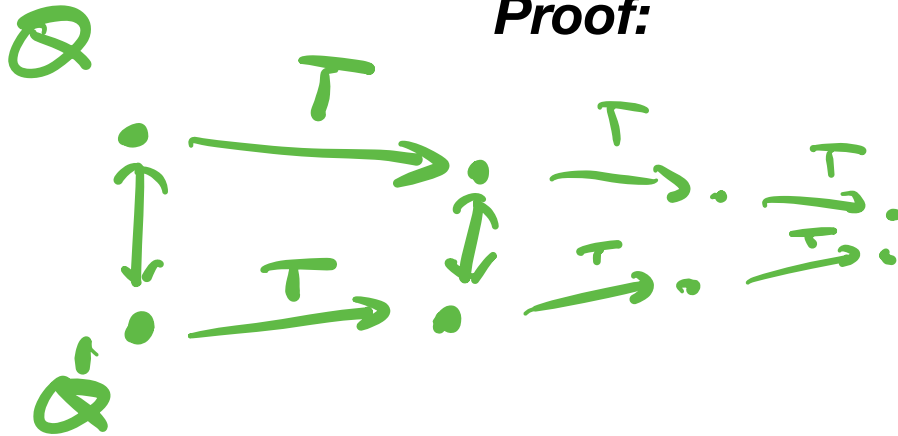
$$\|Q^2 - Q^*\|_\infty = \frac{1}{2}$$

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:

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

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:

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

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:

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

$E \in \max_{s'}$

$\underbrace{\hspace{10em}}_{\equiv \|Q - Q'\|_\infty}$

Convergence of Value Iteration:

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

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

Q $\xrightarrow{\mathcal{T}}$ \cdot
Q' $\xrightarrow{\mathcal{T}}$ \cdot

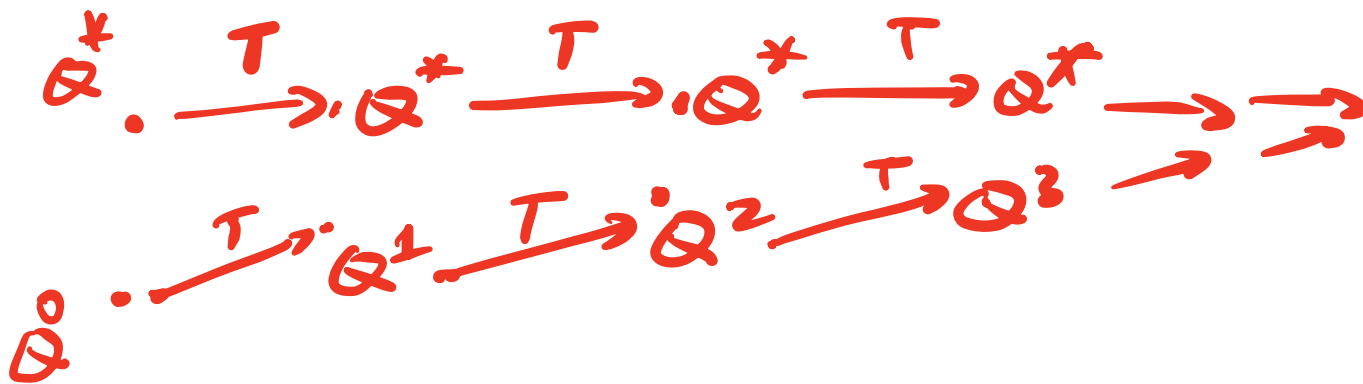
Proof:

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

Convergence of Value Iteration:

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

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



Convergence of Value Iteration:

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

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

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

Convergence of Value Iteration:

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

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

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

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

Summary so far:

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

VI (a fix point iteration alg):

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

$$\underline{\underline{Q^*}} \\ \arg \max_a Q^*(s_a)$$

VI convergence (via contraction)

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

Summary so far:

VI (a fix point iteration alg):

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

VI convergence (via contraction)

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

$$t = 1000 \quad Q^{1000} \approx Q^*$$

What about the policy? Ultimately, we do want π^* ...

$$\hat{\pi} \leftarrow \underset{\alpha}{\operatorname{argmax}} Q^{1000}(s, a)$$

