

# **Optimality in Markov Decision Processes**

**Lucas Janson and Sham Kakade**

**CS/Stat 184: Introduction to Reinforcement Learning  
Fall 2022**

# Today

- Recap
  - In the bandit setting, we were learning.
  - Now we are starting with computation (of the optimal policy).
- Today:
  - Is there a simple way to characterize the optimal policy?
    - The Bellman Optimality Equations
  - The state-action visitation distribution

# Recap

# The Objective

- A “stationary” policy  $\pi : S \mapsto A$ 
  - “stationary” means not history dependent
  - we could also consider  $\pi$  to be random and a function of the history
- Sampling a trajectory: from a given policy  $\pi$  starting at state  $s_0$ :
  - 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  $s_{t+1} \sim P(\cdot | s_t, a_t)$
- Objective: given state starting state  $s$ ,  
find a policy  $\pi$  that maximizes our expected, discounted future reward:

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

# 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:

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

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

# Bellman Consistency Equations:

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

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

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

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

# Notation

- For a distribution  $D$  over a finite set  $\mathcal{X}$ ,

$$E_{x \sim D}[f(x)] = \sum_{x \in \mathcal{X}} D(x)f(x)$$

- $P(\cdot | s, a)$  is a distribution, where  $P(s' | s, a)$  specifies the probability of the transition  $(s, a) \rightarrow s'$

- We will use notation:

$$E_{s' \sim P(\cdot | s, a)}[f(s')] = \sum_{s' \in S} P(s' | s, a)f(s')$$

And, if we are short on space and when it is clear, sometimes:

$$E_{s' \sim P(s, a)}[f(s')] = \sum_{s' \in S} P(s' | s, a)f(s')$$

# Today:

## Optimality in Markov Decision Processes



# Property 1 of an Optimal Policy $\pi^\star$

Even if we consider policies which are randomized and history dependent, the policy which optimizes the the value (starting from any state  $s$ ) is deterministic and memoryless.

- Defs:
  - “NonStat+Rand”: the set of all non-stationary (history dependent), randomized policies.
  - “Stat+Det”: the set of all deterministic, stationary (memoryless), policies.
- For any  $s$ , we have that:

$$\max_{\pi \in \text{NonStat+Rand}} V^\pi(s) = \max_{\pi \in \text{Stat+Det}} V^\pi(s)$$

[see theorem 1.7 in AJKS—no need to understand the proof]

- Part of the reason why: the transition function  $P(s_{t+1} \mid s_t, a_t)$  is no a function of  $t$ .  
So knowledge of  $s_t$  implies that using the history doesn't alter the next state distribution.
- (Until we say otherwise) **we limit ourselves to only consider det. stationary policies.**

## Property 2 of an Optimal Policy $\pi^\star$

- The optimal value at state  $s$  is defined as:

$$V^\star(s) = \max_{\pi} V^\pi(s)$$

Note the above permits the optimizing policy to be a function of the starting state  $s$ .

- There always exists a deterministic policy  $\pi^\star$  such that, for all states  $s$ ,

$$V^{\pi^\star}(s) = V^\star(s)$$

[see theorem 1.7 in AJKS—no need to understand the proof]

- There is an optimal policy that simultaneously dominates all  $\pi$ , from any starting state.

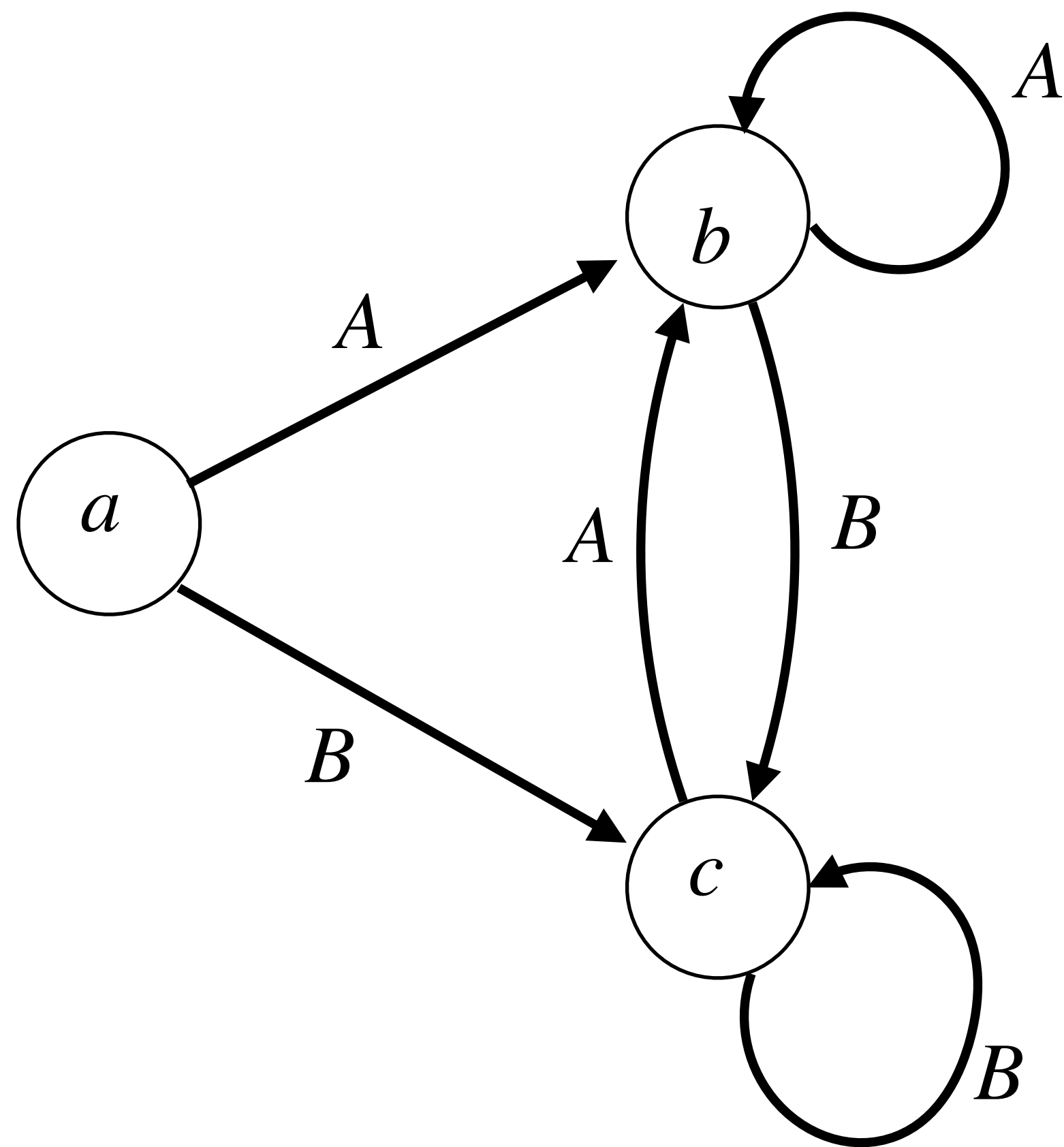
- Intuition:

$$\begin{aligned} V^\pi(s) &= r(s, \pi(s)) + \gamma E_{s' \sim P(\cdot | s, \pi(s))} [V^\pi(s')] \\ &\leq r(s, \pi(s)) + \gamma E_{s' \sim P(\cdot | s, \pi(s))} \left[ \max_{\tilde{\pi}} V^{\tilde{\pi}}(s') \right] \end{aligned}$$

( $\implies$  after reaching any state  $s'$ , we can ignore how we got to  $s'$  and instead choose the next action at  $s'$  to optimize the long term future only as a function of  $s'$ )

# Example of Optimal Policy $\pi^\star$

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



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

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

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

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

$$V^\pi(a) = 0, V^\pi(b) = 0, V^\pi(c) = 0$$

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

## Summary so far:

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

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

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

# Bellman Optimality Equations

**Theorem 1:**  $V^\star$  satisfies the following **Bellman Equations**:

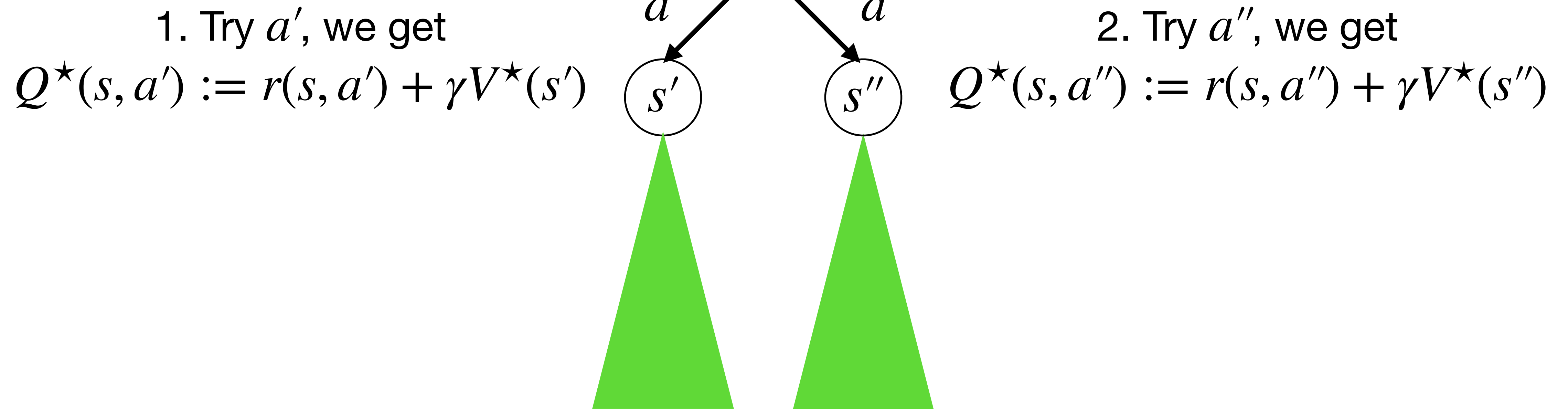
$$V^\star(s) = \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} V^\star(s') \right], \forall s$$

Also, if  $\hat{\pi}(s) = \arg \max_a Q^\star(s, a)$ , then  $\hat{\pi}$  is an optimal policy.

## Intuition for the Bellman Equations

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

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



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

# Proof of the Bellman Equations

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

Proof:

- Denote:
$$\hat{\pi}(s) := \arg \max_a [r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} V^\star(s')]$$
$$= \arg \max_a Q^\star(s, a)$$
- It suffices to show  $V^\star(s) \leq V^{\hat{\pi}}(s)$ , which would complete the proof.
- To see this completes the proof,
  - optimality of  $V^\star$  implies  $V^{\hat{\pi}}(s) \leq V^\star(s)$ .
  - and so:
$$V^\star(s) \leq V^{\hat{\pi}}(s) \leq V^\star(s).$$
- Thus  $V^{\hat{\pi}}(s) = V^\star(s)$  and  $\hat{\pi}$  is optimal.

# Completing the proof: showing $V^\star(s) \leq V^{\hat{\pi}}(s)$

- Recall:  $\hat{\pi}(s) := \arg \max_a [r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} V^\star(s')]$
- We have:
$$\begin{aligned} V^\star(s) &= r(s, \pi^\star(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^\star(s))} V^\star(s') \\ &\leq \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(s,a)} V^\star(s') \right] \\ &= r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} [V^\star(s')] \end{aligned}$$
- Proceeding recursively,
$$\begin{aligned} &\leq r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[ r(s', \hat{\pi}(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \hat{\pi}(s'))} V^\star(s'') \right] \\ &\leq r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[ r(s', \hat{\pi}(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \hat{\pi}(s'))} \left[ r(s'', \hat{\pi}(s'')) + \gamma \mathbb{E}_{s''' \sim P(s'', \hat{\pi}(s''))} V^\star(s''') \right] \right] \\ &\leq \mathbb{E} \left[ r(s, \hat{\pi}(s)) + \gamma r(s', \hat{\pi}(s')) + \dots \mid \hat{\pi} \right] = V^{\hat{\pi}}(s) \end{aligned}$$



## Summary so far:

### Theorem 1: Bellman Optimality

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

Next:

Any function  $V$  that satisfies Bellman Optimality, MUST be equal to  $V^{\star}$

# Bellman Equations, Claim 2

**Theorem 2:** For any  $V : S \rightarrow \mathbb{R}$ , if

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

Bellman Opt allows us to focus on just one step,  
i.e., to check if  $V = V^\star$ , we only need to check if the above equation holds.

# Proving Theorem 2

- Define the “maximal component distance” between  $V$  and  $V^\star$ :

$$\|V - V^\star\|_\infty = \max_s |V(s) - V^\star(s)|$$

- For  $V$  which satisfies the Bellman equations,  
suppose we could show that  $\|V - V^\star\|_\infty \leq \gamma \|V - V^\star\|_\infty$ .

$\implies$  the proof is complete because

$$\|V - V^\star\|_\infty \leq \gamma \|V - V^\star\|_\infty \leq \gamma^2 \|V - V^\star\|_\infty \leq \dots \leq \lim_{k \rightarrow \infty} \gamma^k \|V - V^\star\|_\infty = 0$$

# Proof Continued...

- For  $V$  which satisfies the Bellman equations, we want to show  $\|V - V^\star\|_\infty \leq \gamma \|V - V^\star\|_\infty$ .
- Technial observation:  $\left| \max_x f(x) - \max_x g(x) \right| \leq \max_x |f(x) - g(x)|$
- Using that  $V$  satisfies the Bellman equations, we have, for any  $s$ ,

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

# Summary Today

1.  $V^\star$  satisfies Bellman Optimality:

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

2. If  $V$  satisfies Bellman Optimality Equations,  $V(s) = \max_a \left[ r(s, a) + \mathbb{E}_{s' \sim P(s, a)} V(s') \right]$ ,  
then  $V = V^\star$ .

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

