

Value Iteration

Lucas Janson and Sham Kakade

CS/Stat 184: Introduction to Reinforcement Learning
Fall 2022

Today

- Recap
- Today:
 - An Iterative Algorithm: Value Iteration
 - Visitation distributions

Recap

Infinite horizon Discounted Setting

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

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

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

Quantities that allow us to reason policy's long-term effect:

Value function $V^\pi(s) = \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid s_0 = s, \pi \right]$

Q function $Q^\pi(s, a) = \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid (s_0, a_0) = (s, a), \pi \right]$

Bellman Consistency Equations:

$$V^\pi(s) = r(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi(s))} V^\pi(s')$$

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

Summary so far:

Every discounted MDP has some deterministic optimal policy π^* ,
that dominates all other policies, everywhere

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

So we have, $V^* = V^{\pi^*}$ and $Q^* = Q^{\pi^*}$.

Bellman (Optimality) Equations

- V^* satisfies Bellman Equations:

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

- If V satisfies the Bellman Equations,

$$V(s) = \max_a [r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s')],$$

then $V = V^*$.

- The optimal policy is:

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

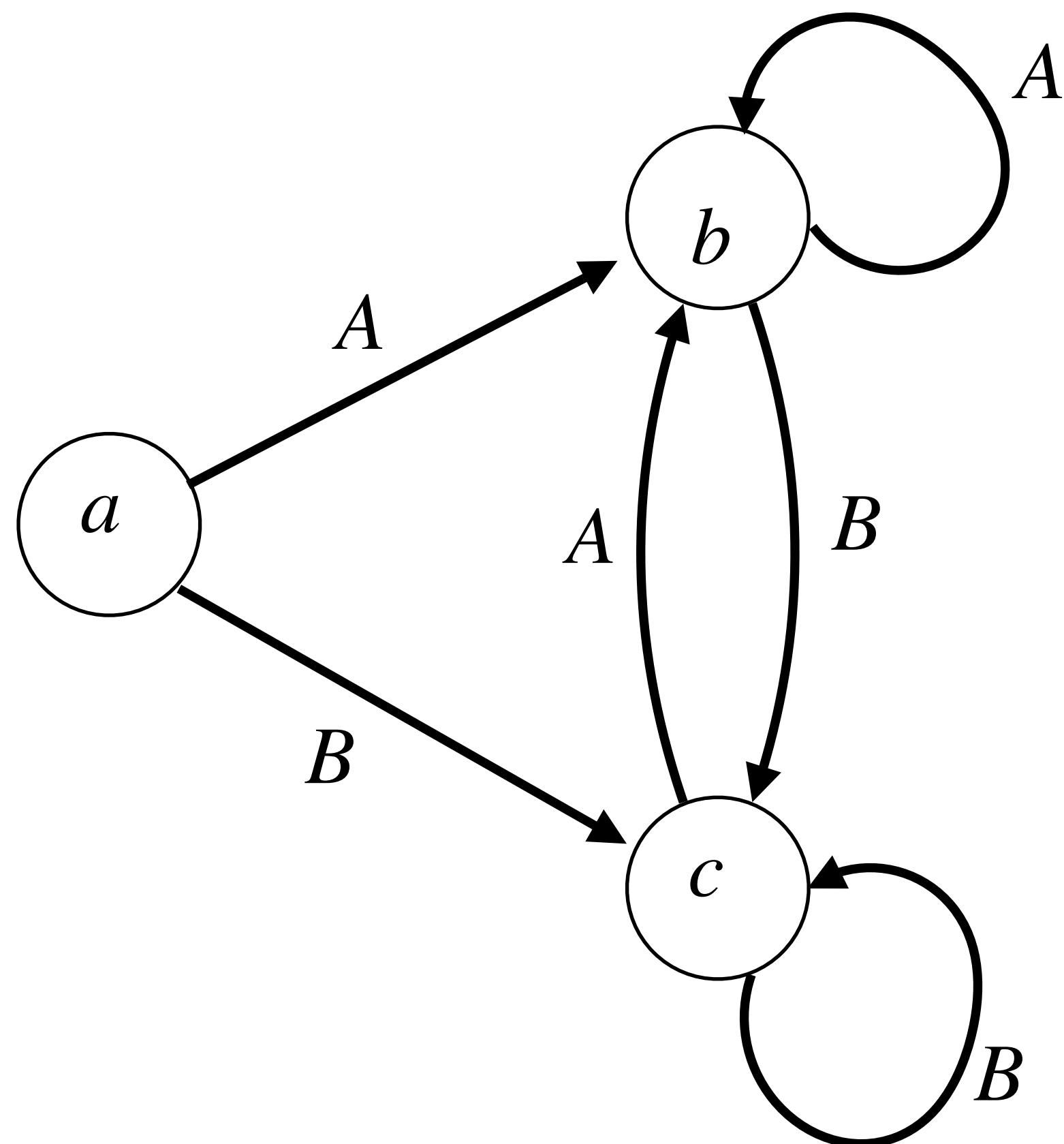
Today:
Value Iteration

Question for Today:

Given an MDP $\mathcal{M} = (S, A, P, r, \gamma)$,
how can we (approximately) find π^* ?

Example of Optimal Policy π^*

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

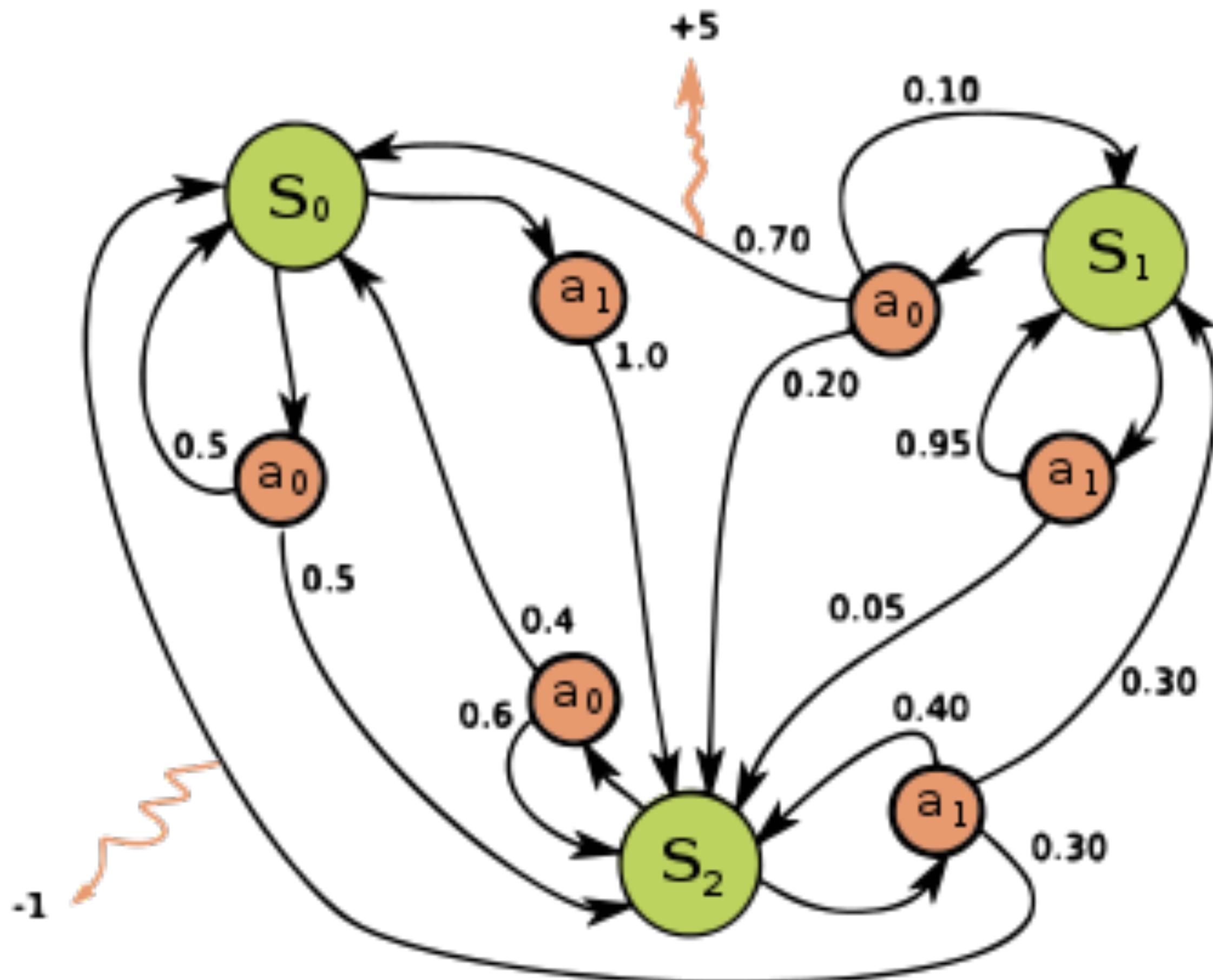


$$\pi^*(s) = A, \forall s$$

$$V^*(a) = \frac{\gamma}{1 - \gamma}, V^*(b) = \frac{1}{1 - \gamma}, V^*(c) = \frac{\gamma}{1 - \gamma}$$

Reward: $r(b, A) = 1, \& 0$ everywhere else

What about this one...



Let's design an algorithm that computes V^*/Q^* for any given MDP

Can we efficiently compute in a optimal policy? (polynomial in $|S|$, $|A|$, and other relevant quantities)

- Suppose we can efficiently compute $V^\pi(s)$ for any given $\pi : S \mapsto A$.
 - Brute force search would to find π^* would still take $|A|^{|S|}$ time.
- Can we construct an iterative algorithm based on the BEs?
 - Will it converge?
 - What is the computation time to get an approximate solution?

Detour: fix-point solution

Consider $x^* = f(x^*)$, $f: [a, b] \mapsto [a, b]$

A naive approach to find x^* :

Initialize $x^0 \in [a, b]$, repeat: $x^{t+1} = f(x^t)$

If f is a contraction mapping,
i.e., $\forall x, x', |f(x) - f(x')| \leq \gamma |x - x'|$, for some $\gamma \in [0, 1)$, then:
 $x^t \rightarrow x^*$, as $t \rightarrow \infty$

Define Bellman Operator \mathcal{T} :

Bellman Equations: $V(s) = \max_a [r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s')]$

- Any function $V: S \mapsto \mathbb{R}$ can also be viewed as a vector in $V \in \mathbb{R}^{|S|}$.

- Define $\mathcal{T}: \mathbb{R}^{|S|} \mapsto \mathbb{R}^{|S|}$, where

$$(\mathcal{T}V)(s) := \max_a [r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s')]$$

- Bellman equations in terms of \mathcal{T} :

$$\mathcal{T}V = V$$

Value Iteration Algorithm:

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

Guarantee of VI:

We will see this fix-point iteration converges, i.e., $V^t \rightarrow V^\star$, as $t \rightarrow \infty$

Alternative Version: Bellman Operator \mathcal{T} on \mathcal{Q}

(HW2 Q2 is the Q-version of the Bellman Equations)

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

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

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

What is the Per-Iteration Computational Complexity?

- Making the update $V^{t+1} \leftarrow \mathcal{T}V^t$ explicit:

- Define Q^{t+1} :

$$\forall s, a \quad Q^{t+1}(s, a) = r(s, a) + \gamma \sum_{s' \in S} P(s' | s, a) V^t(s')$$

- Set V^{t+1} :

$$\forall s \quad V^{t+1}(s) = \max_a Q^{t+1}(s, a)$$

- What is the order of the number of basic arithmetic operations?

$$O(|S|^2 |A|)$$

With matrix multiplication?

- Making the update $V^{t+1} \leftarrow \mathcal{T}V^t$ explicit:

- Define Q^{t+1} :

$$\forall s, a \quad Q^{t+1}(s, a) = r(s, a) + \gamma \sum_{s' \in S} P(s' | s, a) V^t(s')$$

- Set V^{t+1} :

$$\forall s \quad V^{t+1}(s) = \max_a Q^{t+1}(s, a)$$

- In terms of matrix multiplication,
let us view r as a vector, $r \in \mathbb{R}^{|S| \cdot |A|}$ and P as a matrix $P \in \mathbb{R}^{|S| \cdot |A| \times |S|}$

$$Q^{t+1} = r + \gamma P V^t$$

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 V, V' , we have:

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

Convergence of Value Iteration:

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

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

Proof:

$$\begin{aligned} |(\mathcal{T}V)(s) - (\mathcal{T}V')(s)| &= \left| \max_a \left\{ r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s') \right\} - \max_a \left\{ r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V'(s') \right\} \right| \\ &\leq \max_a \left| r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s') - \left(r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V'(s') \right) \right| \\ &= \gamma \max_a \left| \mathbb{E}_{s' \sim P(s, a)} [V(s') - V'(s')] \right| \\ &\leq \gamma \max_a \mathbb{E}_{s' \sim P(s, a)} [|V(s') - V'(s')|] \\ &\leq \gamma \max_s |V(s') - V'(s')| = \gamma \|V - V'\|_\infty \end{aligned}$$

Convergence of Value Iteration:

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

$$\|V^t - V^\star\|_\infty \leq \gamma^t \|V^0 - V^\star\|_\infty$$

Proof:

$$\|V^t - V^\star\|_\infty = \|\mathcal{T}V^{t-1} - \mathcal{T}V^\star\|_\infty \leq \gamma \|V^{t-1} - V^\star\|_\infty$$

$$\dots \leq \gamma^t \|V^0 - V^\star\|_\infty$$

Computational Complexity of VI

VI will return a V^t s.t. $\|V^t - V^*\|_\infty \leq \epsilon$ in no more than,

$$\frac{\ln(\|V^0 - V^*\|_\infty / \epsilon)}{\ln(1/\gamma)} \leq \frac{\ln(\|V^0 - V^*\|_\infty / \epsilon)}{(1 - \gamma)} \text{ iterations.}$$

So the computational complexity for an ϵ -accurate solution is

$$O\left(\frac{|S|^2 |A|}{1 - \gamma} \ln\left(\frac{1}{\epsilon(1 - \gamma)}\right)\right)$$

But what about the policy we find with VI?

Theorem: For any V , let $\pi(s) = \arg \max_a [r(s, a) + \mathbb{E}_{s' \sim P(s, a)} V(s')]$, then

$$V^\pi(s) \geq V^*(s) - \frac{2\gamma}{1-\gamma} \|V - V^*\|_\infty$$

But what about the policy we find with VI?

Theorem: For any V , let $\pi(s) = \arg \max_a [r(s, a) + \mathbb{E}_{s' \sim P(s, a)} V(s')]$, then

$$V^\pi(s) \geq V^*(s) - \frac{2\gamma}{1-\gamma} \|V - V^*\|_\infty$$

Runtime: After $\frac{\ln(2/((1-\gamma)^2\epsilon))}{1-\gamma}$ iterations of PI, we have: $V^{\pi^t}(s) \geq V^*(s) - \epsilon$,

and, the total runtime of VI is:

$$O\left(\frac{|S|^2 |A|}{1-\gamma} \ln\left(1/((1-\gamma)^2\epsilon)\right)\right)$$

We replace $\epsilon \leftarrow (1-\gamma)\epsilon/2$, then VI will return V^t s.t. $\|V^t - V^*\|_\infty \leq (1-\gamma)\epsilon/2$.

Thus, $V^{\pi^t}(s) \geq V^*(s) - \frac{2\gamma}{1-\gamma} \|V^t - V^*\|_\infty \geq V^*(s) - \epsilon$

$$\text{Proof of } \|V^* - V^\pi\| \leq \frac{2\gamma}{1-\gamma} \|V - V^*\|_\infty$$

Proof:

- To prove the theorem, it suffice to show that:

$$\|V^\pi - V^*\| \leq \frac{2\gamma}{1-\gamma} \|V - V^*\|_\infty$$

- We have:

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

Term 1

Term 2

- For term 2,

$$\begin{aligned} Q^*(s, \pi(s)) - Q^\pi(s, \pi(s)) &= \gamma \mathbb{E}_{s' \sim P(s, \pi(s))} (V^*(s') - V^\pi(s')) \\ &\leq \gamma \|V^* - V^\pi\|_\infty \end{aligned}$$

Continuing....

- Let $Q(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s')$, which implies that $Q(s, \pi^*(s)) \leq Q(s, \pi^t(s))$
- Using this for term 1,

$$\begin{aligned} Q^*(s, \pi^*(s)) - Q^*(s, \pi^t(s)) &\leq Q^*(s, \pi^*(s)) - Q(s, \pi^*(s)) + Q(s, \pi(s)) - Q^*(s, \pi(s)) \\ &\leq \gamma \mathbb{E}_{s' \sim P(s, \pi^*(s))} [V^*(s') - V(s')] - \gamma \mathbb{E}_{s' \sim P(s, \pi(s))} [V(s') - V^*(s')] \\ &\leq 2\gamma \|V - V^*\|_\infty \end{aligned}$$

- Therefore, we have:

$$V^*(s) - V^{\pi^t}(s) \leq 2\gamma \|V - V^*\|_\infty + \gamma \|V^* - V^{\pi^t}\|_\infty$$

- which complete the proof because we have shown:

$$\|V^* - V^{\pi^t}\|_\infty \leq 2\gamma \|V - V^*\|_\infty + \gamma \|V^* - V^{\pi^t}\|_\infty$$

Understanding the sampling

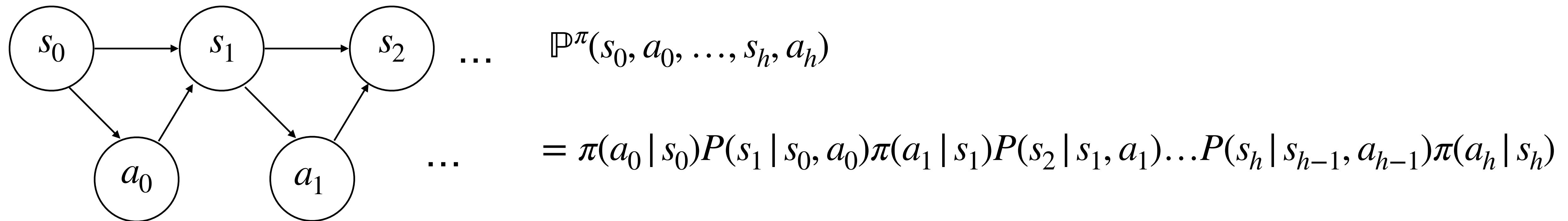
“Occupancy measures” are a helpful concept

Discounted State (action) Occupancy Measures

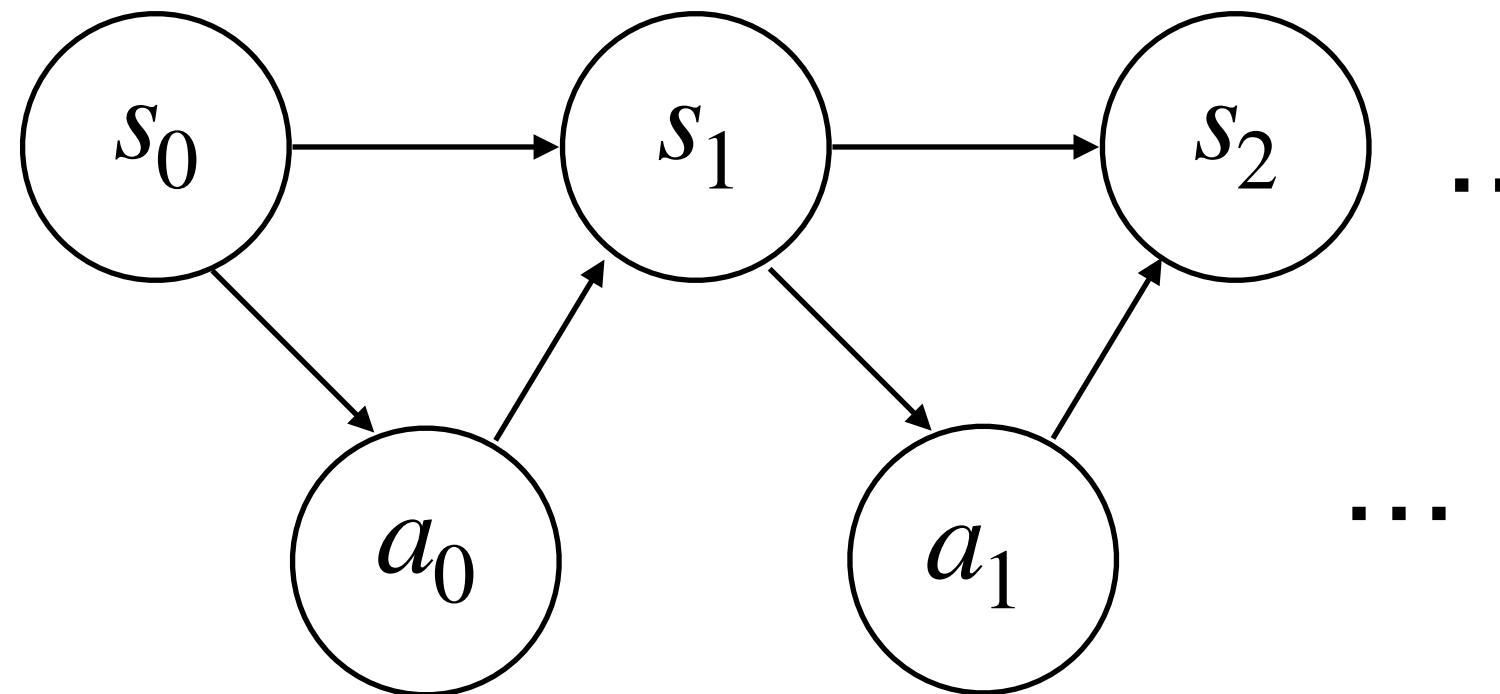
Assume we start at s_0 , following π to step h , what's probability of seeing a trajectory:

$$(s_0, a_0, s_1, a_1, \dots, s_h, a_h)?$$

Let's write π as a delta distribution, i.e., $\pi(a | s) = \begin{cases} 1, & a = \pi(s), \\ 0, & \text{else} \end{cases}$



State-action distribution at time step h



$$\mathbb{P}^\pi(a_0, \dots, s_h, a_h | s_0, \pi)$$

$$= \pi(a_0 | s_0) P(s_1 | s_0, a_0) \pi(a_1 | s_1) P(s_2 | s_1, a_1) \dots P(s_h | s_{h-1}, a_{h-1}) \pi(a_h | s_h)$$

Q: what's the probability of π visiting state (s, a) at time step h ?

$$\mathbb{P}_h^\pi(s_h, a_h | s_0, \pi) = \sum_{a_0, s_1, a_1, \dots, s_{h-1}, a_{h-1}} \mathbb{P}^\pi(a_0, \dots, s_{h-1}, a_{h-1}, s_h = s, a_h = a | s_0, \pi)$$

Discounted Average State-action distribution

Probability of π visiting (s, a) at h , starting from s_0

$$\mathbb{P}_h^\pi(s_h, a_h | s_0, \pi) = \sum_{a_0, s_1, a_1, \dots, s_{h-1}, a_{h-1}} \mathbb{P}^\pi(s_0, a_0, \dots, s_{h-1}, a_{h-1} | s_h = s, a_h = a)$$

$$d_{s_0}^\pi(s, a) = (1 - \gamma) \sum_{h=0}^{\infty} \gamma^h \mathbb{P}_h^\pi(s, a; s_0)$$

Can you show that this is a valid distribution?

$$V^\pi(s_0) = \frac{1}{1 - \gamma} \sum_{s,a} d_{s_0}^\pi(s, a) r(s, a)$$

HW0 questions!

Can you show the above is true?

Summary Today

- Value iteration: an iterative algorithm with a “linear” convergence rate.
- The concept of an “occupancy measure”.

1-minute feedback form: <https://bit.ly/3RHtIxy>

