# Infinite Horizon MDPs: Value and Policy Iteration

## Lucas Janson and Sham Kakade

**CS/Stat 184: Introduction to Reinforcement Learning**
**Fall 2023**

# Today

✓ • Recap

• Infinite Horizon MDPs

    • Optimality & the Bellman Equations

    • Value Iteration

    • Policy Iteration

- HW 1 is posted.

- HW1 is long. Please start early.

# Recap

# Infinite Horizon MDPs:

- An MDP: $\mathcal{M} = \{\mu, S, A, P, r, \gamma\}$
  - $\mu$, $S$, $A$, $P : S \times A \mapsto \Delta(S)$, $r : S \times A \to [0,1]$ same as before
  - instead of finite horizon $H$, we have a discount factor $\gamma \in [0,1)$

- Objective: find policy $\pi$ that maximizes our expected, discounted future reward:
$$\max_{\pi} \mathbb{E}\left[r(s_0, a_0) + \gamma r(s_1, a_1) + \gamma^2 r(s_2, a_2) + \ldots\ldots \,\Big|\, \pi\right]$$

# Value function and Q functions:

# Value function and Q functions:

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

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

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

# Value function and Q functions:

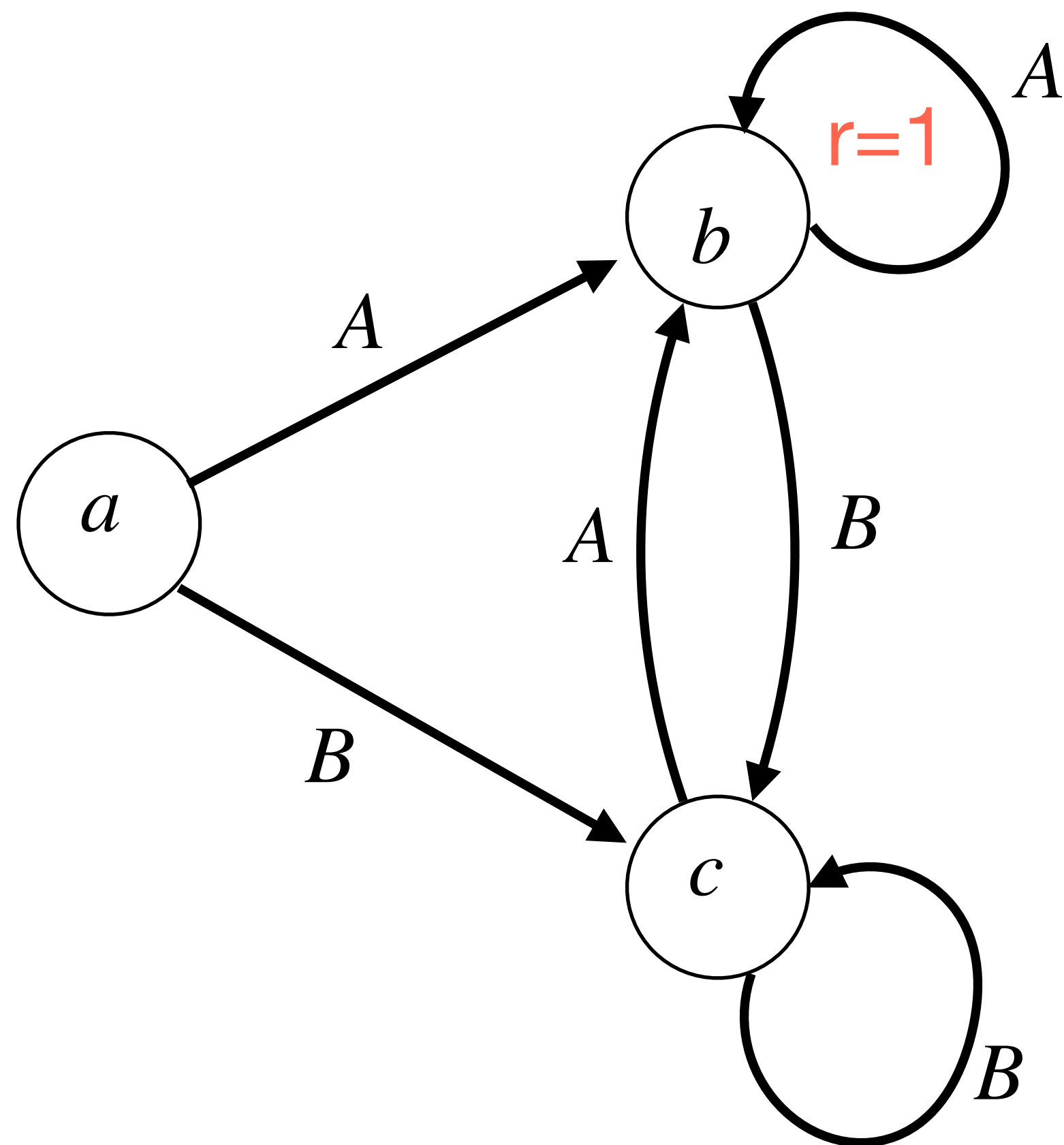- Quantities that allow us to reason about the policy's long-term effect:

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

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

- What are upper and lower bounds on $V^{\pi}$ and $Q^{\pi}$?
$$0 \leq V^{\pi}(s), Q^{\pi}(s, a) \leq 1/(1 - \gamma)$$

# Example of Policy Evaluation (e.g. computing $V^\pi$ and $Q^\pi$)

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



- Consider the policy
  $\pi(a) = B, \pi(b) = A, \pi(c) = A$
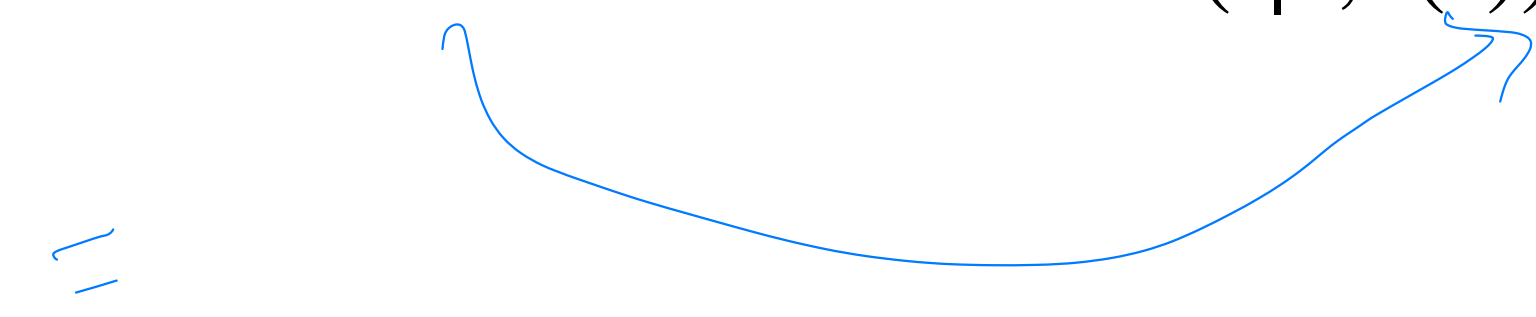- What is $V^\pi$?
  $V^\pi(a) = \gamma^2/(1 - \gamma)$

  $0 + 0 + \gamma^2 + \gamma^3 \ldots$

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

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

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

6

# Bellman Consistency (theorem)

- Consider a fixed policy, $\pi : S \mapsto A$.

- By definition, $V^{\pi}(s) = Q^{\pi}(s, \pi(s))$

- Bellman consistency conditions:

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

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

$r(s, a, s')$

# Computation of $V^\pi$

# Computation of $V^\pi$

- 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 $\pi$

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

- How do we use this to find a solution?
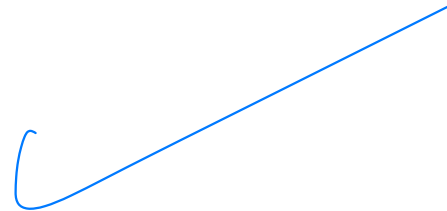
  Solve Linear.

- What is the time complexity?

  $O(|S|^3)$

# Computation of $V^\pi$

- 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 $\pi$

$$V^\pi(s) = r(s, \pi(s)) + \gamma \sum_{s'} P(s' \,|\, s, \pi(s)) V^\pi(s')$$
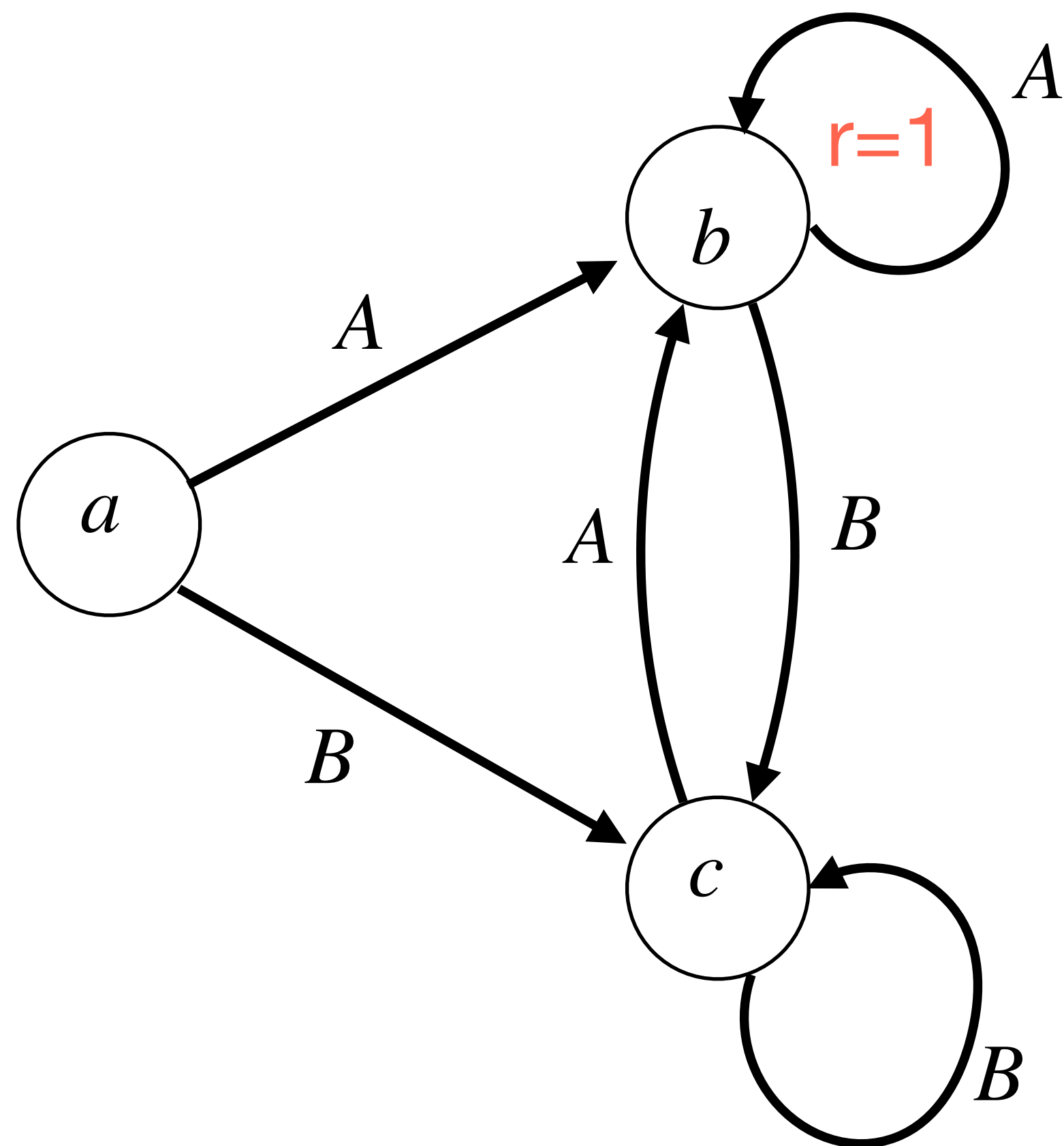
- How do we use this to find a solution?

- What is the time complexity?

- Do you see how to write this with matrix algebra?

# Let's use Bellman Consistency for computing $V^\pi$

$\pi[a] = B$  $\pi(b), \pi(c) = A$

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



$\dfrac{\gamma^2}{1-\gamma}$

$V^\pi(a) = 0 + \gamma \cdot V^\pi(c)$

$V^\pi(b) = 1 + \gamma \cdot V^\pi(c)$ $\Rightarrow V^\pi(b)$

$= \dfrac{1}{1-\gamma}$

$V^\pi(c) = 0 + \gamma V^\pi(b)$

$\dfrac{\gamma}{1-\gamma}$

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

9

# Properties of an Optimal Policy $\pi^\star$

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

  - i.e. there exists a policy $\pi^\star : S \mapsto A$ such that
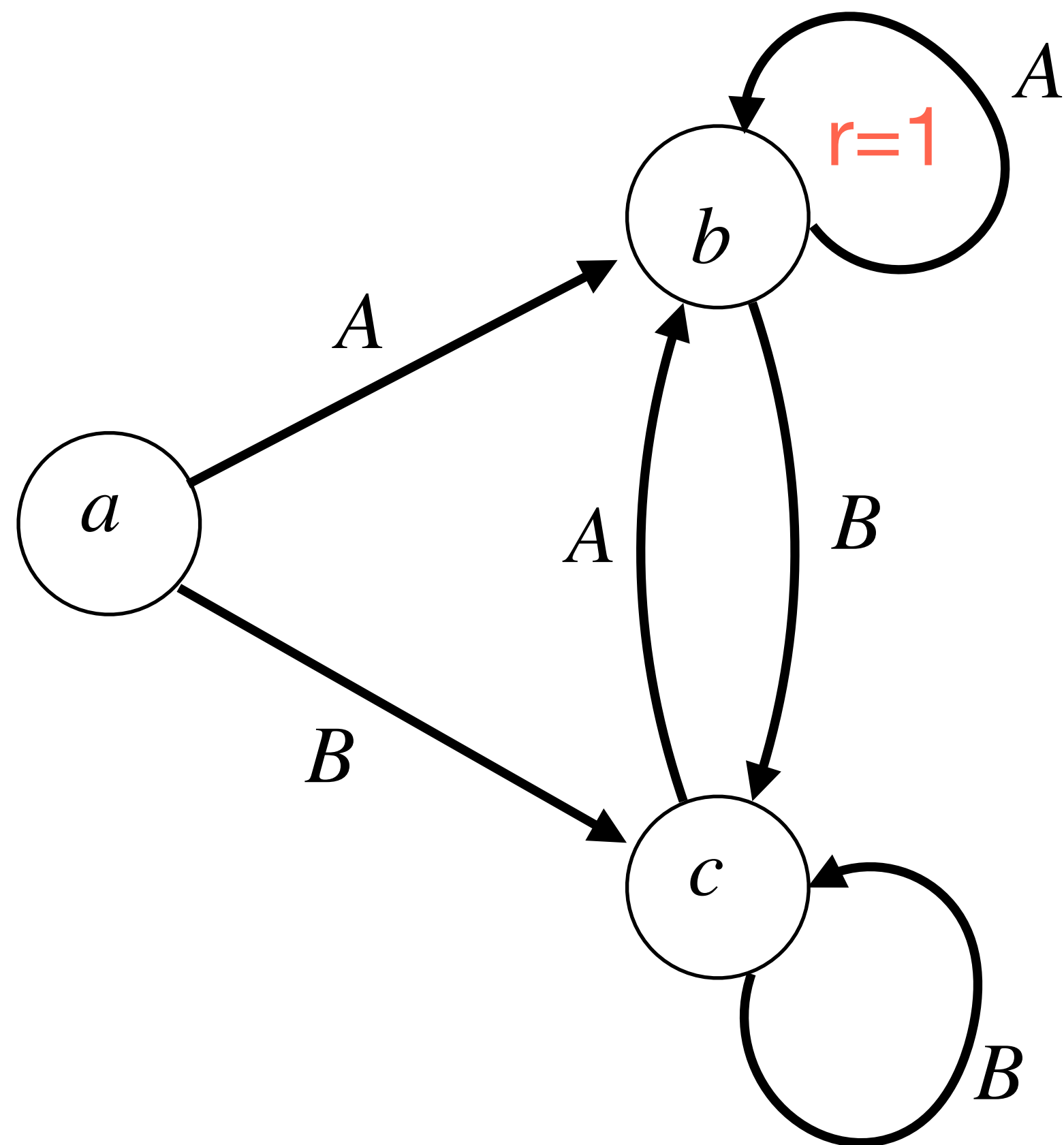    $$V^{\pi^\star}(s) \geq V^\pi(s) \quad \forall s, \ \forall \pi \in \Pi$$

    (again $\Pi$ is the set of all time dependent, history dependent, stochastic policies)

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

# Example of Optimal Policy $\pi^\star$, discounted case

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



- What's the optimal policy?
$\pi^\star(s) = A, \forall s$

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

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

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

# How do we compute $\pi^\star$ 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 polices are there?  $|A|^{|S|}$

# Today

- Recap

- Infinite Horizon MDPs

  ✓ - Optimality & the Bellman Equations

  - Value Iteration

  - Policy Iteration

# The Bellman Equations

# The Bellman Equations

- A function $V : S \to R$ satisfies the Bellman equations if

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

# The Bellman Equations

- A function $V : S \to R$ satisfies the Bellman equations if

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

- **Theorem:**

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

# The Bellman Equations

- A function $V : S \to R$ satisfies the Bellman equations if

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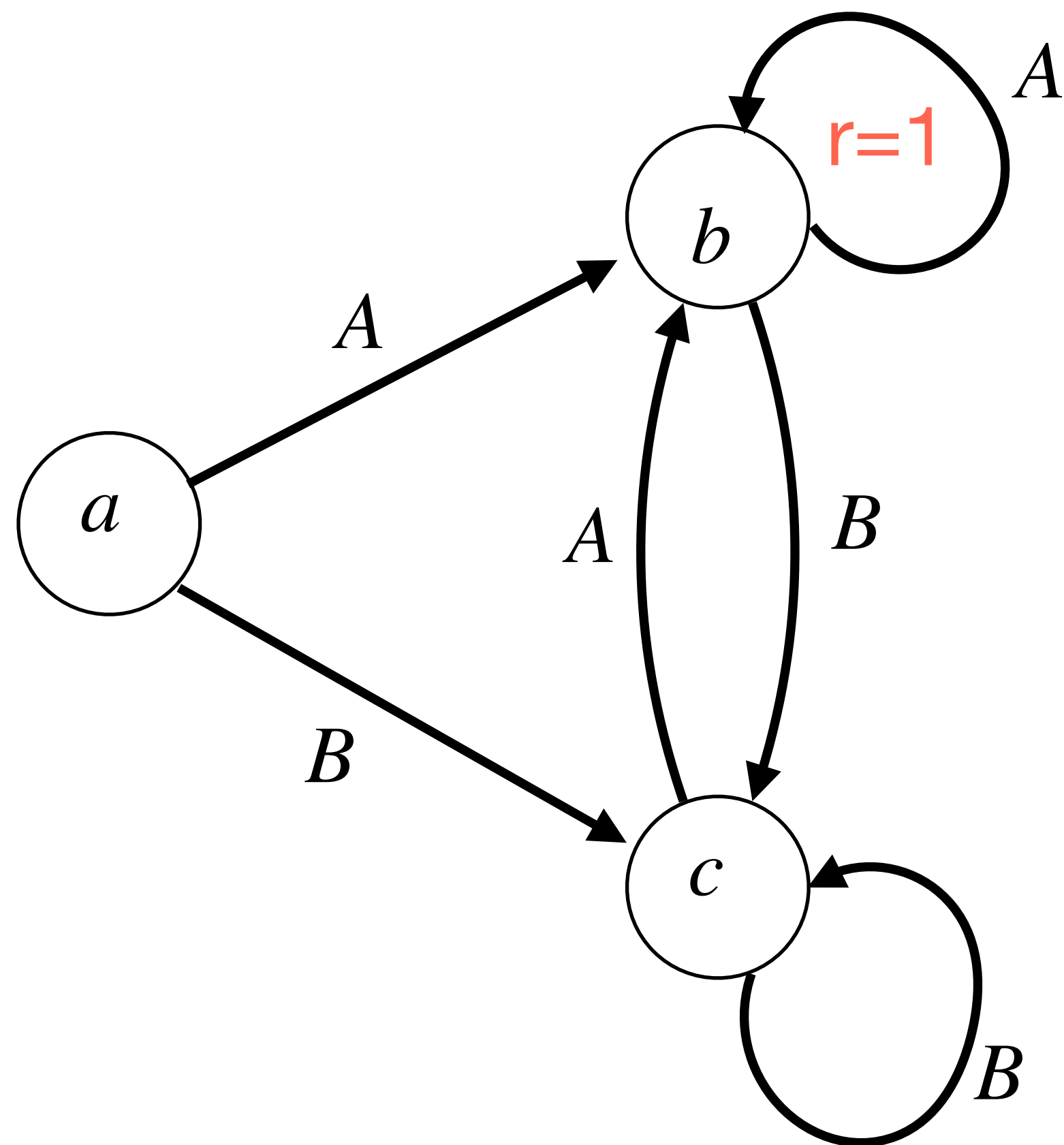
- **Theorem:**

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

  - The optimal policy is: $\pi^\star(s) = \arg\max_a \left\{ r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[ V^\star(s') \right] \right\}.$

# Exercise: use the BE to the purported $\pi^\star$ is optimal

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



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

# Exercise: use the BE to the purported $\pi^\star$ is optimal

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



- What's the optimal policy?
$\pi^\star(s) = A, \forall s$

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

# Exercise: use the BE to the purported $\pi^\star$ is optimal

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



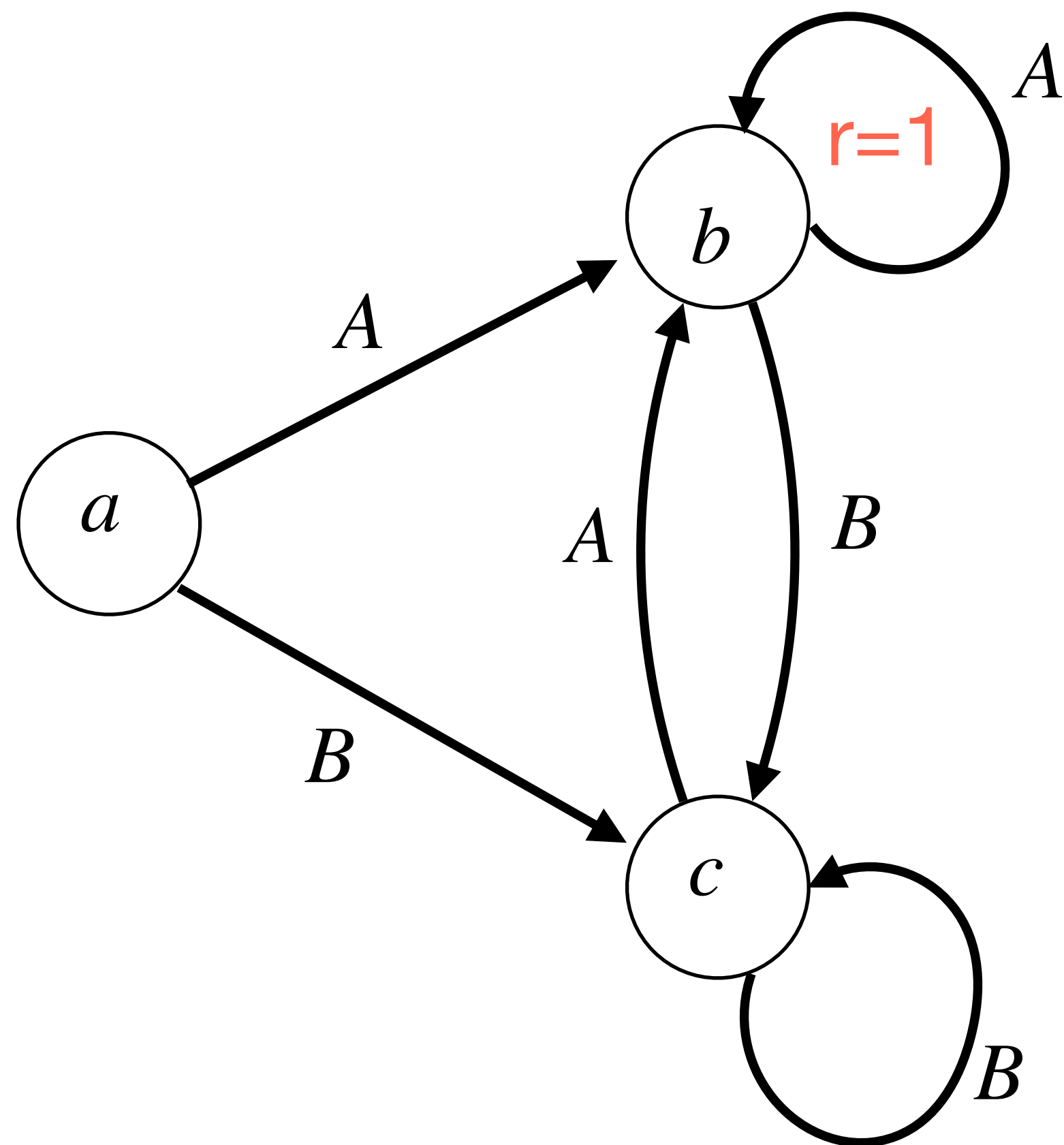- What's the optimal policy?
  $\pi^\star(s) = A, \forall s$

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

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

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

15

# Today

- Recap

- Infinite Horizon MDPs

    - Optimality & the Bellman Equations

    ✓ Value Iteration

    - Policy Iteration

# Detour: fix-point solution

# Detour: fix-point solution

- Suppose we want to find an $x^\star$ s.t. $x^\star = f(x^\star), \quad f : [a, b] \mapsto [a, b]$

# Detour: fix-point solution

- Suppose we want to find an $x^\star$ s.t. $x^\star = f(x^\star), \quad f : [a, b] \mapsto [a, b]$

- A naive approach to find $x^\star$ :

# Detour: fix-point solution

- Suppose we want to find an $x^\star$ s.t. $x^\star = f(x^\star)$, $\quad f : [a, b] \mapsto [a, b]$

- A naive approach to find $x^\star$ :

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

# Detour: fix-point solution

- Suppose we want to find an $x^\star$ s.t. $x^\star = f(x^\star), \quad f : [a, b] \mapsto [a, b]$

- A naive approach to find $x^\star$ :

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

- Suppose $f$ is a contraction mapping: $\forall x, x', \ |f(x) - f(x')| \leq \gamma |x - x'|$, for $\gamma \in [0, 1)$. Then it converges, i.e. $x^t \to x^\star$, as $t \to \infty$.

# Detour: fix-point solution

- Suppose we want to find an $x^\star$ s.t. $x^\star = f(x^\star), \quad f : [a, b] \mapsto [a, b]$

- A naive approach to find $x^\star$ :

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

- Suppose $f$ is a contraction mapping: $\forall x, x', \ |f(x) - f(x')| \leq \gamma |x - x'|,$ for $\gamma \in [0, 1)$. Then it converges, i.e. $x^t \to x^\star$, as $t \to \infty$.

- Observe $|x^t - x^\star| = |f(x^{t-1}) - f(x^\star)| \leq \gamma |x^{t-1} - x^\star|$

# Detour: fix-point solution

- Suppose we want to find an $x^\star$ s.t. $x^\star = f(x^\star)$, $\quad f : [a, b] \mapsto [a, b]$

- A naive approach to find $x^\star$ :

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

- Suppose $f$ is a contraction mapping: $\forall x, x', \ |f(x) - f(x')| \leq \gamma |x - x'|$, for $\gamma \in [0,1)$. Then it converges, i.e. $x^t \to x^\star$, as $t \to \infty$.

$$\leq \gamma^2 |x^{t-2} - x^\star|$$

- Observe $|x^t - x^\star| = |f(x^{t-1}) - f(x^\star)| \leq \gamma |x^{t-1} - x^\star|$ $\quad \cdots \leq \gamma^t |x^0 - x^\star|$

- If we want $|x^t - x^\star| \leq \epsilon$, then how should we set $t$?

$$\leq \gamma^t (b - a)$$

17

# Detour: fix-point solution

- Suppose we want to find an $x^\star$ s.t. $x^\star = f(x^\star), \quad f : [a, b] \mapsto [a, b]$

- A naive approach to find $x^\star$ :

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

- Suppose $f$ is a contraction mapping: $\forall x, x', \ |f(x) - f(x')| \le \gamma |x - x'|, \ $ for $\gamma \in [0,1)$. Then it converges, i.e. $x^t \to x^\star, \ $ as $t \to \infty$.

- Observe $|x^t - x^\star| = |f(x^{t-1}) - f(x^\star)| \le \gamma |x^{t-1} - x^\star|$

- If we want $|x^t - x^\star| \le \epsilon$, then how should we set $t$?

  - Want $t$ such that $\gamma^t (b - a) \le \epsilon$

sufficient

$t = \dfrac{\ln\left(\frac{\epsilon}{b-a}\right)}{\ln \gamma}$

$\approx -\ln\left(\frac{b-a}{\epsilon}\right)$

$\dfrac{}{\ln \gamma}$

$\le \dfrac{\ln(b-a/\epsilon)}{1-\gamma} \leftarrow$ making

$t$ larger only is better

$\circledast$ since $1+x \le e^x$

$\Rightarrow$ ly $(1+x) \le x$

replacing $\gamma - 1 \le x$

$\Rightarrow -\dfrac{1}{\log \gamma} \le \dfrac{1}{1-\gamma}$

using $\circledast$

17

# Detour: fix-point solution

- Suppose we want to find an $x^\star$ s.t. $x^\star = f(x^\star), \quad f : [a, b] \mapsto [a, b]$

- A naive approach to find $x^\star$ :

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

- Suppose $f$ is a contraction mapping: $\forall x, x', \ |f(x) - f(x')| \leq \gamma |x - x'|$, for $\gamma \in [0,1)$. Then it converges, i.e. $x^t \to x^\star$, as $t \to \infty$.

- Observe $|x^t - x^\star| = |f(x^{t-1}) - f(x^\star)| \leq \gamma |x^{t-1} - x^\star|$

- If we want $|x^t - x^\star| \leq \epsilon$, then how should we set $t$?

  - Want $t$ such that $\gamma^t (b - a) \leq \epsilon$

  - $\implies t \geq \ln\big((b - a)/\epsilon\big)/(1 - \gamma)$

$t = \log \frac{b-a}{\epsilon}$

using

$\log\left(\frac{1}{\gamma}\right)$

also ok but

is a more interpretable

expression.

# Value Iteration Algorithm:

# Value Iteration Algorithm:

1. Initialization: $V^0(s) = 0, \ \forall s$

2. For $t = 0, \ldots T - 1$

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

3. Return: $V^T(s)$

# Value Iteration Algorithm:

1. Initialization: $V^0(s) = 0, \ \forall s$

2. For $t = 0, \ldots T - 1$

$$V^{t+1}(s) = \max_a \left\{ r(s, a) + \gamma \sum_{s' \in S} P(s' \mid s, a) V^t(s') \right\}, \ \forall s$$

3. Return: $V^T(s)$

$$\pi(s) = \arg\max_a \left\{ r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot \mid s, a)} V^T(s') \right\}$$

$$V^t \in \mathbb{R}^{|S|}$$

# Value Iteration Algorithm:

*fixing* $s, a$

$O(|S|)$

1. Initialization: $V^0(s) = 0, \ \forall s$
2. For $t = 0, \dots T - 1$

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

3. Return: $V^T(s)$

$$\pi(s) = \arg \max_a \left\{ r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} V^T(s') \right\}$$

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

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

18

# Value Iteration Algorithm:

1. Initialization: $V^0(s) = 0, \ \forall s$
2. For $t = 0, \ldots T - 1$

$$V^{t+1}(s) = \max_a \left\{ r(s, a) + \gamma \sum_{s' \in S} P(s' \mid s, a) V^t(s') \right\}, \ \forall s$$

3. Return: $V^T(s)$

$$\pi(s) = \arg\max_a \left\{ r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot \mid s, a)} V^T(s') \right\}$$

- What is the per iteration computational complexity of VI? (assume scalar $+, -, \times, \div$ are $O(1)$ operations)
- Guarantee: VI is fix-point iteration, which contracts, so $V^t \to V^\star$, as $t \to \infty$

# Define Bellman Operator $\mathcal{T}$ :

# Define Bellman Operator $\mathscr{T}$ :

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

# Define Bellman Operator $\mathscr{T}$:

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

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

# Define Bellman Operator $\mathscr{T}$:

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

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

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

- Bellman equations: $V = \mathscr{T}V$

# Define Bellman Operator $\mathscr{T}$:

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

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

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

- Bellman equations: $V = \mathscr{T}V$

- Value iteration: $V^{t+1} \leftarrow \mathscr{T}V^t$

# Convergence of Value Iteration:

# Convergence of Value Iteration:

- The "infinity norm": For any vector $x \in R^d$, define $|x|_\infty = \max_i |x_i|$

# Convergence of Value Iteration:

- The "infinity norm": For any vector $x \in R^d$, define $|x|_\infty = \max_i |x_i|$

- Theorem: Given any $V, V'$, we have: $\|\mathscr{T}V - \mathscr{T}V'\|_\infty \leq \gamma \|V - V'\|_\infty$

# Convergence of Value Iteration:

- The "infinity norm": For any vector $x \in R^d$, define $|x|_\infty = \max_i |x_i|$

- Theorem: Given any $V, V'$, we have: $\|\mathscr{T}V - \mathscr{T}V'\|_\infty \leq \gamma \|V - V'\|_\infty$

- Corollary: If we set $T = \dfrac{1}{1 - \gamma} \ln\left(\dfrac{1}{\epsilon(1 - \gamma)}\right)$ iterations,

  VI will return a value $V^T$ s.t. $\|V^T - V^\star\|_\infty \leq \epsilon$.

# Convergence of Value Iteration:

- The "infinity norm": For any vector $x \in R^d$, define $|x|_\infty = \max_i |x_i|$

- Theorem: Given any $V, V'$, we have: $\|\mathscr{T}V - \mathscr{T}V'\|_\infty \leq \gamma\|V - V'\|_\infty$

- Corollary: If we set $T = \dfrac{1}{1-\gamma}\ln\left(\dfrac{1}{\epsilon(1-\gamma)}\right)$ iterations,

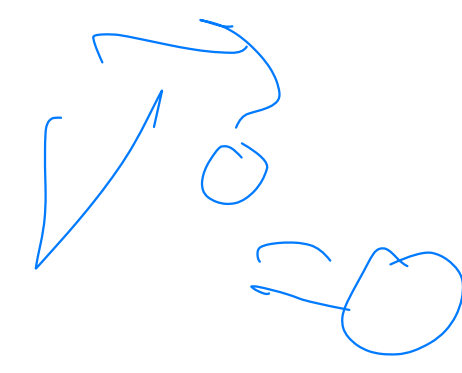  VI will return a value $V^T$ s.t. $\|V^T - V^\star\|_\infty \leq \epsilon$.

  - VI then has computational complexity $O(|S|^2|A|T)$.

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

# Today

- Recap

- Infinite Horizon MDPs

    - Optimality & the Bellman Equations

    - Value Iteration

    ✓ • Policy Iteration

# Policy Iteration (PI)

- Initialization: choose a policy $\pi^0 : S \mapsto A$

- For $t = 0,1,\ldots T - 1$

  1. **Policy Evaluation**:     given $\pi^t$, compute $Q^{\pi^t}(s, a)$:

  2. **Policy Improvement**:   set $\pi^{t+1}(s) := \arg\max_a Q^{\pi^t}(s, a)$

# Policy Iteration (PI)

- Initialization: choose a policy $\pi^0 : S \mapsto A$

- For $t = 0,1,\ldots T - 1$

  1. **Policy Evaluation**: given $\pi^t$, compute $Q^{\pi^t}(s, a)$:

  2. **Policy Improvement**: set $\pi^{t+1}(s) := \arg\max_a Q^{\pi^t}(s, a)$

- What's the computational complexity per iteration?
  Let's do this in parts:

# Policy Iteration (PI)

- Initialization: choose a policy $\pi^0 : S \mapsto A$

- For $t = 0,1,\dots T-1$

  1. **Policy Evaluation**:  given $\pi^t$, compute $Q^{\pi^t}(s,a)$:

  2. **Policy Improvement**:  set $\pi^{t+1}(s) := \arg\max_a Q^{\pi^t}(s,a)$

- What's the computational complexity per iteration?
  Let's do this in parts:

- Computing $V^{\pi^t}$:    solve lin system    with $\pi^t$    $O(|S|^3)$

- Computing $Q^{\pi^t}$ with $V^{\pi^t}$:    $Q^{\pi^t}(s,a) = r(s,a) + \gamma \sum_{s'} P(s'|s,a) V^{\pi^t}(s')$

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

step 1 $\{$

# Policy Iteration (PI)

- Initialization: choose a policy $\pi^0 : S \mapsto A$

- For $t = 0, 1, \ldots T - 1$

    1. **Policy Evaluation**:    given $\pi^t$, compute $Q^{\pi^t}(s, a)$:

    2. **Policy Improvement**:   set $\pi^{t+1}(s) := \arg\max_a Q^{\pi^t}(s, a)$

- What's the computational complexity per iteration?
  Let's do this in parts:

    - Computing $V^{\pi^t}$:

    - Computing $Q^{\pi^t}$ with $V^{\pi^t}$:

    - Computing $\pi^t$ with $Q^{\pi^t}$:    $O\left(|S||A|\right)$

step

# Policy Iteration (PI)

- Initialization: choose a policy $\pi^0 : S \mapsto A$

- For $t = 0,1,\dots T-1$

  1. **Policy Evaluation**:  given $\pi^t$, compute $Q^{\pi^t}(s,a)$:

  2. **Policy Improvement**:  set $\pi^{t+1}(s) := \arg\max_a Q^{\pi^t}(s,a)$

- What's the computational complexity per iteration?
  Let's do this in parts:

  - Computing $V^{\pi^t}$:

  - Computing $Q^{\pi^t}$ with $V^{\pi^t}$:

  - Computing $\pi^{t+1}$ with $Q^{\pi^t}$:

  Per iteration complexity:

# Policy Iteration (PI)

- Initialization: choose a policy $\pi^0 : S \mapsto A$

- For $t = 0, 1, \ldots T - 1$

  1. **Policy Evaluation**:      given $\pi^t$, compute $Q^{\pi^t}(s, a)$:

  2. **Policy Improvement**:   set $\pi^{t+1}(s) := \arg \max_a Q^{\pi^t}(s, a)$

- What's the computational complexity per iteration?
  Let's do this in parts:

  - Computing $V^{\pi^t}$:

  - Computing $Q^{\pi^t}$ with $V^{\pi^t}$:

  - Computing $\pi^{t+1}$ with $Q^{\pi^t}$:

Per iteration complexity:

- What about convergence?  $O\left( |S|^3 + |S|^2 |A| \right)$

22

# Convergence of Policy Iteration:

# Convergence of Policy Iteration:

- Theorem: PI has two properties:

# Convergence of Policy Iteration:

- <span style="color:#29abe2">Theorem:</span> PI has two properties:

  - montone improvement: $V^{\pi^{t+1}}(s) \geq V^{\pi^t}(s)$

# Convergence of Policy Iteration:

- Theorem: PI has two properties:

  - montone improvement: $V^{\pi^{t+1}}(s) \geq V^{\pi^t}(s)$

  - "contraction": $\|V^{\pi^{t+1}} - V^\star\|_\infty \leq \gamma \|V^{\pi^t} - V^\star\|_\infty$

# Convergence of Policy Iteration:

- Theorem: PI has two properties:

  - montone improvement: $V^{\pi^{t+1}}(s) \geq V^{\pi^t}(s)$

  - "contraction": $\|V^{\pi^{t+1}} - V^\star\|_\infty \leq \gamma \|V^{\pi^t} - V^\star\|_\infty$

- Corollary: If we set $T = \dfrac{1}{1-\gamma} \ln\left(\dfrac{1}{\epsilon(1-\gamma)}\right)$ iterations,

  PI will return a policy $\pi^{t+1}$ s.t. $\|V^{\pi^{t+1}} - V^\star\|_\infty \leq \epsilon$

# Convergence of Policy Iteration:

- Theorem: PI has two properties:

  - montone improvement: $V^{\pi^{t+1}}(s) \geq V^{\pi^t}(s)$

  - "contraction": $\|V^{\pi^{t+1}} - V^\star\|_\infty \leq \gamma \|V^{\pi^t} - V^\star\|_\infty$

- Corollary: If we set $T = \dfrac{1}{1-\gamma} \ln\left(\dfrac{1}{\epsilon(1-\gamma)}\right)$ iterations,

  PI will return a policy $\pi^{t+1}$ s.t. $\|V^{\pi^{t+1}} - V^\star\|_\infty \leq \epsilon$

  - with total computational complexity $O\left(\left(|S|^3 + |S|^2 |A|\right)T\right).$

# Summary:

- **Discounted infinite horizon MDP:**
  - Key Concepts: Bellman equations; Value Iteration; Policy Iteration

Attendance:
bit.ly/3RcTC9T



Feedback:
bit.ly/3RHtlxy