

Dynamic Programming & Infinite Horizons

Lucas Janson

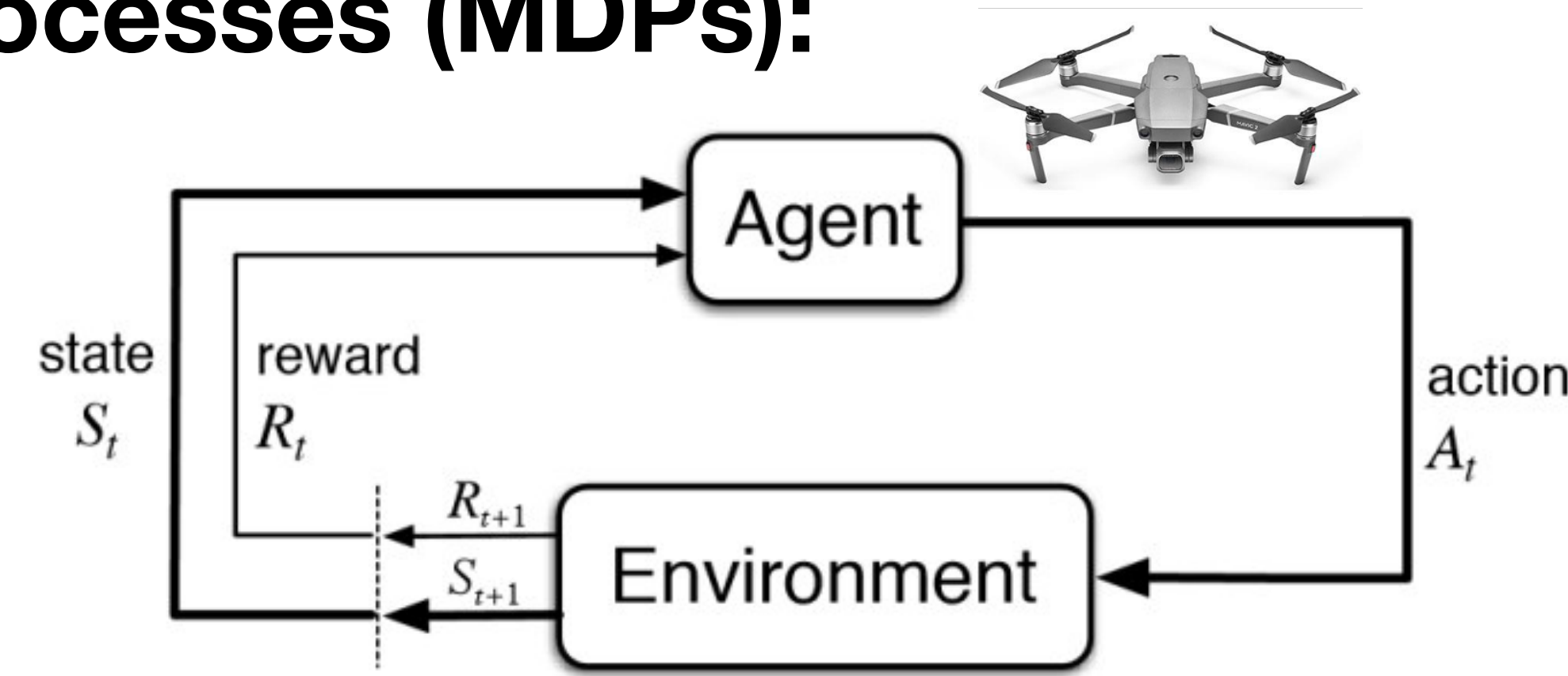
**CS/Stat 184(0): Introduction to Reinforcement Learning
Fall 2024**

Today

- Recap
- Optimality
- The Bellman Equations & Dynamic Programming
- Infinite Horizons

Finite Horizon Markov Decision Processes (MDPs):

- An MDP: $\mathcal{M} = \{\mu, S, A, P, r, H\}$
 - μ is a distribution over initial states
(sometimes we assume we start a given state s_0)
 - S a set of states
 - A a set of actions
 - $P : S \times A \mapsto \Delta(S)$ specifies the dynamics model,
i.e. $P(s' | s, a)$ is the probability of transitioning to s' from state s via action a
 - $r : S \times A \rightarrow [0,1]$
 - For now, let's assume this is a deterministic function
 - (sometimes we use a cost $c : S \times A \rightarrow [0,1]$)
 - A time horizon $H \in \mathbb{N}$



Policy Evaluation = Computing Value function and/or Q function

We evaluate policies via quantities that allow us to reason about the policy's long-term effect:

- **Value function** $V_h^\pi(s) = \mathbb{E} \left[\sum_{t=h}^{H-1} r(s_t, a_t) \mid s_h = s \right]$

- **Q function** $Q_h^\pi(s, a) = \mathbb{E} \left[\sum_{t=h}^{H-1} r(s_t, a_t) \mid (s_h, a_h) = (s, a) \right]$

- For deterministic policy π , **Bellman consistency**:

- $V_h^\pi(s) = r(s, \pi_h(s)) + \mathbb{E}_{s' \sim P(\cdot | s, \pi_h(s))} [V_{h+1}^\pi(s')]$

- $Q_h^\pi(s, a) = r(s, a) + \mathbb{E}_{s' \sim P(\cdot | s, a)} [V_{h+1}^\pi(s')]$

- **DP:**

- Initialize: $V_H^\pi(s) = 0, \forall s \in S$

- For $h = H - 1, \dots, 0$, set:

$$V_h^\pi(s) = r(s, \pi_h(s)) + \mathbb{E}_{s' \sim P(\cdot | s, \pi_h(s))} [V_{h+1}^\pi(s')], \forall s \in S$$

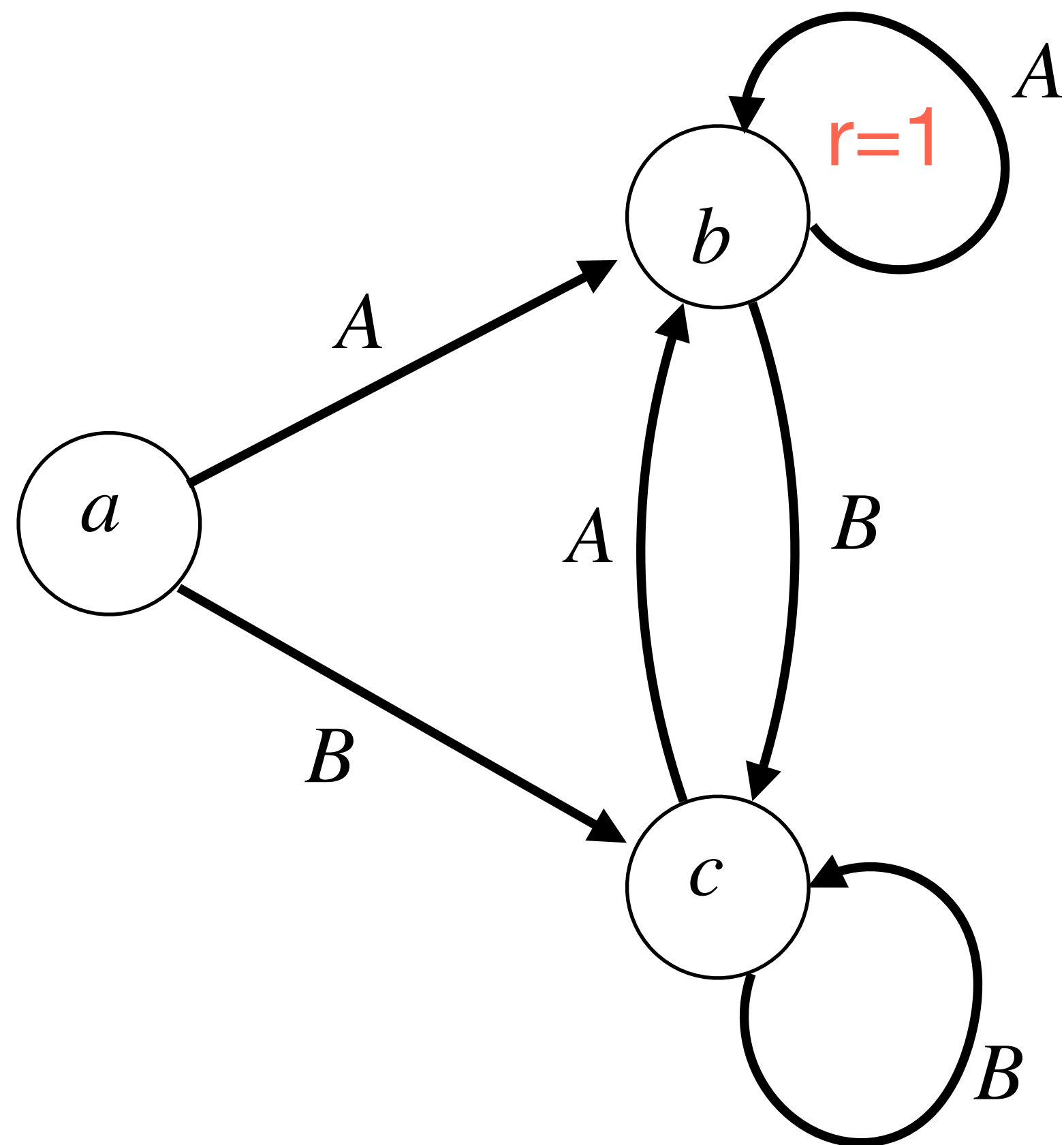
Today



- Recap
- Optimality
- The Bellman Equations & Dynamic Programming
- Infinite Horizons

Example of Optimal Policy π^\star

Consider the following **deterministic** MDP w/ 3 states & 2 actions, with $H = 3$



- What's the optimal policy?

$$\pi_h^\star(s) = A, \quad \forall s, h$$

- What is optimal value function, $V^{\pi^\star} = V^\star$?

$$V_2^\star(a) = 0, \quad V_2^\star(b) = 1, \quad V_2^\star(c) = 0$$

$$V_1^\star(a) = 1, \quad V_1^\star(b) = 2, \quad V_1^\star(c) = 1$$

$$V_0^\star(a) = 2, \quad V_0^\star(b) = 3, \quad V_0^\star(c) = 2$$

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

How do we compute π^\star and V^\star ?

- Naively, we could compute the value of all policies and take the best one.
- Suppose $|S|$ states, $|A|$ actions, and horizon H .
How many different policies there are?
- Can we do better?

Properties of an Optimal Policy π^\star

- Let Π be the set of all time dependent, history dependent, stochastic policies.
- **Theorem:** Every finite horizon MDP has a **deterministic, history-independent** optimal policy, that **dominates all other policies, everywhere**.
 - i.e. there exists a deterministic policy $\pi^\star := \{\pi_0^\star, \pi_1^\star, \dots, \pi_{H-1}^\star\}$, $\pi_h^\star : S \mapsto A$ such that

$$V_h^{\pi^\star}(s) \geq V_h^\pi(s) \quad \forall s, h, \forall \pi \in \Pi$$

- \implies we can write: $V_h^\star = V_h^{\pi^\star}$ and $Q_h^\star = Q_h^{\pi^\star}$.
- $\implies \pi^\star$ doesn't depend on the initial state distribution μ .

What's the Proof Intuition?

- **Theorem:** Every finite horizon MDP has a **deterministic, history-independent** optimal policy, that **dominates all other policies, everywhere.**
- What's the Proof Intuition?
 - “Only the state matters”: how got here doesn't matter to where we go next, conditioned on the action.
 - This explains both determinism and history-independence
- Caveat: some legitimate reward functions are not additive/linear (so, naively, not an MDP). (But, RL is general: think about redefining the state so you can do these.)

Today

- ✓ • Recap
- ✓ • Optimality
- The Bellman Equations & Dynamic Programming
- Infinite Horizons

The Bellman Equations

- A function $V = \{V_0, \dots, V_{H-1}\}$, $V_h : S \rightarrow R$ satisfies the **Bellman equations** if

$$V_h(s) = \max_a \left\{ r(s, a) + \mathbb{E}_{s' \sim P(\cdot | s, a)} [V_{h+1}(s')] \right\}, \forall s$$

(assume $V_H = 0$).

- **Theorem:**

- V satisfies the Bellman equations **if and only if** $V = V^*$.

- The optimal policy is: $\pi_h^*(s) = \arg \max_a \left\{ r(s, a) + \mathbb{E}_{s' \sim P(\cdot | s, a)} [V_{h+1}^*(s')] \right\}$.

Computation of V^\star with Dynamic Programming

- **Theorem:** the following **Dynamic Programming** algorithm computes π^\star and V^\star
Prf: the Bellman equations directly lead to this backwards induction.

- Initialize: $V_H^\pi(s) = 0 \quad \forall s \in S$
For $t = H - 1, \dots, 0$, set:
 - $V_h^\star(s) = \max_a \left[r(s, a) + \mathbb{E}_{s' \sim P(\cdot | s, a)} [V_{h+1}^\star(s')] \right], \quad \forall s \in S$
 - $\pi_h^\star(s) = \arg \max_a \left[r(s, a) + \mathbb{E}_{s' \sim P(\cdot | s, a)} [V_{h+1}^\star(s')] \right], \quad \forall s \in S$

- What is the per iteration computational complexity of DP?
(assume scalar $+$, $-$, \times , \div are $O(1)$ operations)
- What is the total computational complexity of DP?

Today

- ✓ • Recap
- ✓ • Optimality
- ✓ • The Bellman Equations & Dynamic Programming
- Infinite Horizons

Infinite Horizon MDPs:

- An MDP: $\mathcal{M} = \{\mu, S, A, P, r, \gamma\}$
 - $\mu, S, A, P : S \times A \mapsto \Delta(S)$, $r : S \times A \rightarrow [0,1]$ same as before
 - instead of finite horizon H , we have a **discount factor** $\gamma \in [0,1)$
- **Objective:** find policy π that maximizes our expected, discounted future reward:
$$\max_{\pi} \mathbb{E} \left[r(s_0, a_0) + \gamma r(s_1, a_1) + \gamma^2 r(s_2, a_2) + \dots \dots \mid s_0 \right]$$

The Setting and Our Objective

- Consider a deterministic, **stationary policy** $\pi : S \mapsto A$
 - stationary means not history or time dependent
- **Sampling a trajectory τ on an episode:** for a given policy π
 - Sample an initial state $s_0 \sim \mu$:
 - For $t = 0, 1, 2, \dots, \infty$
 - Take action $a_t = \pi(s_t)$
 - Observe reward $r_t = r(s_t, a_t)$
 - Transition to (and observe) s_{t+1} where $s_{t+1} \sim P(\cdot | s_t, a_t)$
- The infinite trajectory: $\tau = \{s_0, a_0, r_0, s_1, a_1, r_1, \dots, \}$

Value function and Q functions:

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

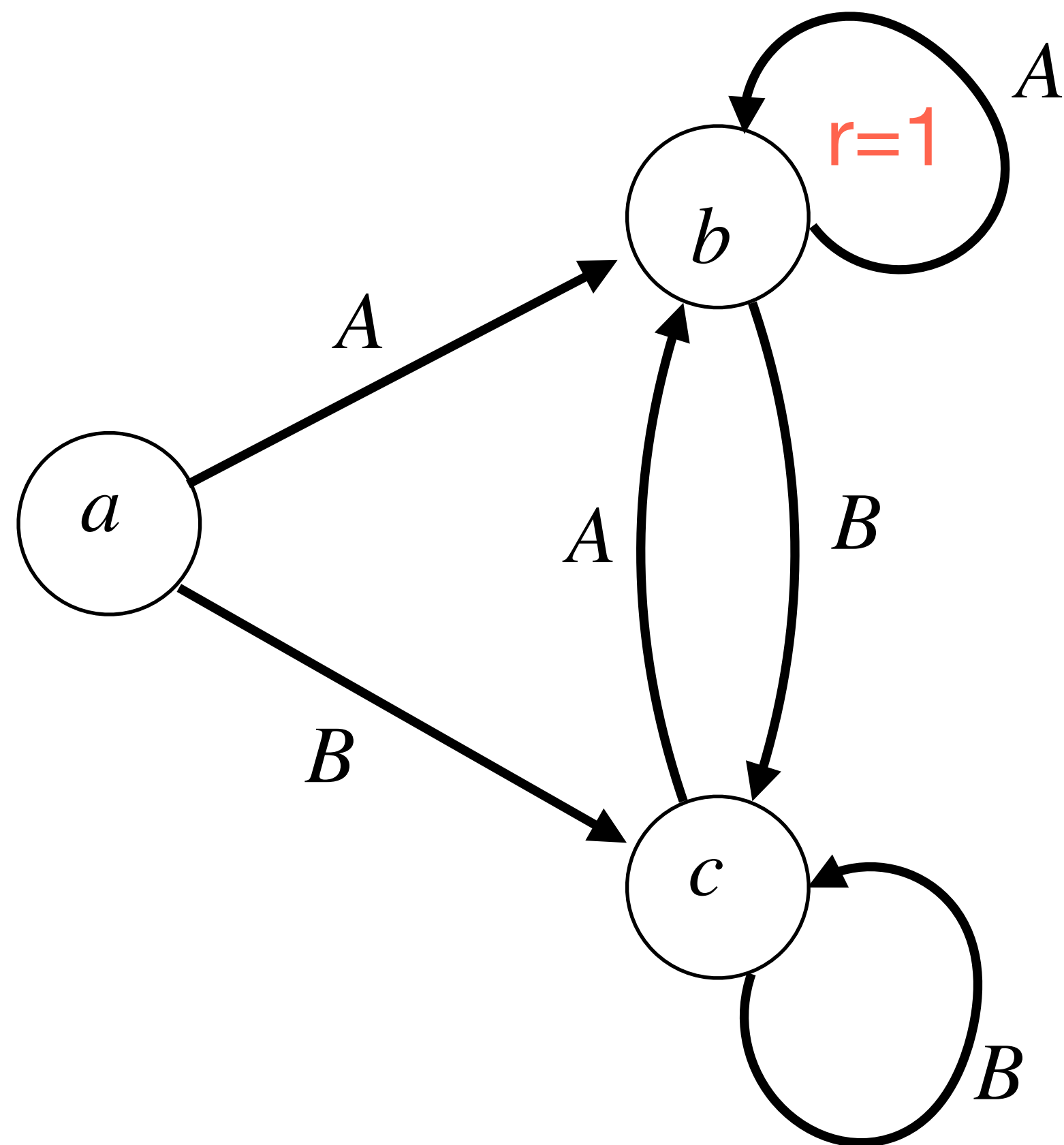
- Value function $V^\pi(s) = \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid s_0 = s \right]$

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

- What are upper and lower bounds on V^π and Q^π

Example of Policy Evaluation (e.g. computing V^π and Q^π)

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



- Consider the policy
 $\pi(a) = B, \pi(b) = A, \pi(c) = A$

- What is V^π ?

$$V^\pi(a) =$$

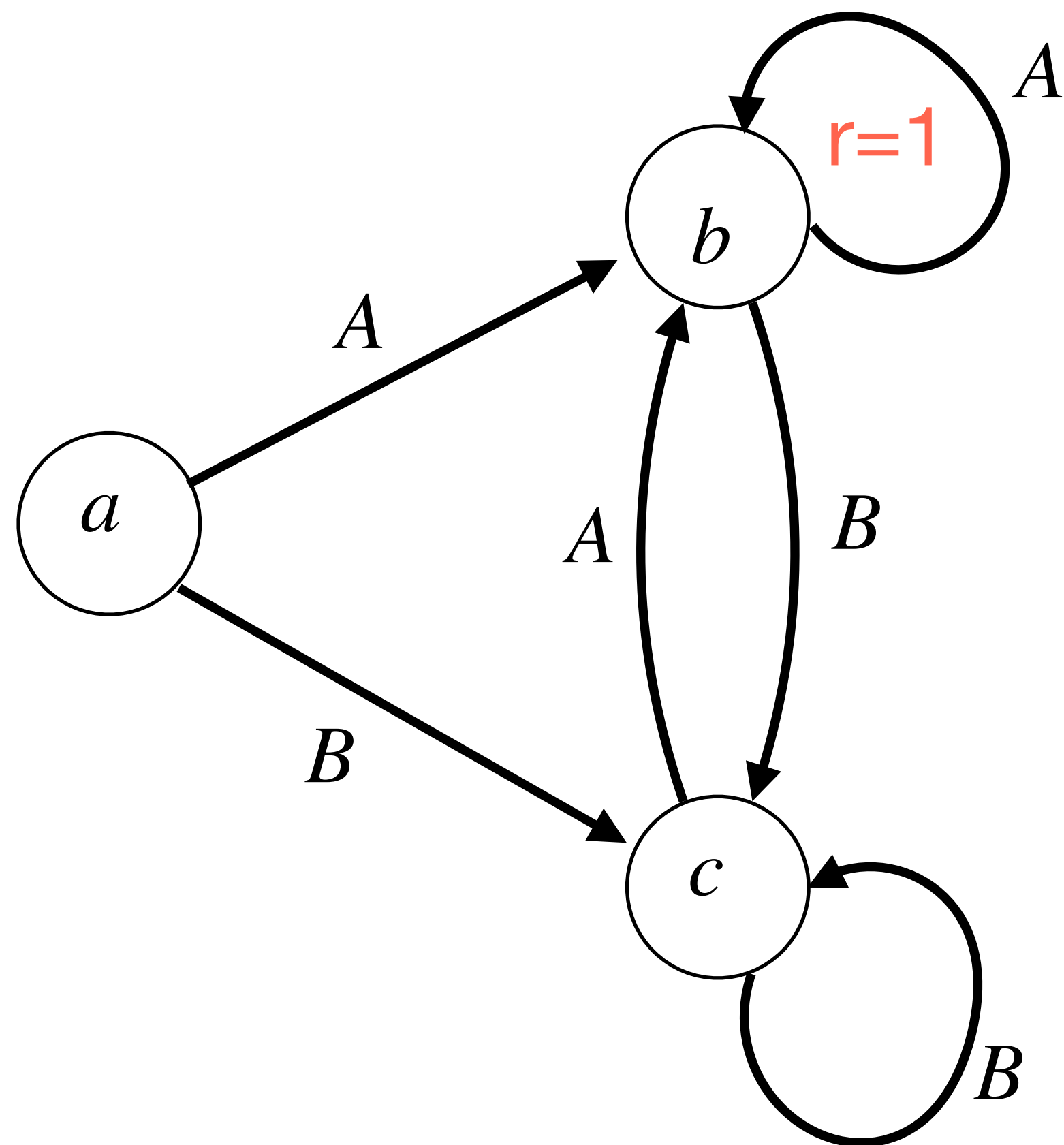
$$V^\pi(b) =$$

$$V^\pi(c) =$$

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

Example of Policy Evaluation (e.g. computing V^π and Q^π)

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



- Consider the policy
 $\pi(a) = B, \pi(b) = A, \pi(c) = A$

- What is V^π ?

$$V^\pi(a) = \gamma^2 / (1 - \gamma)$$

$$V^\pi(b) = 1 / (1 - \gamma)$$

$$V^\pi(c) = \gamma / (1 - \gamma)$$

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

Bellman Consistency (theorem)

- Consider a fixed policy, $\pi : S \mapsto A$.
- By definition, $V^\pi(s) = Q^\pi(s, \pi(s))$
- Bellman consistency conditions:
 - $V^\pi(s) = r(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi(s))} [V^\pi(s')]$
 - $Q^\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [V^\pi(s')]$

(Optional) Proof: Bellman Consistency for V-function:

- By definition and by the “tower” property of conditional expectations:

$$\begin{aligned} V^\pi(s) &= \mathbb{E} \left[r(s_0, a_0) + \gamma r(s_1, a_1) + \gamma^2 r(s_2, a_2) + \dots \mid s_0 = s \right] \\ &= \mathbb{E} \left[r(s_0, a_0) + \mathbb{E} \left[\gamma r(s_1, a_1) + \gamma^2 r(s_2, a_2) + \dots \mid s_0 = s, a_0, s_1 \right] \mid s_0 = s \right] \end{aligned}$$

- By the Markov property:

$$\begin{aligned} &= \mathbb{E} \left[r(s_0, a_0) + \gamma \mathbb{E} \left[r(s_1, a_1) + \gamma r(s_2, a_2) + \dots \mid s_1 \right] \mid s_0 = s \right] \\ &= \mathbb{E} \left[r(s_0, a_0) + \gamma V^\pi(s_1) \mid s_h = s \right] \\ &= r(s, \pi(s)) + \gamma \sum_{s'} P(s' \mid s, \pi(s)) V^\pi(s') \end{aligned}$$

Computation of V^π

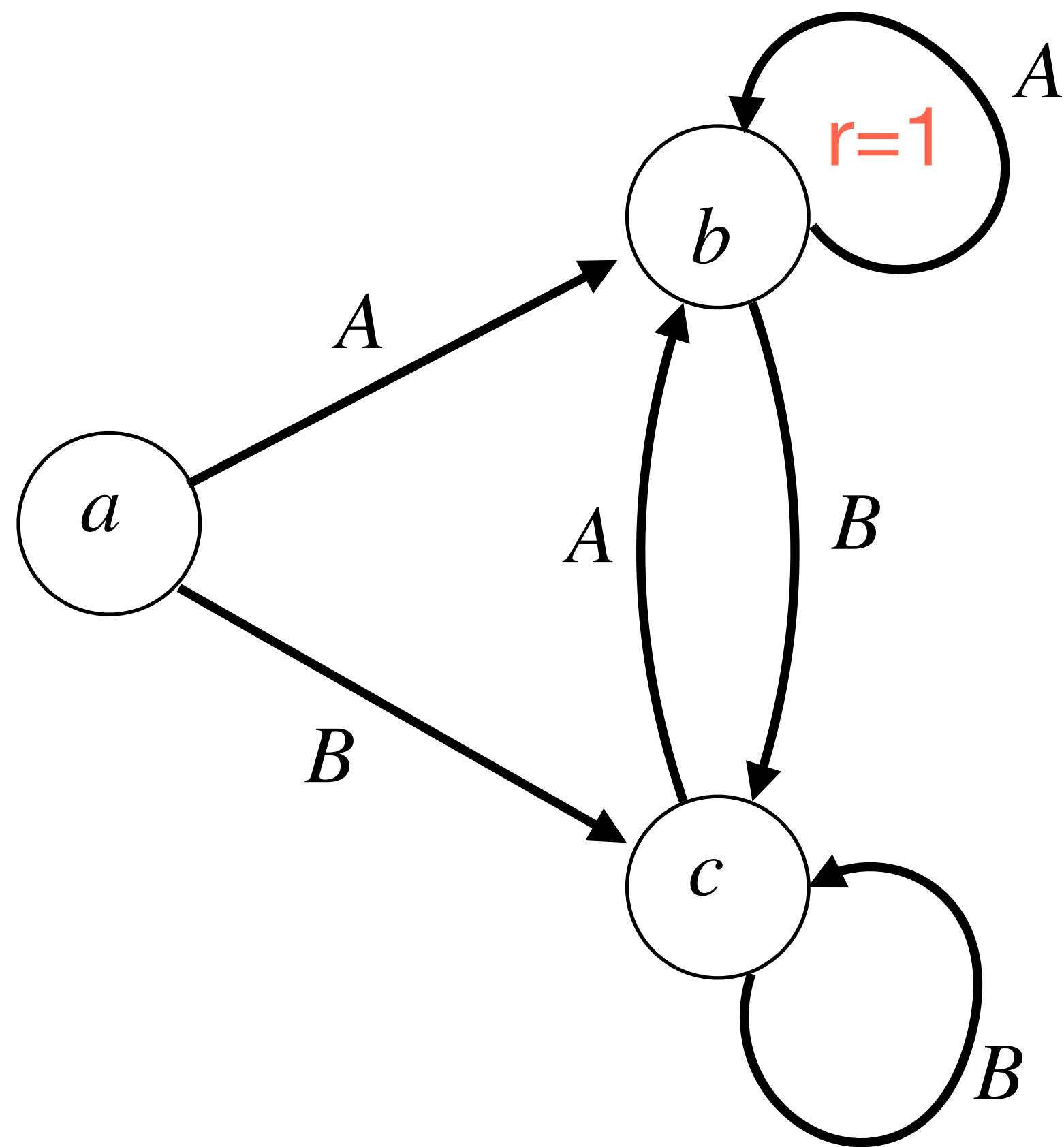
- For a fixed policy, $\pi : S \mapsto A$, let's compute its V (and Q) value functions.
- We have the Bellman consistency conditions, for a given policy π

$$V^\pi(s) = r(s, \pi(s)) + \gamma \sum_{s'} P(s' | s, \pi(s)) V^\pi(s')$$

- How do we use this to find a solution?
- What is the time complexity?

Example of Optimal Policy π^* , discounted case

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



- What's the optimal policy?

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

- What is optimal value function, $V^{\pi^*} = V^*$?

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

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

How do we compute π^\star and V^\star ?

- Naively, we could compute the value of all policies and take the best one.
- Suppose $|S|$ states, $|A|$ actions.
How many different stationary policies are there?

Properties of an Optimal Policy π^\star

- **Theorem:** Every infinite horizon MDP has a **stationary, history independent, deterministic** optimal policy, that **dominates all other policies, everywhere.**

- i.e. there exists a policy $\pi^\star : S \mapsto A$ such that

$$V^{\pi^\star}(s) \geq V^\pi(s) \quad \forall s, \forall \pi \in \Pi$$

(again Π is the set of all time dependent, history dependent, stochastic policies)

- \implies we can write: $V^\star = V^{\pi^\star}$ and $Q^\star = Q^{\pi^\star}$.

Summary:

- **Dynamic Programming lets us efficiently compute optimal policies.**
 - We remember the results on “sub-problems”
 - Optimal policies are history independent.
- Discounted infinite horizon MDP analogous to finite-horizon case

Attendance:

bit.ly/3RcTC9T



Feedback:

bit.ly/3RHtlxy

