

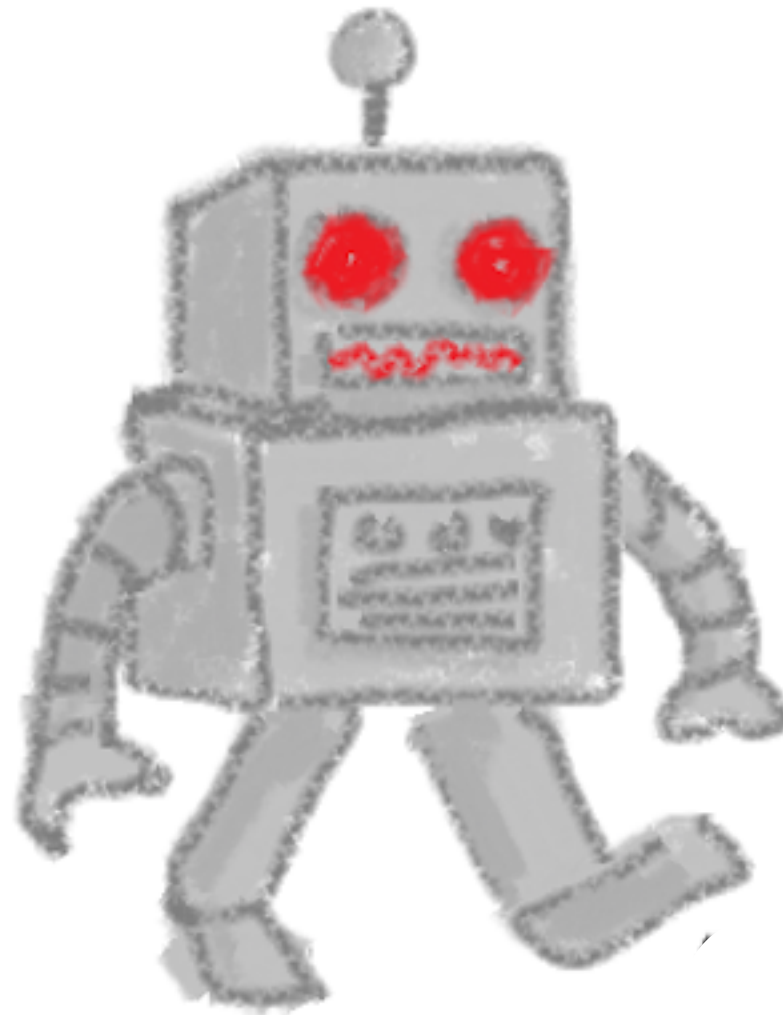
Basics of Markov Decision Process

Announcement:

1. For wait list, we will need to prioritize CS students and seniors (cs-course-enroll@cornell.edu)
2. Clarification on the attendance bonus (5%)

Recap:

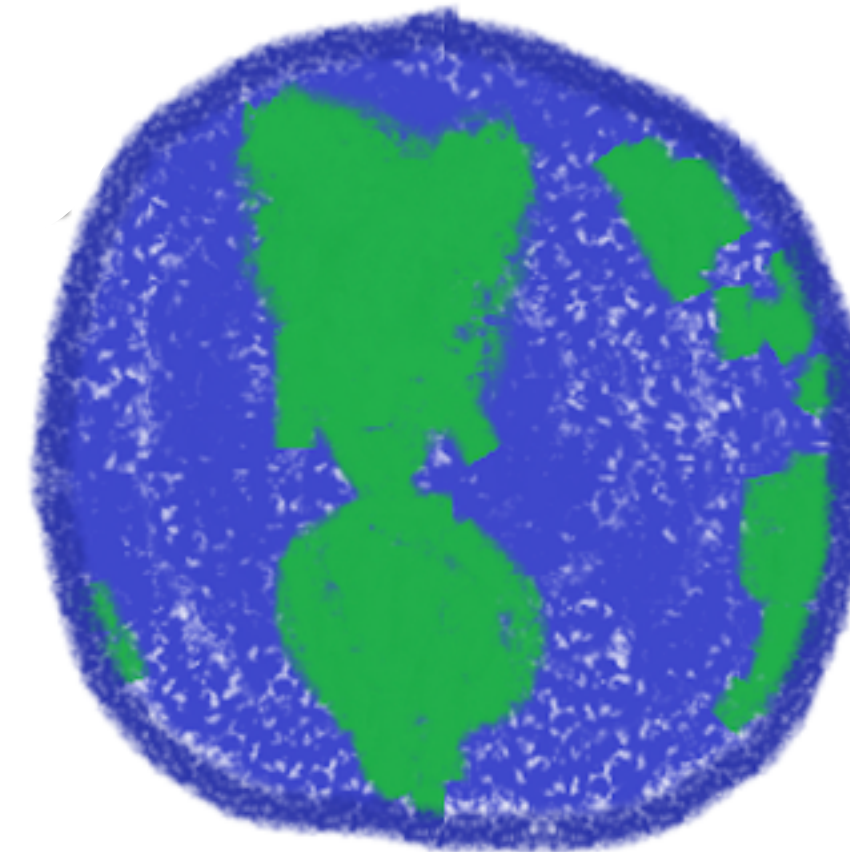
Learning Agent



$$\pi(s) \rightarrow a$$

Policy: determine **action** based on **state**

Environment



Multiple Steps

Send **reward** and **next state** from a Markovian transition dynamics

$$r(s, a), s' \sim P(\cdot | s, a)$$

Recap:

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

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

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

Recap:

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

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

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

$$\text{Value function } V^\pi(s) = \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid s_0 = s, a_h = \pi(s_h), s_{h+1} \sim P(\cdot \mid s_h, a_h) \right]$$

Recap:

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

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

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

$$\text{Value function } V^\pi(s) = \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid s_0 = s, a_h = \pi(s_h), s_{h+1} \sim P(\cdot \mid s_h, a_h) \right]$$

$$\text{Q function } Q^\pi(s, a) = \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid (s_0, a_0) = (s, a), a_h = \pi(s_h), s_{h+1} \sim P(\cdot \mid s_h, a_h) \right]$$

Bellman Equation for V/Q-function:

$$V^\pi(s) = \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid s_0 = s, a_h = \pi(s_h), s_{h+1} \sim P(\cdot \mid s_h, a_h) \right]$$

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

Bellman Equation for V/Q-function:

$$V^\pi(s) = \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid s_0 = s, a_h = \pi(s_h), s_{h+1} \sim P(\cdot \mid s_h, a_h) \right]$$

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

$$Q^\pi(s, a) = \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid (s_0, a_0) = (s, a), a_h = \pi(s_h), s_{h+1} \sim P(\cdot \mid s_h, a_h) \right]$$

What is the BE for Q function??

Bellman Equation for V/Q-function:

$$V^\pi(s) = \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid s_0 = s, a_h = \pi(s_h), s_{h+1} \sim P(\cdot \mid s_h, a_h) \right]$$

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

$$Q^\pi(s, a) = \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid (s_0, a_0) = (s, a), a_h = \pi(s_h), s_{h+1} \sim P(\cdot \mid s_h, a_h) \right]$$

What is the BE for Q function??

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

Bellman Equation for Q-function:

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

Today:

1. We have A^S many policies, which one is the optimal policy π^\star ?
2. Key property of the optimal policy π^\star
3. State-action distributions

Definition of 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]

Definition of 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!

Definition of Optimal Policy π^\star

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

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

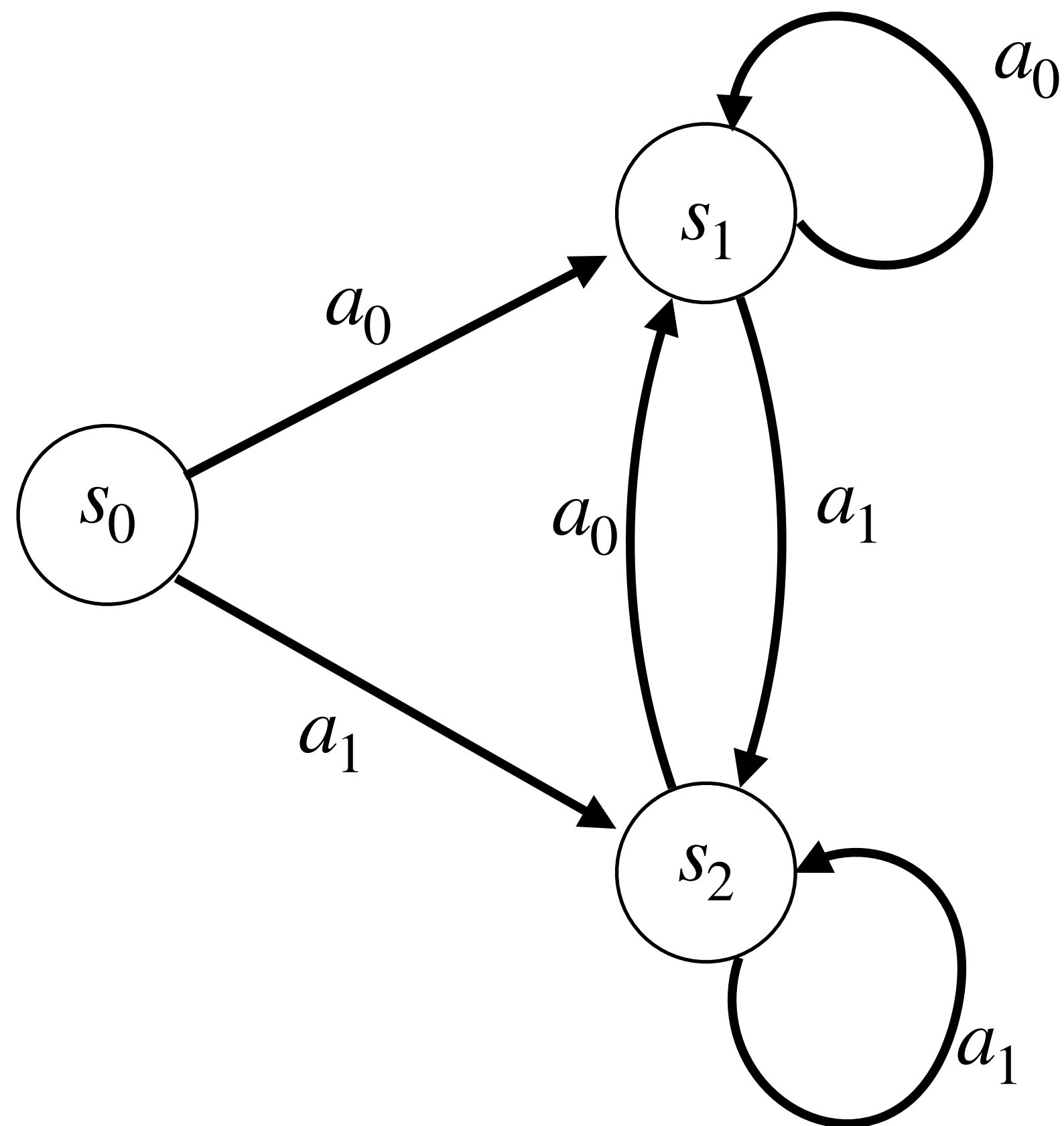
[Puterman 94 chapter 6, also see theorem 1.7 in the RL monograph—no need to understand the proof]

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

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

Example of Optimal Policy π^\star

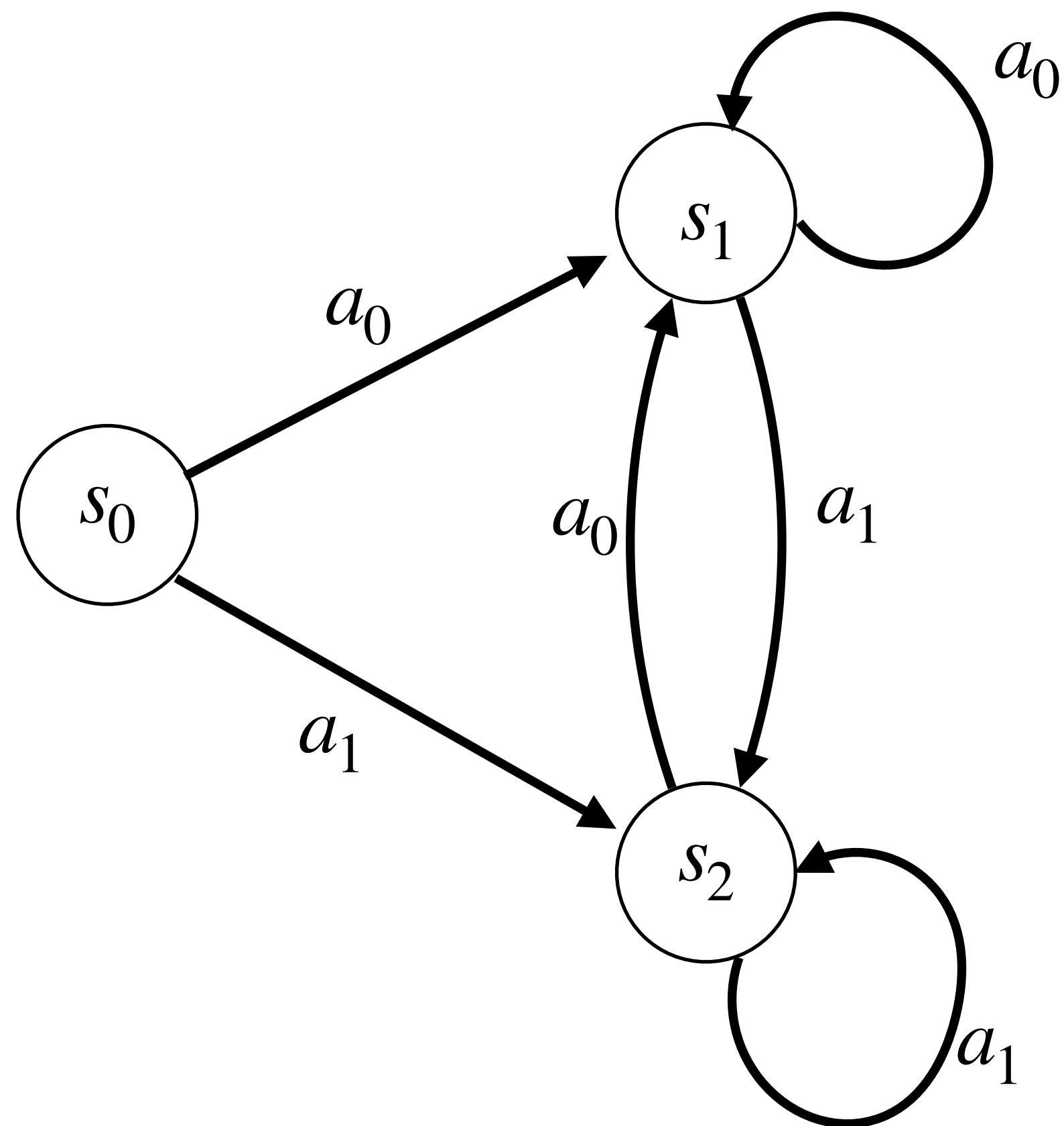
Consider the following **deterministic** MDP w/ 3 states & 2 actions



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

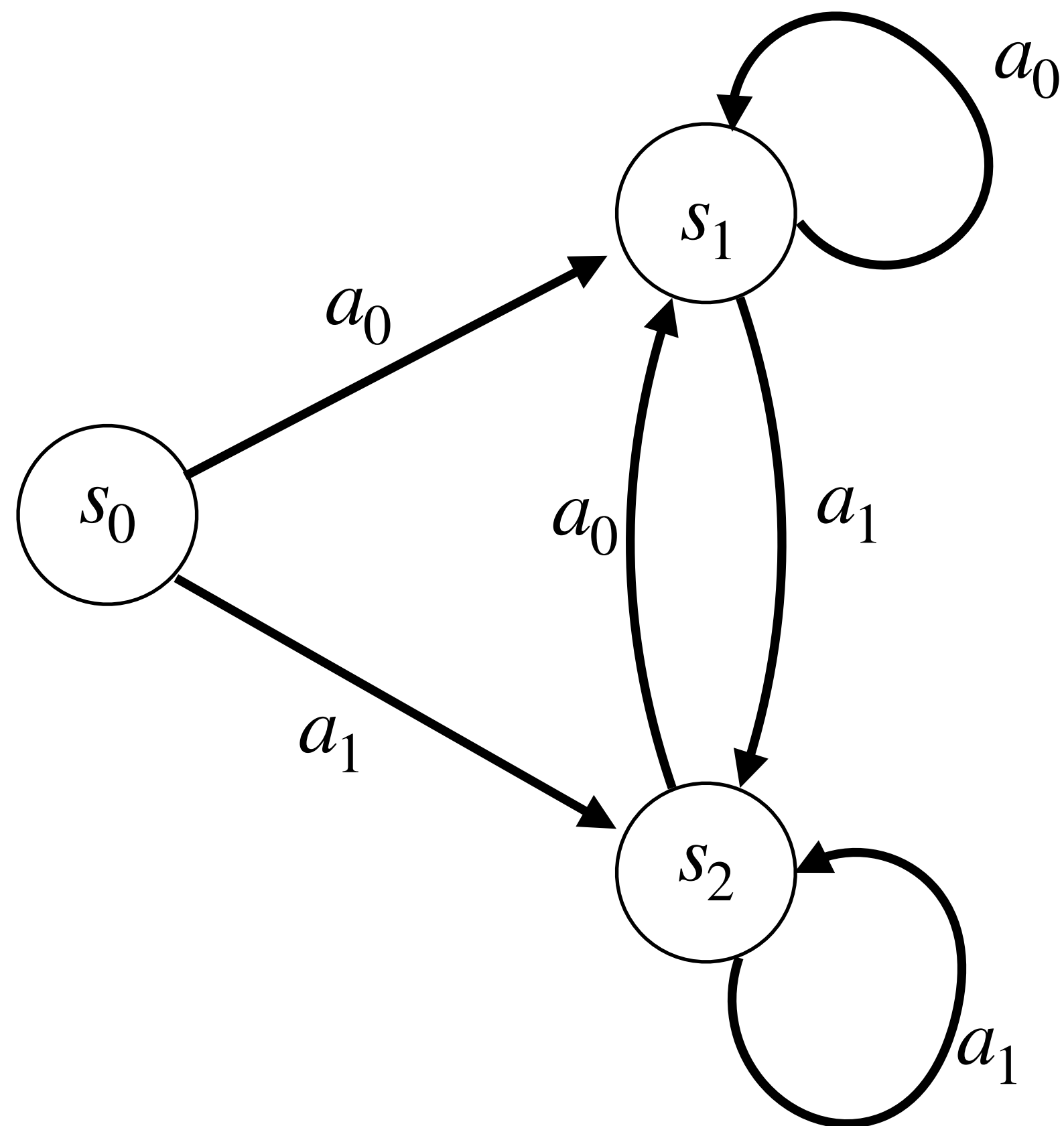


Let's say $\gamma \in (0,1)$
What's the optimal policy?

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



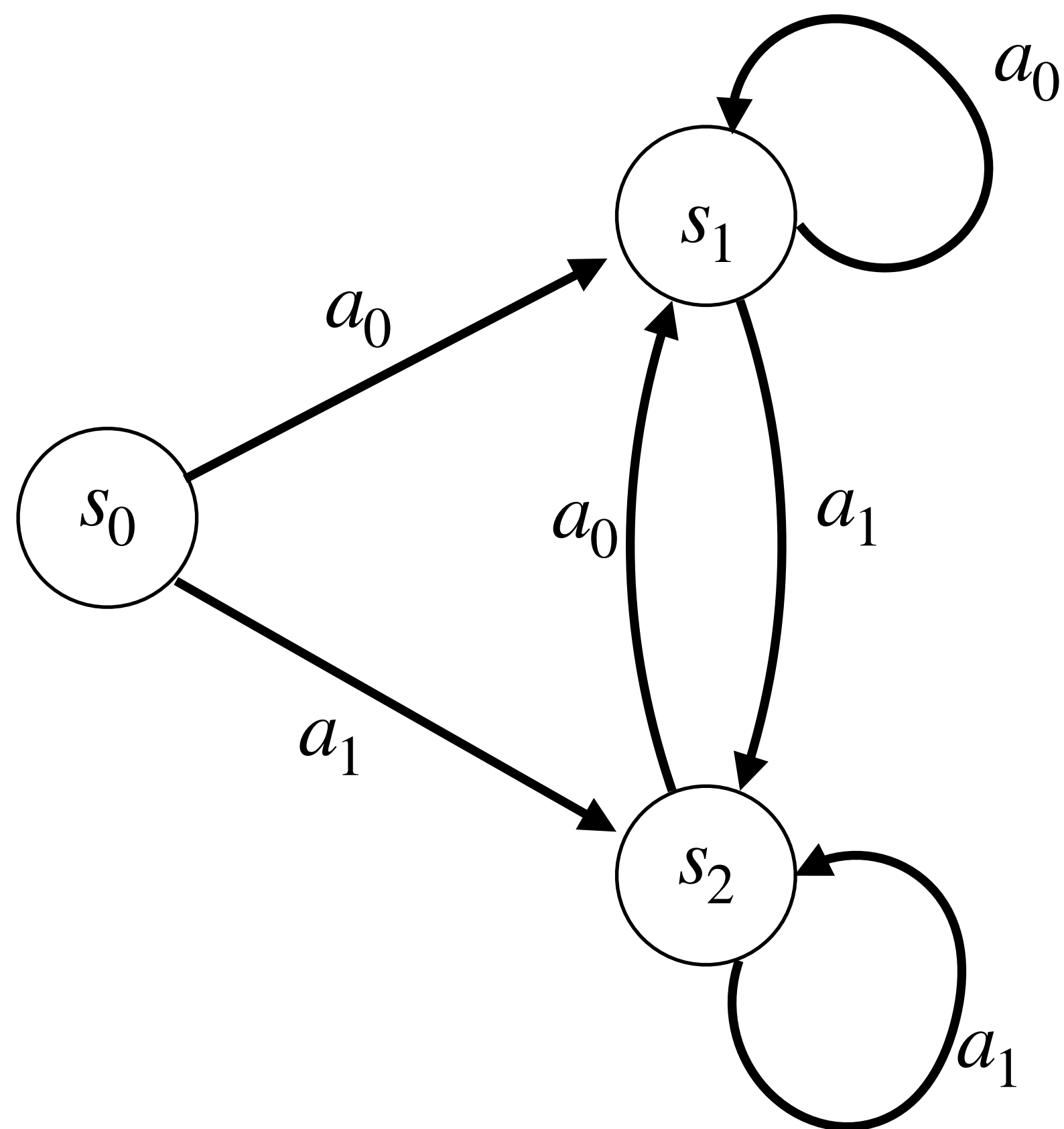
Let's say $\gamma \in (0,1)$
What's the optimal policy?

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

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



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

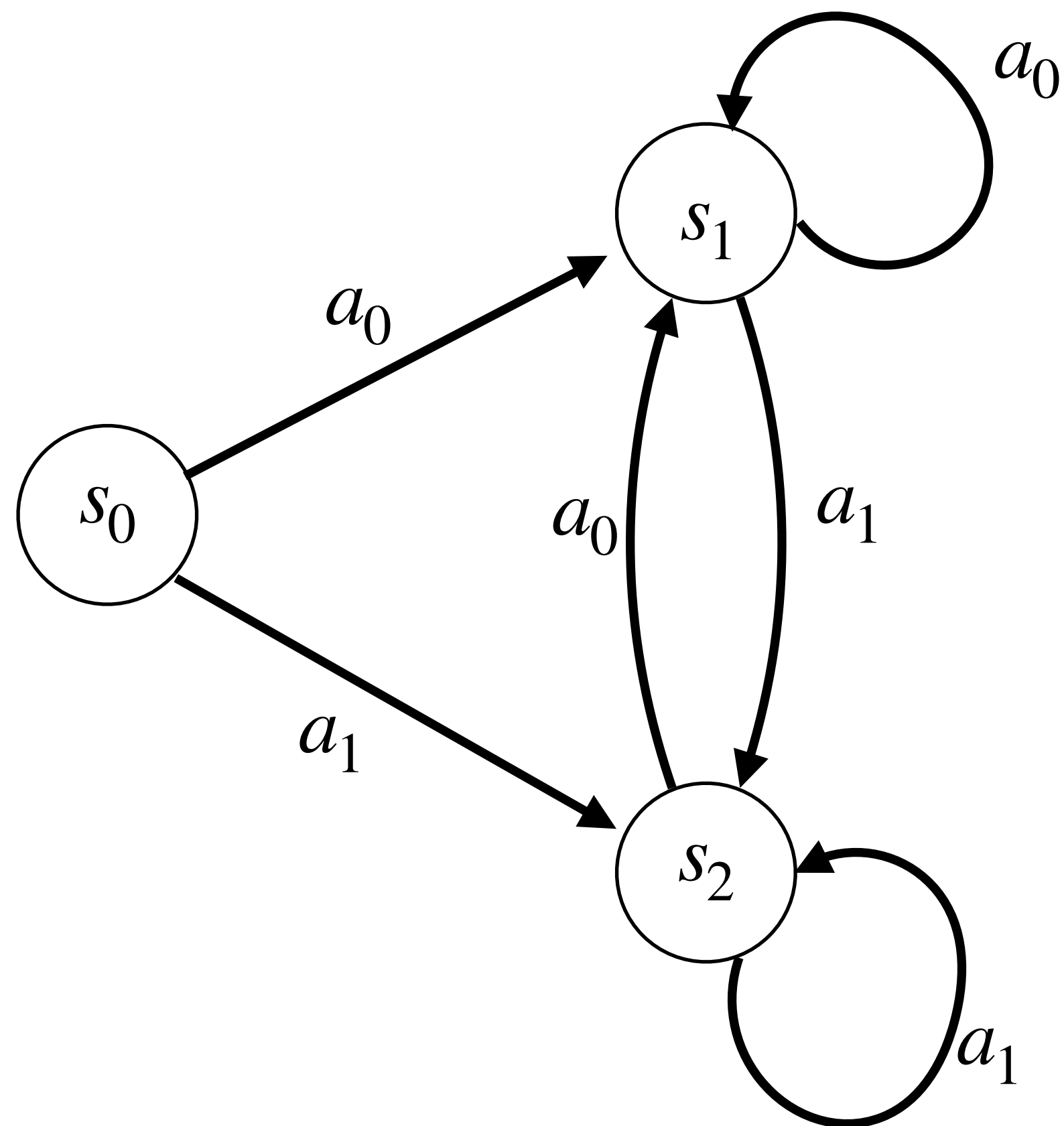
$$\pi^\star(s) = a_0, \forall s$$

$$V^\star(s_0) = \frac{\gamma}{1-\gamma}, V^\star(s_1) = \frac{1}{1-\gamma}, V^\star(s_2) = \frac{\gamma}{1-\gamma}$$

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



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

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

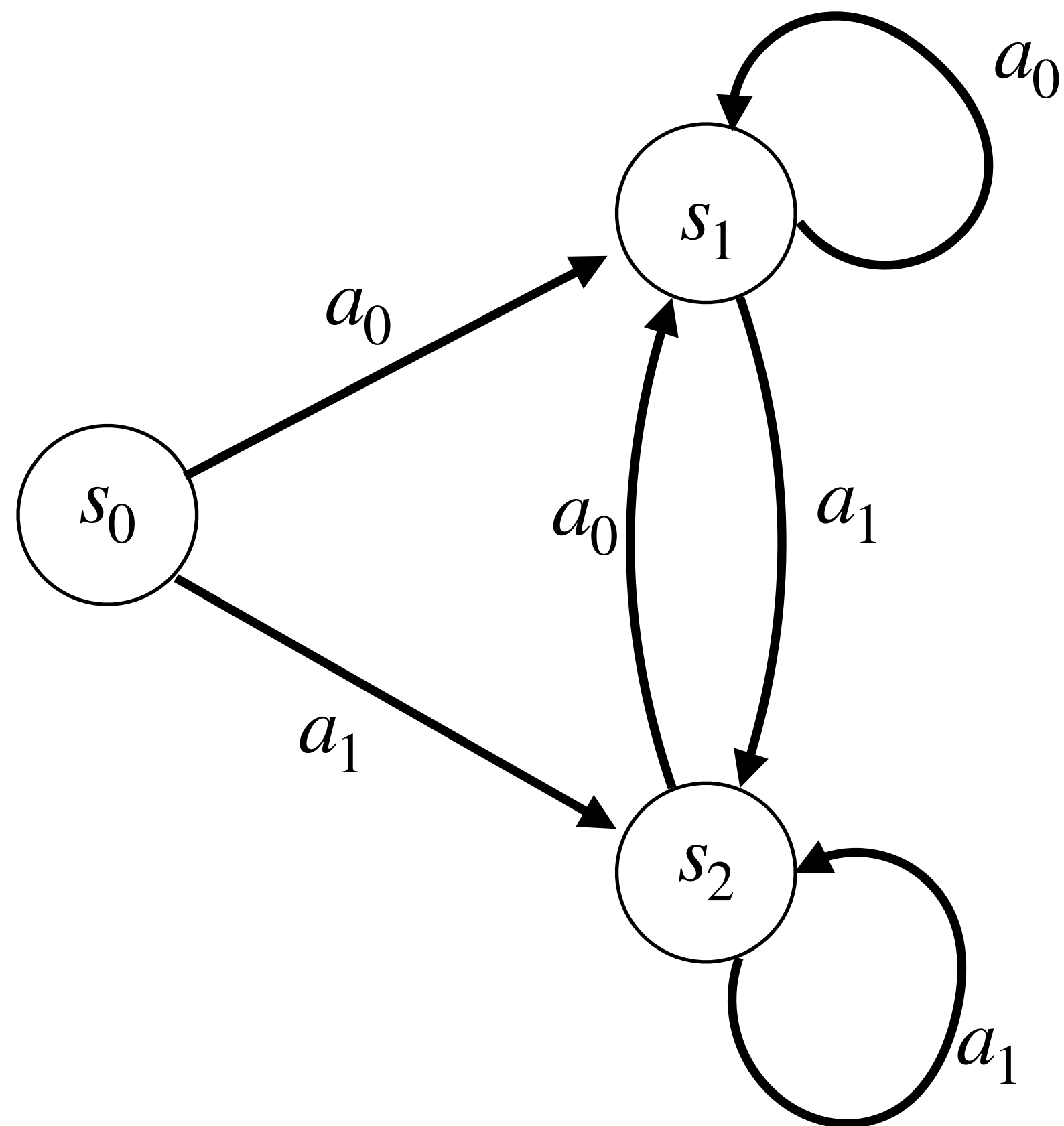
$$V^\star(s_0) = \frac{\gamma}{1-\gamma}, V^\star(s_1) = \frac{1}{1-\gamma}, V^\star(s_2) = \frac{\gamma}{1-\gamma}$$

What about policy $\pi(s) = a_1, \forall s$

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



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

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

$$V^\star(s_0) = \frac{\gamma}{1-\gamma}, V^\star(s_1) = \frac{1}{1-\gamma}, V^\star(s_2) = \frac{\gamma}{1-\gamma}$$

What about policy $\pi(s) = a_1, \forall s$

$$V^\pi(s_0) = 0, V^\pi(s_1) = 0, V^\pi(s_2) = 0$$

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

Summary so far:

Every discounted MDP has a deterministic optimal policy, that **dominates other policies everywhere** (proof is out of the scope)

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

Outline

✓ 1. We have A^S many policies, which one is the optimal policy π^\star ?

2. Key property of the optimal policy π^\star

3. State-action distributions

Bellman Optimality

Theorem 1: Bellman Optimality

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

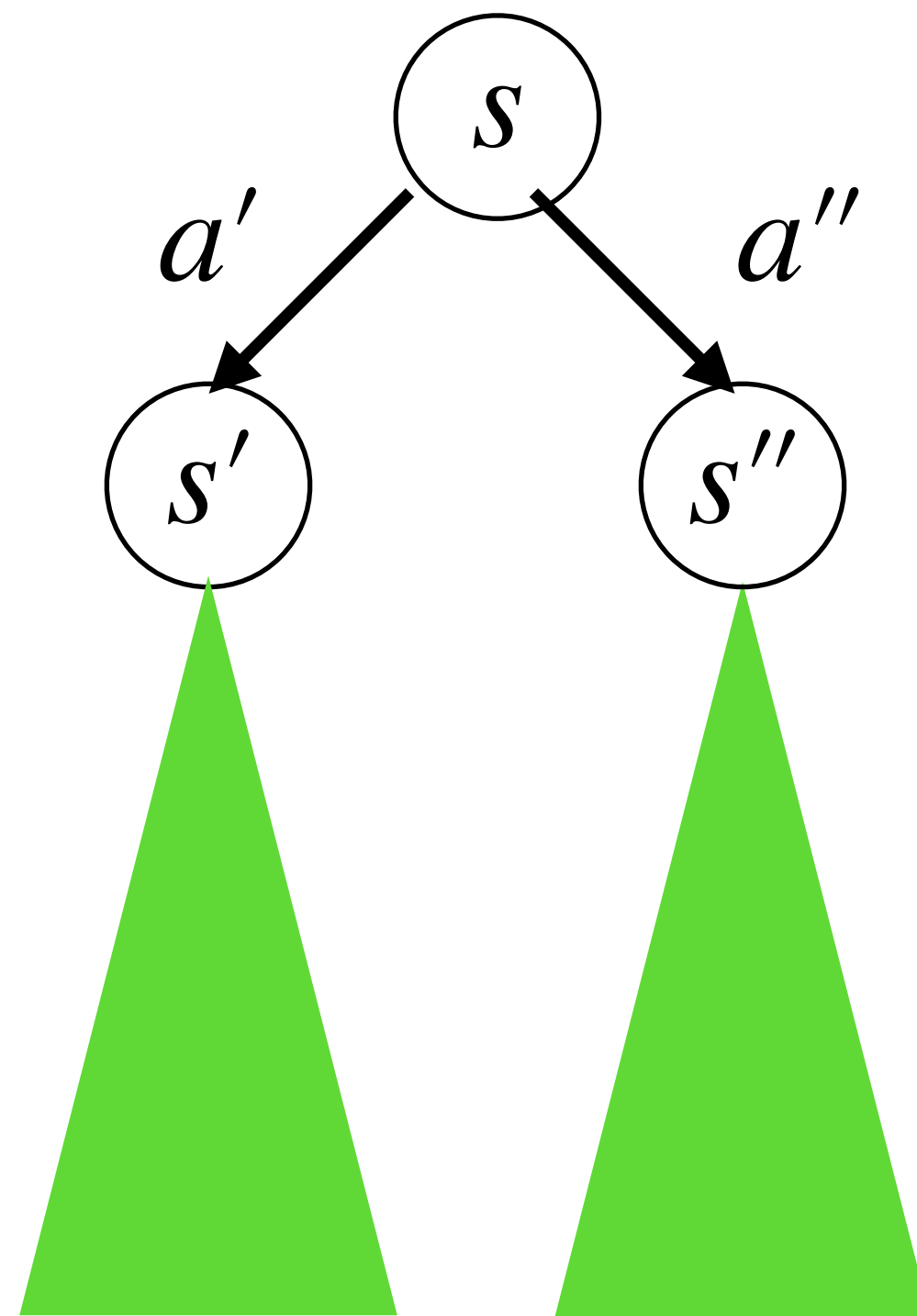
Understanding Bellman Optimality

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

Understanding Bellman Optimality

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

Q: If we know the optimal value at s' , s'' , i.e., $V^*(s')$, $V^*(s'')$, what we do at s ?

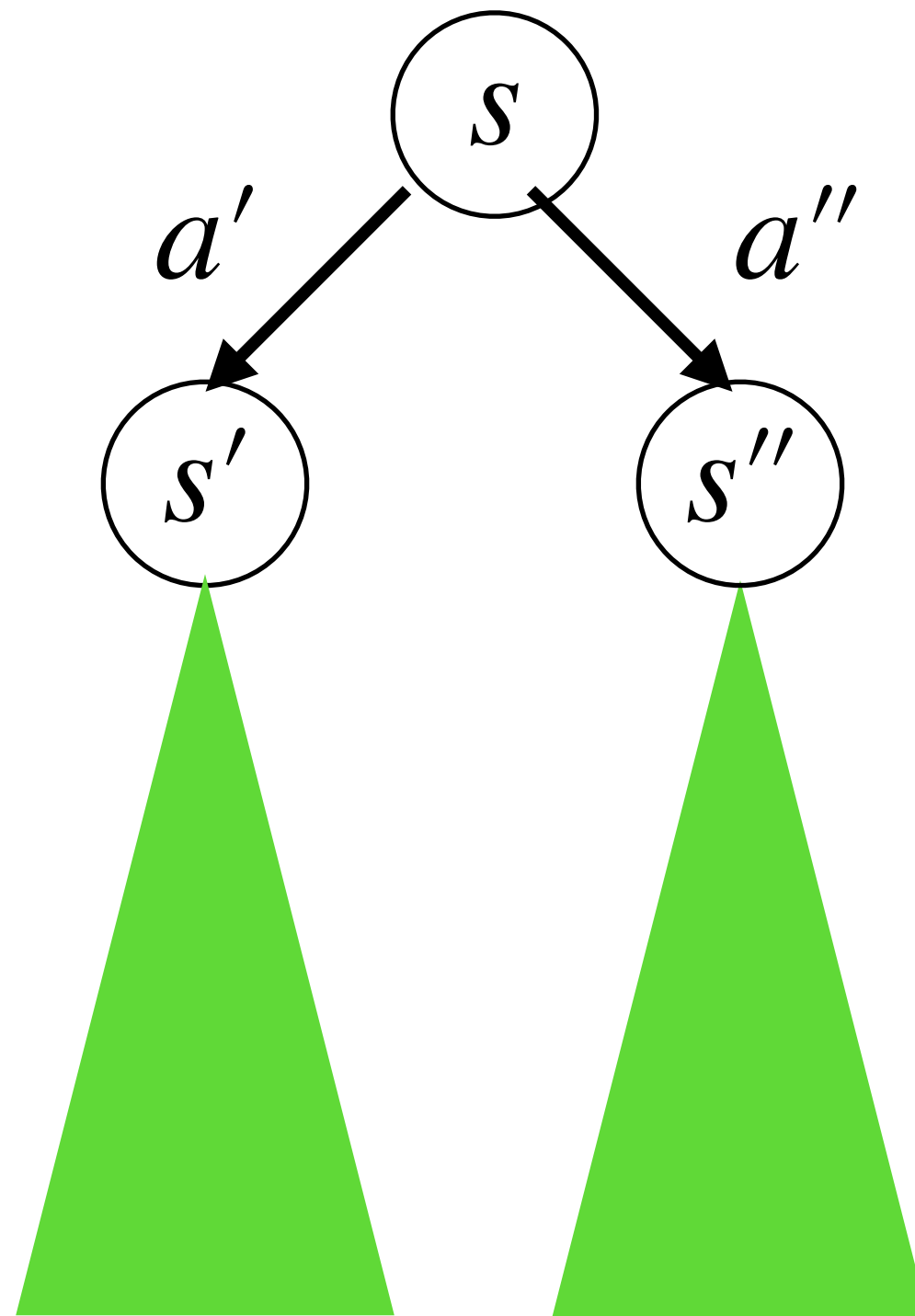


Understanding Bellman Optimality

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

Q: If we know the optimal value at s' , s'' , i.e., $V^*(s')$, $V^*(s'')$, what we do at s ?

1. Try a' , we get
 $Q^*(s, a') := r(s, a') + \gamma V^*(s')$

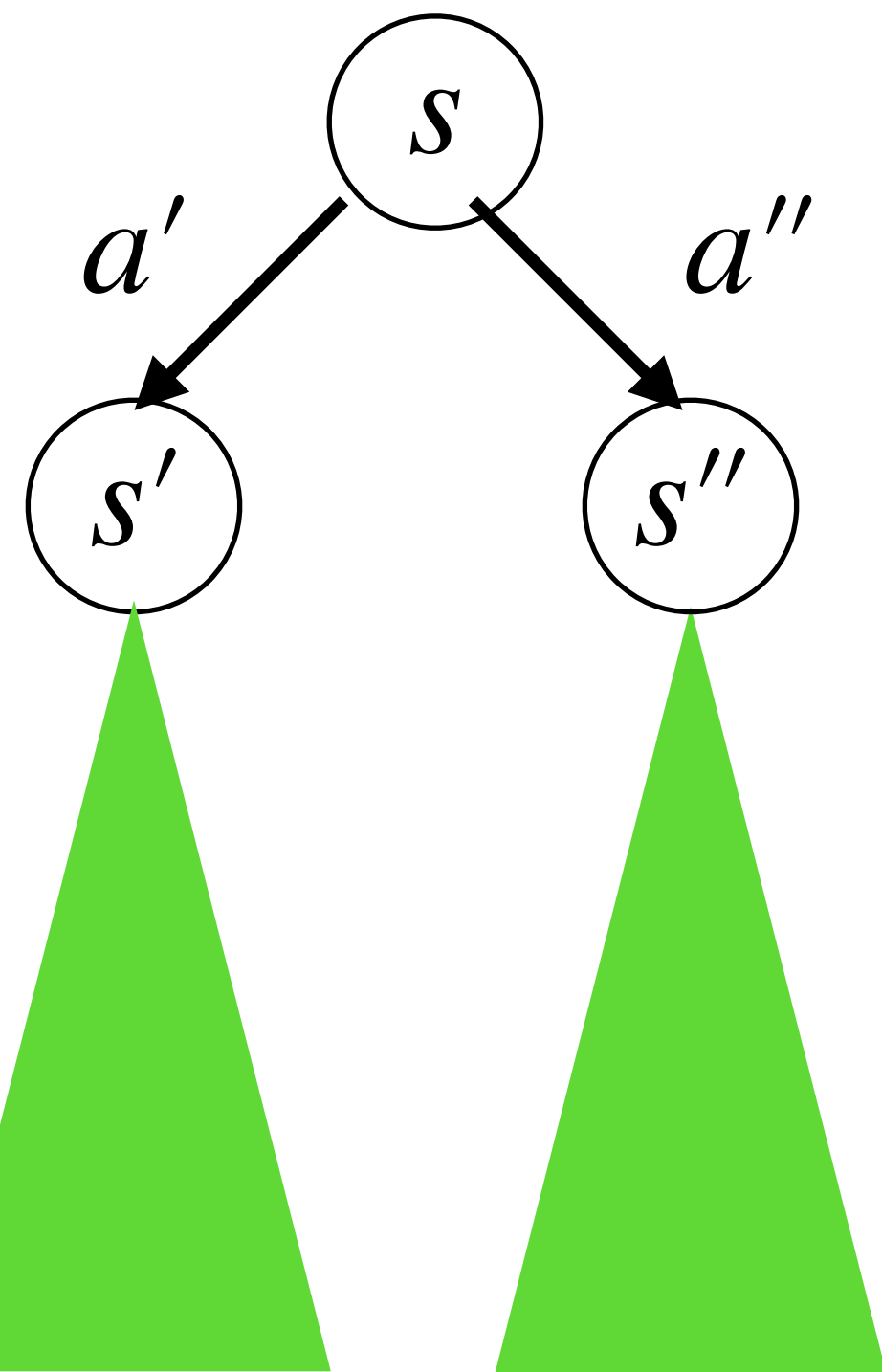


Understanding Bellman Optimality

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

Q: If we know the optimal value at s' , s'' , i.e., $V^*(s')$, $V^*(s'')$, what we do at s ?

1. Try a' , we get
 $Q^*(s, a') := r(s, a') + \gamma V^*(s')$



2. Try a'' , we get
 $Q^*(s, a'') := r(s, a'') + \gamma V^*(s'')$

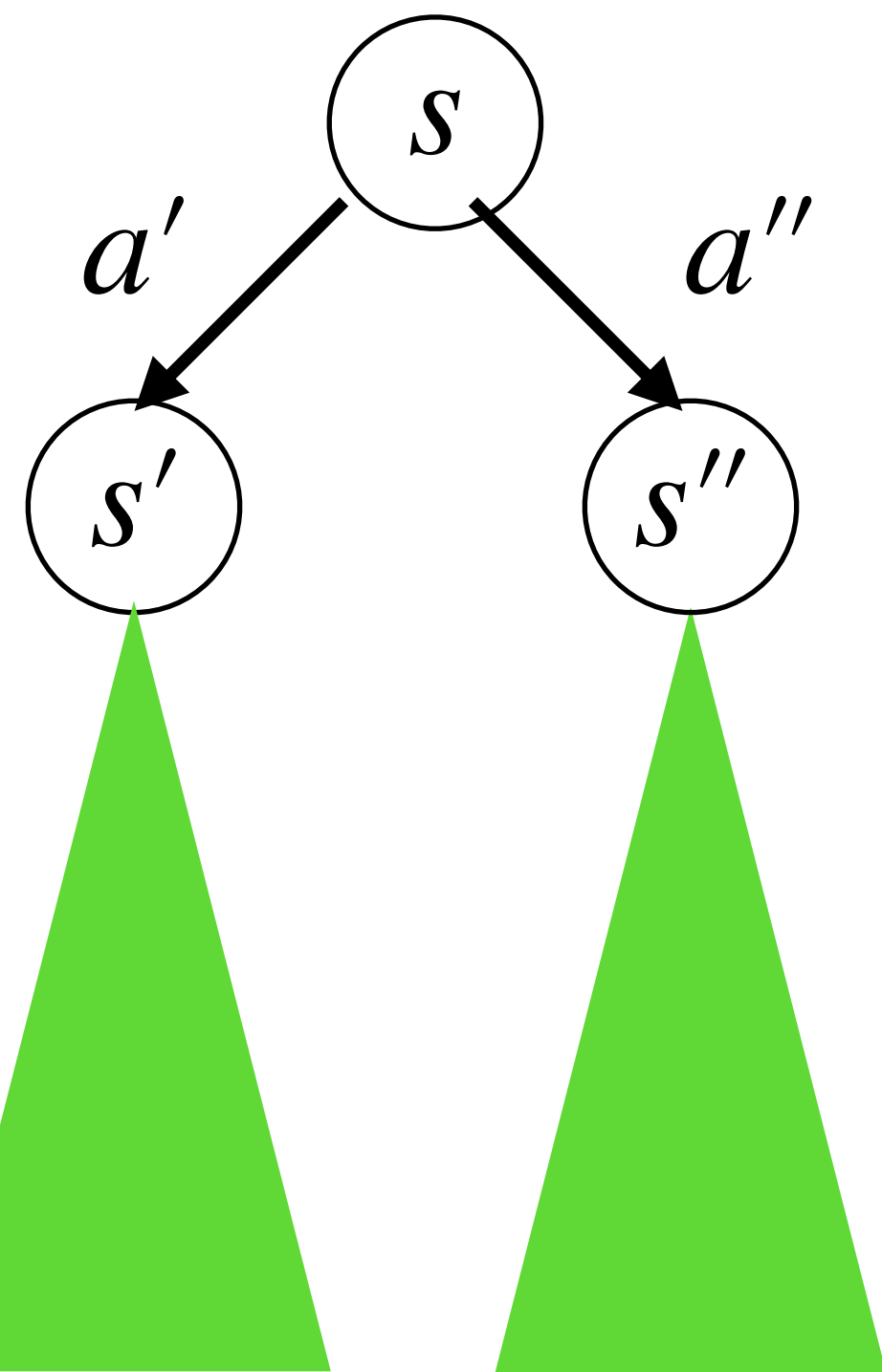
Understanding Bellman Optimality

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

Q: If we know the optimal value at s' , s'' , i.e., $V^*(s')$, $V^*(s'')$, what we do at s ?

1. Try a' , we get

$$Q^*(s, a') := r(s, a') + \gamma V^*(s')$$



2. Try a'' , we get

$$Q^*(s, a'') := r(s, a'') + \gamma V^*(s'')$$

$$V^*(s) = \max_{a', a''} \{ Q^*(s, a'), Q^*(s, a'') \}$$

Proof of Bellman Optimality

Theorem 1: Bellman Optimality

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

Proof of Bellman Optimality

Theorem 1: Bellman Optimality

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

Denote $\hat{\pi}(s) := \arg \max_a Q^*(s, a)$, we will prove $V^{\hat{\pi}}(s) = V^*(s), \forall s$

Proof of Bellman Optimality

Theorem 1: Bellman Optimality

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

Denote $\hat{\pi}(s) := \arg \max_a Q^*(s, a)$, we will prove $V^{\hat{\pi}}(s) = V^*(s), \forall s$

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

Proof of Bellman Optimality

Theorem 1: Bellman Optimality

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

Denote $\hat{\pi}(s) := \arg \max_a Q^*(s, a)$, we will prove $V^{\hat{\pi}}(s) = V^*(s), \forall s$

$$\begin{aligned} V^*(s) &= r(s, \pi^*(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^*(s))} V^*(s') \\ &\leq \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right] = r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} V^*(s') \end{aligned}$$

Proof of Bellman Optimality

Theorem 1: Bellman Optimality

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

Denote $\hat{\pi}(s) := \arg \max_a Q^*(s, a)$, we will prove $V^{\hat{\pi}}(s) = V^*(s), \forall s$

$$\begin{aligned} V^*(s) &= r(s, \pi^*(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^*(s))} V^*(s') \\ &\leq \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right] = r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} V^*(s') \\ &= r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[r(s', \pi^*(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \pi^*(s'))} V^*(s'') \right] \end{aligned}$$

Proof of Bellman Optimality

Theorem 1: Bellman Optimality

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

Denote $\hat{\pi}(s) := \arg \max_a Q^*(s, a)$, we will prove $V^{\hat{\pi}}(s) = V^*(s), \forall s$

$$\begin{aligned} V^*(s) &= r(s, \pi^*(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^*(s))} V^*(s') \\ &\leq \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right] = r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} V^*(s') \\ &= r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[r(s', \pi^*(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \pi^*(s'))} V^*(s'') \right] \\ &\leq r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[r(s', \hat{\pi}(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \hat{\pi}(s'))} V^*(s'') \right] \end{aligned}$$

Proof of Bellman Optimality

Theorem 1: Bellman Optimality

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

Denote $\hat{\pi}(s) := \arg \max_a Q^*(s, a)$, we will prove $V^{\hat{\pi}}(s) = V^*(s), \forall s$

$$\begin{aligned} V^*(s) &= r(s, \pi^*(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^*(s))} V^*(s') \\ &\leq \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right] = r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} V^*(s') \\ &= r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[r(s', \pi^*(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \pi^*(s'))} V^*(s'') \right] \\ &\leq r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[r(s', \hat{\pi}(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \hat{\pi}(s'))} V^*(s'') \right] \\ &\leq r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[r(s', \hat{\pi}(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \hat{\pi}(s'))} \left[r(s'', \hat{\pi}(s'')) + \gamma \mathbb{E}_{s''' \sim P(s'', \hat{\pi}(s''))} V^*(s''') \right] \right] \end{aligned}$$

Proof of Bellman Optimality

Theorem 1: Bellman Optimality

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

Denote $\hat{\pi}(s) := \arg \max_a Q^*(s, a)$, we will prove $V^{\hat{\pi}}(s) = V^*(s), \forall s$

$$\begin{aligned} V^*(s) &= r(s, \pi^*(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^*(s))} V^*(s') \\ &\leq \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right] = r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} V^*(s') \\ &= r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[r(s', \pi^*(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \pi^*(s'))} V^*(s'') \right] \\ &\leq r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[r(s', \hat{\pi}(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \hat{\pi}(s'))} V^*(s'') \right] \\ &\leq r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[r(s', \hat{\pi}(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \hat{\pi}(s'))} \left[r(s'', \hat{\pi}(s'')) + \gamma \mathbb{E}_{s''' \sim P(s'', \hat{\pi}(s''))} V^*(s''') \right] \right] \\ &\leq \mathbb{E} \left[r(s, \hat{\pi}(s)) + \gamma r(s', \hat{\pi}(s')) + \dots \right] = V^{\hat{\pi}}(s) \end{aligned}$$

Proof of Bellman Optimality

Theorem 1: Bellman Optimality

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

Denote $\hat{\pi}(s) := \arg \max_a Q^*(s, a)$, we just proved $V^{\hat{\pi}}(s) = V^*(s), \forall s$

Proof of Bellman Optimality

Theorem 1: Bellman Optimality

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

Denote $\hat{\pi}(s) := \arg \max_a Q^*(s, a)$, we just proved $V^{\hat{\pi}}(s) = V^*(s), \forall s$

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

Proof of Bellman Optimality

Theorem 1: Bellman Optimality

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

Denote $\hat{\pi}(s) := \arg \max_a Q^*(s, a)$, we just proved $V^{\hat{\pi}}(s) = V^*(s), \forall s$

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

Q: why?

Summary so far:

Bellman Optimality and DP

Theorem 1: Bellman Optimality

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

Summary so far:

Bellman Optimality and DP

Theorem 1: Bellman Optimality

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

Next:

Any function $V(s)$ that satisfies Bellman Optimality, MUST be equal to V^*

Bellman Optimality

Theorem 2:

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

Bellman Optimality

Theorem 2:

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

Bellman Opt allows us to focus on just one step,

i.e., to check if $V = V^*$,

we only need to check if $\left| V(s) - \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V(s') \right] \right| = 0, \forall s,$

Bellman Optimality

Theorem 2:

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

Bellman Optimality

Theorem 2:

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

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

Bellman Optimality

Theorem 2:

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

$$\begin{aligned} |V(s) - V^*(s)| &= \left| \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s')) - \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s')) \right| \\ &\leq \max_a \left| (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s')) - (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s')) \right| \end{aligned}$$

Bellman Optimality

Theorem 2:

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

$$\begin{aligned} |V(s) - V^*(s)| &= \left| \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s')) - \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s')) \right| \\ &\leq \max_a \left| (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s')) - (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s')) \right| \\ &\leq \max_a \gamma \mathbb{E}_{s' \sim P(s, a)} |V(s') - V^*(s')| \end{aligned}$$

Bellman Optimality

Theorem 2:

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

$$\begin{aligned} |V(s) - V^*(s)| &= \left| \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s')) - \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s')) \right| \\ &\leq \max_a \left| (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s')) - (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s')) \right| \\ &\leq \max_a \gamma \mathbb{E}_{s' \sim P(s, a)} \left| V(s') - V^*(s') \right| \\ &\leq \max_a \gamma \mathbb{E}_{s' \sim P(s, a)} \left(\max_{a'} \gamma \mathbb{E}_{s'' \sim P(s', a')} \left| V(s'') - V^*(s'') \right| \right) \end{aligned}$$

Bellman Optimality

Theorem 2:

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

$$\begin{aligned} |V(s) - V^*(s)| &= \left| \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s')) - \max_a (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s')) \right| \\ &\leq \max_a \left| (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s')) - (r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s')) \right| \\ &\leq \max_a \gamma \mathbb{E}_{s' \sim P(s, a)} \left| V(s') - V^*(s') \right| \\ &\leq \max_a \gamma \mathbb{E}_{s' \sim P(s, a)} \left(\max_{a'} \gamma \mathbb{E}_{s'' \sim P(s', a')} \left| V(s'') - V^*(s'') \right| \right) \\ &\leq \max_{a_1, a_2, \dots, a_{k-1}} \gamma^k \mathbb{E}_{s_k} \left| V(s_k) - V^*(s_k) \right| \end{aligned}$$

Summary so far:

1. V^* satisfies Bellman Optimality:

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

2. Any V that satisfies Bellman Optimality, i.e., $V(s) = \max_a \left[r(s, a) + \mathbb{E}_{s' \sim P(s, a)} V(s') \right], \forall s,$

MUST be that $V(s) = V^*(s),$ for all s

Outline

✓ 1. We have A^S many policies, which one is the optimal policy π^\star ?

✓ 2. Key property of the optimal policy π^\star : **Bellman Optimality**

3. State-action distributions

What's the probability of π visiting a particular state s ?

Discounted State (action) Occupancy Measures

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

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

Discounted State (action) Occupancy Measures

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

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

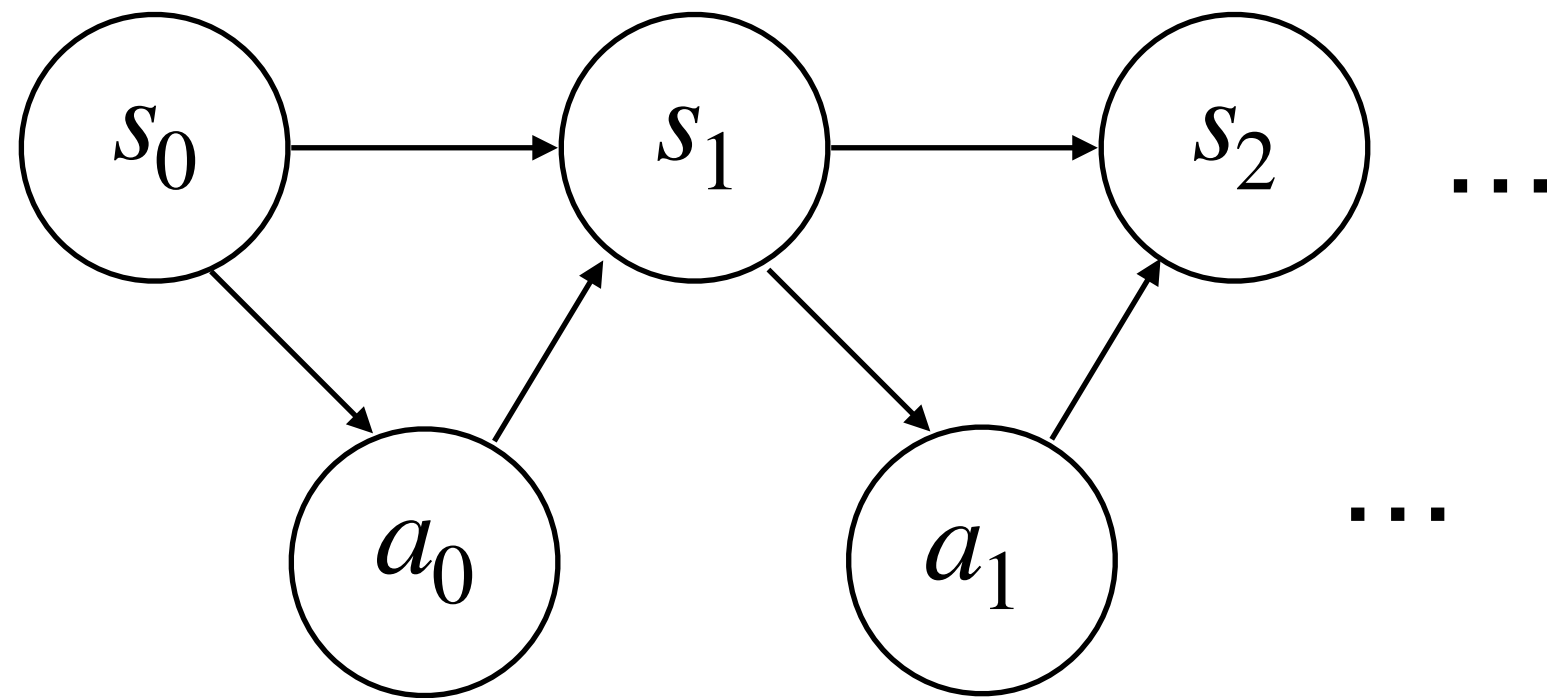
Let's write π as a delta distribution, i.e., $\pi(a | s) = \begin{cases} 1, & a = \pi(s), \\ 0, & \text{else} \end{cases}$

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., $\pi(a | s) = \begin{cases} 1, & a = \pi(s), \\ 0, & \text{else} \end{cases}$

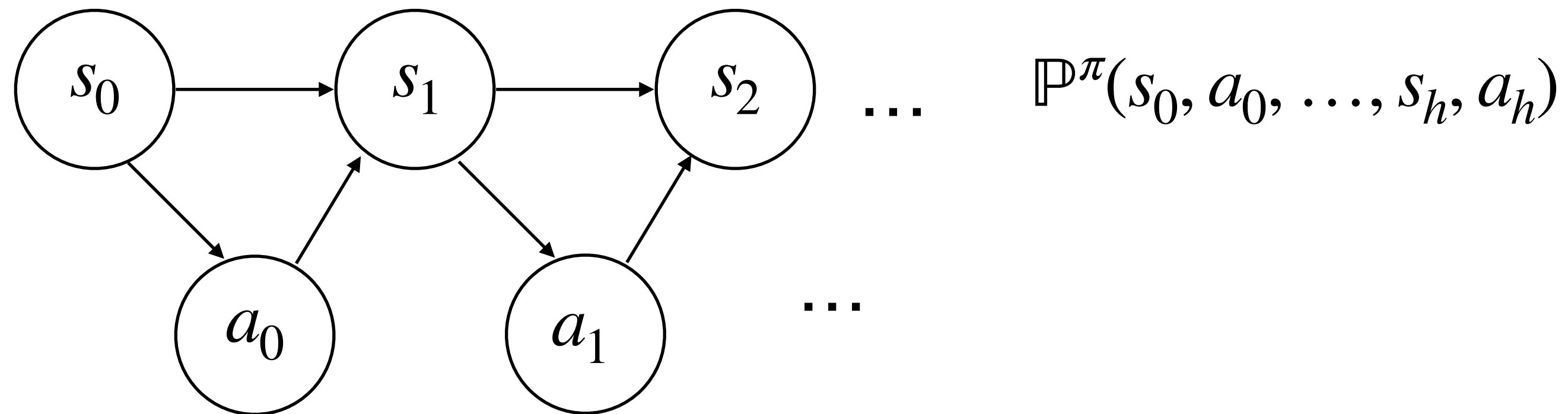


Discounted State (action) Occupancy Measures

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., $\pi(a | s) = \begin{cases} 1, & a = \pi(s), \\ 0, & \text{else} \end{cases}$

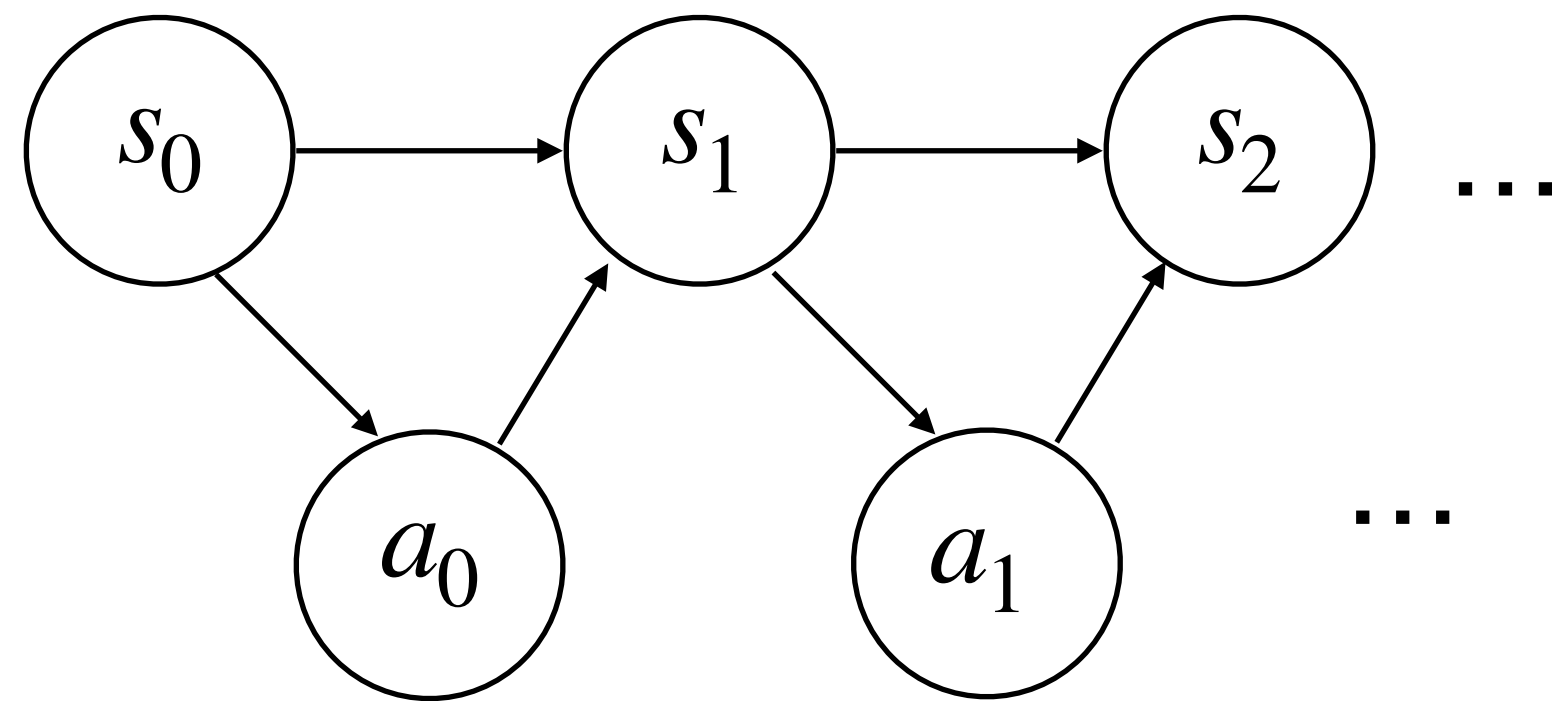


Discounted State (action) Occupancy Measures

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

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

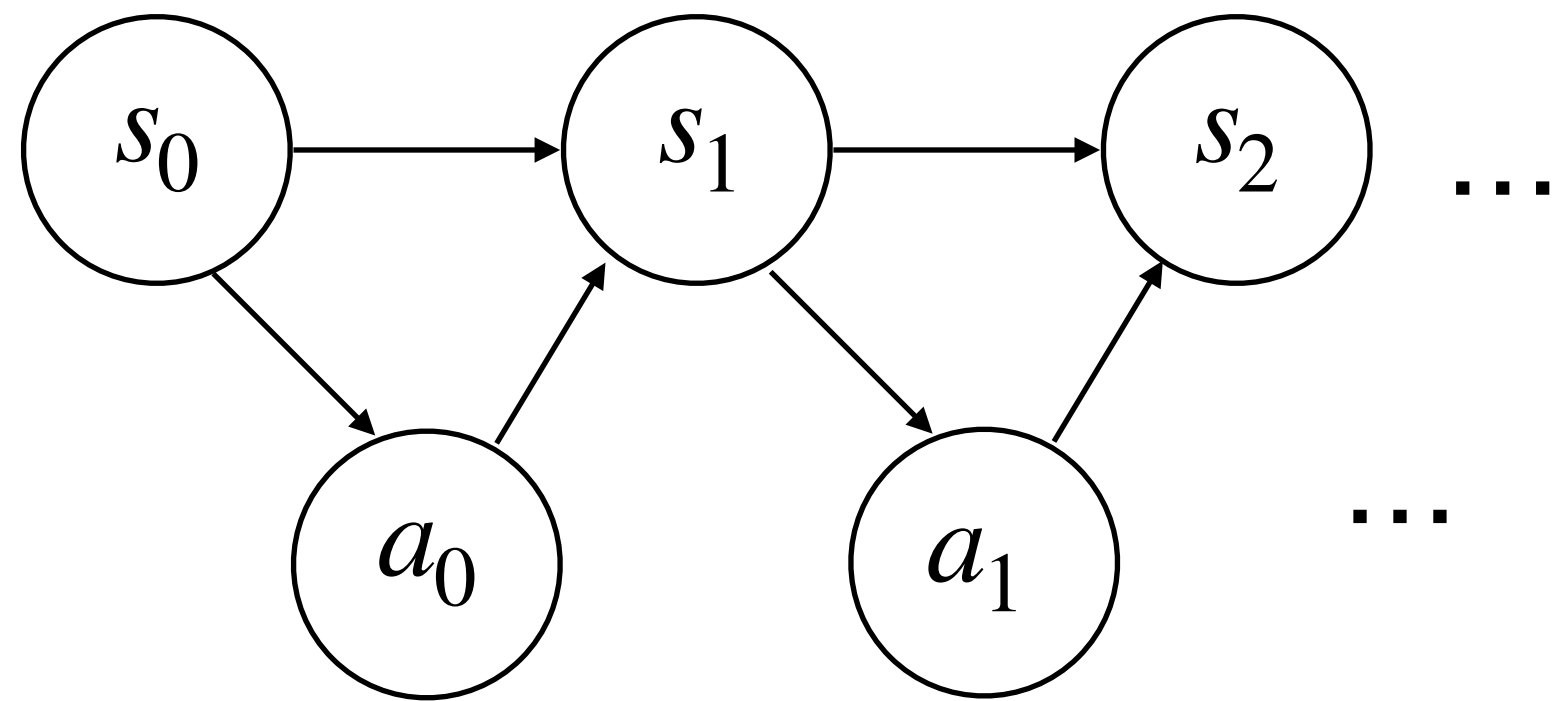
Let's write π as a delta distribution, i.e., $\pi(a | s) = \begin{cases} 1, & a = \pi(s), \\ 0, & \text{else} \end{cases}$



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

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

State-action distribution at time step h

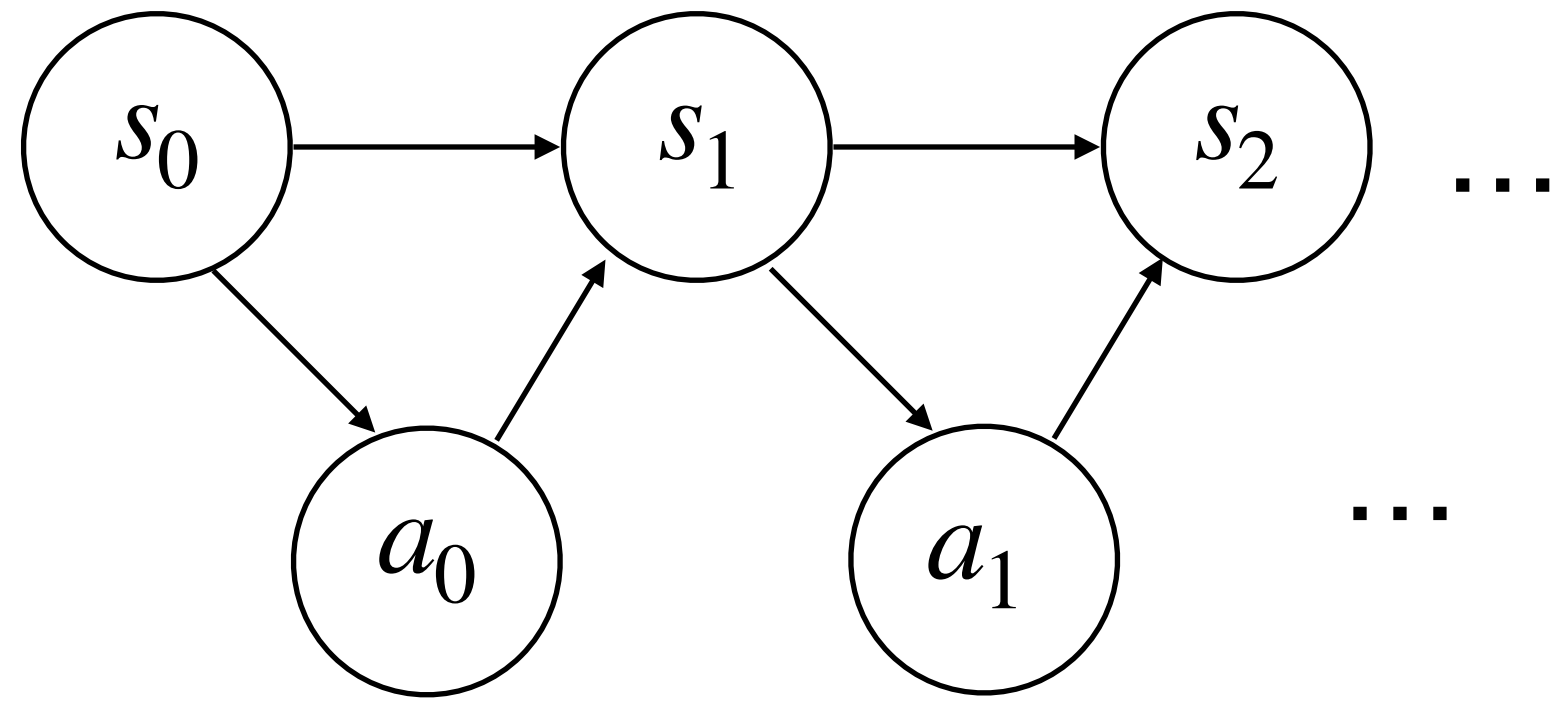


$$\mathbb{P}^{\pi}(s_0, a_0, \dots, s_h, a_h)$$

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

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

State-action distribution at time step h



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

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

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

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

Discounted Average State-action distribution

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

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

Discounted Average State-action distribution

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

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

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

Discounted Average State-action distribution

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

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

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

Can you show that this is a valid distribution?

Discounted Average State-action distribution

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

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

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

Can you show that this is a valid distribution?

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

Discounted Average State-action distribution

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

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

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

Can you show that this is a valid distribution?

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

Can you show the above is true?

Discounted Average State-action distribution

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

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

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

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 for today:

1. π^\star **dominates** other policies, i.e., $V^\star(s) \geq V^\pi(s), \forall s, \pi$
2. Key property of the optimal policy π^\star : **Bellman Optimality**
(BE and B-Opt allow us to focus on one step)
3. State-action distribution: $d_{s_0}^\pi(s, a)$

Summary for today:

1. π^\star **dominates** other policies, i.e., $V^\star(s) \geq V^\pi(s), \forall s, \pi$
2. Key property of the optimal policy π^\star : **Bellman Optimality**
(BE and B-Opt allow us to focus on one step)
3. State-action distribution: $d_{s_0}^\pi(s, a)$

RL is notation heavy! But we will see these over and over again during the semester.

