

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).

Please provide feedback.

- 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 π to be random and a function of the history
- Sampling a trajectory: from a given policy π 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 π that maximizes our expected, discounted future reward:

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

non-stationary policies

π : histories

\rightarrow action

a_{t+1}

$\pi(r_0, s_0, a_1,$

$\dots, r_t, s_t, a_t,$

$s_{t+1})$

Infinite horizon Discounted Setting

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

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

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

$P(s'/s,a)$ is the prob. $s,a \rightarrow s'$

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:

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

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

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

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

Notation

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

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

$$P(s' | s, a)$$

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

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

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

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

Today:

Optimality in Markov Decision Processes

$$\# \text{ policies} = |A|^{|S|}$$

Property 1 of an Optimal Policy π^\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.

Property 1 of an Optimal Policy π^\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]

Property 1 of an Optimal Policy π^\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]

- ~~Intuition: $\Pr(s_{t+1} = s' | s_t, a_t, r_t, s_{t-1}, a_{t-1}, r_{t-1}, \dots) = P(s' | s_t, a_t)$~~
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 π^\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 π^\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]

Property 2 of an Optimal Policy π^\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 π^\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 π , from any starting state.

Property 2 of an Optimal Policy π^\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 π^\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 π , from any starting state.

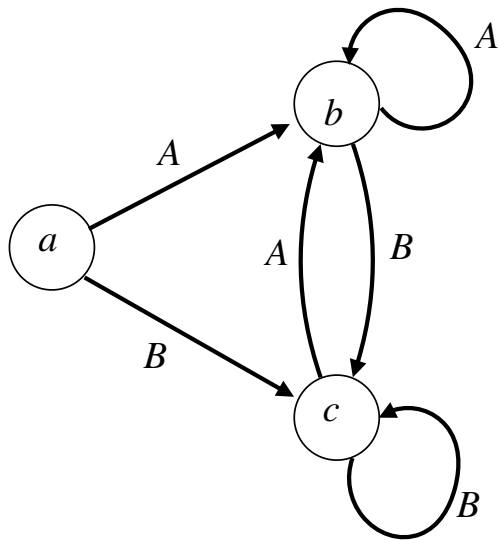
- Intuition:

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

(\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 π^\star

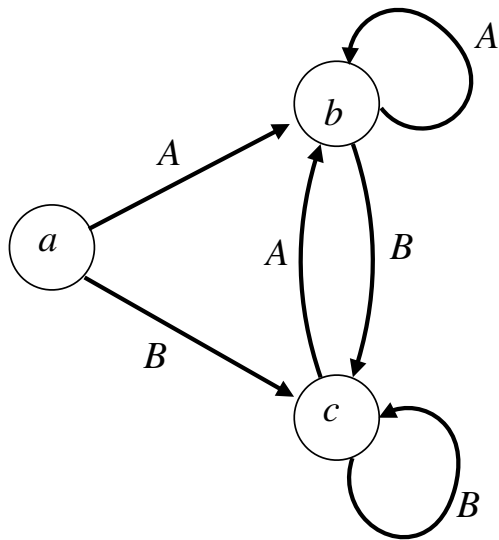
Consider the following **deterministic** MDP w/ 3 states & 2 actions



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

Example of Optimal Policy π^\star

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

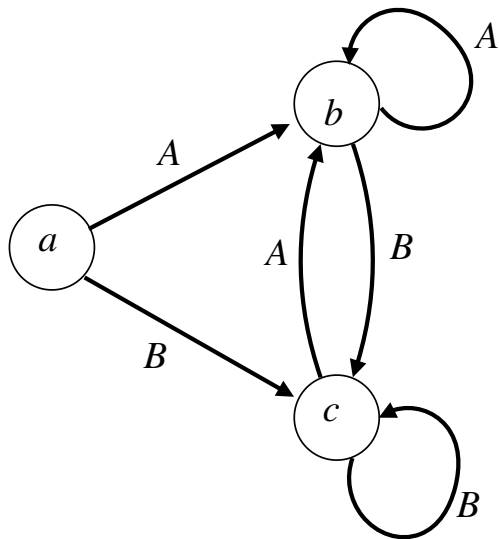


Let's say $\gamma \in (0,1)$
What's the optimal policy?

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

Example of Optimal Policy π^\star

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



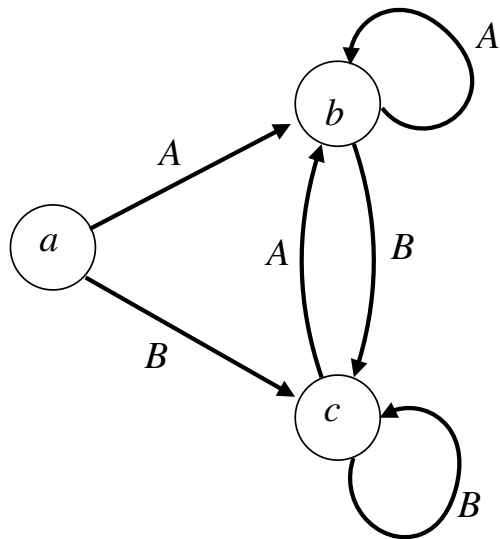
Let's say $\gamma \in (0,1)$
What's the optimal policy?

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

Reward: $r(b, A) = 1$, & 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?

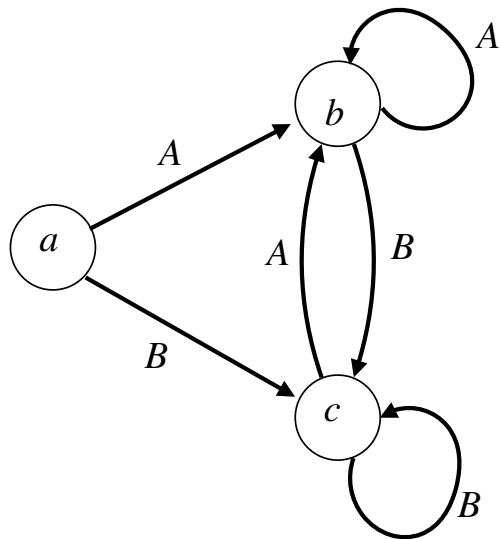
$$\pi^\star(s) = A, \forall s$$

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

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

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?

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

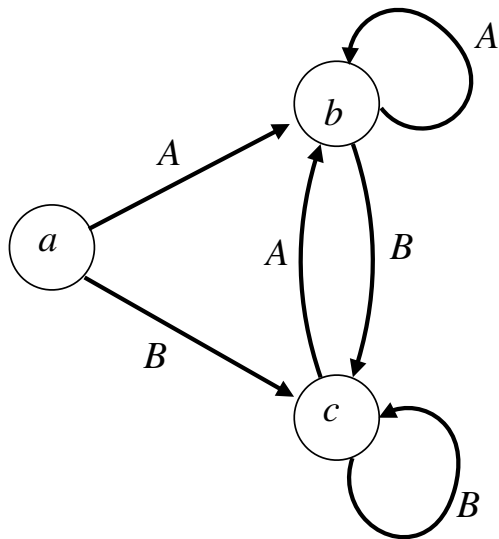
$$V^\star(a) = \frac{\gamma}{1-\gamma}, V^\star(b) = \frac{1}{1-\gamma}, V^\star(c) = \frac{\gamma}{1-\gamma}$$

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

Reward: $r(b, A) = 1$, & 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?

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

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

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

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

Summary so far:

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

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

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

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.

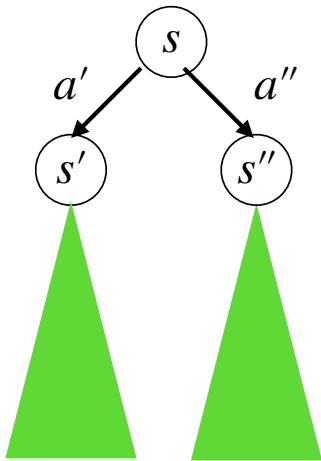
Intuition for the Bellman Equations

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

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 ?



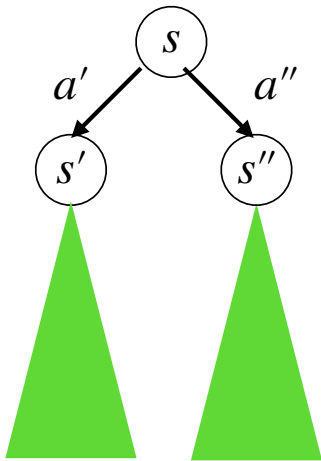
Intuition for the Bellman Equations

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

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

1. Try a' , we get

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



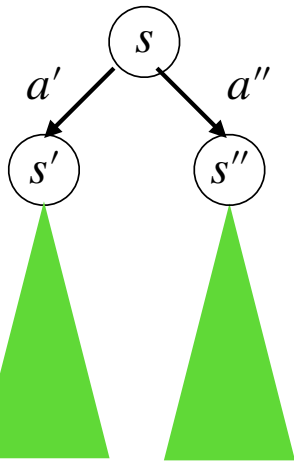
Intuition for the Bellman Equations

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

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

1. Try a' , we get

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



2. Try a'' , we get

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

Intuition for the Bellman Equations

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

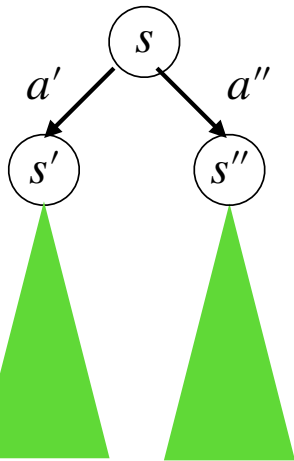
"Dynamic

Programming"

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 the Bellman Equations

Proof of the Bellman Equations

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

Proof of the Bellman Equations

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

Proof:

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:

$$\begin{aligned}\hat{\pi}(s) &:= \arg \max_a [r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} V^\star(s')] \\ &= \arg \max_a Q^\star(s, a)\end{aligned}$$

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.

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:

$$\begin{aligned}\hat{\pi}(s) &:= \arg \max_a [r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} V^\star(s')] \\ &= \arg \max_a Q^\star(s, a)\end{aligned}$$

- 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 $V^\star(s) \leq V^{\hat{\pi}}(s)$

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

- Recall: $\hat{\pi}(s) := \arg \max_a [r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} V^\star(s')]$

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:

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

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

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

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

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:

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

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

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

same argument

- Proceeding recursively,

$$\leq r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[r(s', \hat{\pi}(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \hat{\pi}(s'))} V^\star(s'') \right]$$

$$\leq r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \hat{\pi}(s))} \left[r(s', \hat{\pi}(s')) + \gamma \mathbb{E}_{s'' \sim P(s', \hat{\pi}(s'))} \left[r(s'', \hat{\pi}(s'')) + \gamma \mathbb{E}_{s''' \sim P(s'', \hat{\pi}(s''))} V^\star(s''') \right] \right]$$

$$\leq \mathbb{E} [r(s, \hat{\pi}(s)) + \gamma r(s', \hat{\pi}(s')) + \dots | \hat{\pi}] = V^{\hat{\pi}}(s)$$



Summary so far:

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:

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 that satisfies Bellman Optimality, MUST be equal to V^*

Bellman Equations, Claim 2

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

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

Bellman Equations, Claim 2

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

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

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

Proving Theorem 2

Proving Theorem 2

def $\vec{x} \in \mathbb{R}^d$.

$$\|\vec{x}\|_{\infty} = \max_i |x_i|$$

$$\|\vec{x}\|_2 = \sqrt{\sum_i x_i^2}$$

- Define the “maximal component distance” between V and V^* :

$$\|V - V^*\|_{\infty} = \max_s |V(s) - V^*(s)|$$

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...

Proof Continued...

- For V which satisfies the Bellman equations, we want to show $\|V - V^*\|_\infty \leq \gamma \|V - V^*\|_\infty$.

Proof Continued...

- For V which satisfies the Bellman equations, we want to show $\|V - V^*\|_\infty \leq \gamma \|V - V^*\|_\infty$.

- Technical observation: $|\max_x f(x) - \max_x g(x)| \leq \max_x |f(x) - g(x)|$

$\exists x' - x''$ Case 1: suppose $f(x)$ is positive

$\max_x f(x) - \max_x g(x) = f(x') - g(x'')$

$\leq f(x') - g(x')$

$\leq |f(x') - g(x')|$

$\leq \max_x |f(x) - g(x)|$

(because $g(x') \leq g(x'')$)

$f(x') = \max_x f(x)$
 $g(x'') = \max_x g(x)$

Proof Continued...

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

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

by assumption, tech obs

Proof Continued...

- For V which satisfies the Bellman equations, we want to show $\|V - V^*\|_\infty \leq \gamma \|V - V^*\|_\infty$.
- Technical 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 ,

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

$$|\mathbb{E}[f(x)]| \leq \mathbb{E}[|f(x)|]$$

$$\mathbb{E}[|f(x)|] \leq \max_x |f(x)|$$

take max
on both sides
 \Rightarrow

def.

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>



B.F. want to find $V \pi$.

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

↑
this is a "fixed" point equation.

suppose we want to find \vec{x}
s.t. $\vec{x} = f(\vec{x})$

one "hack": try the algorithm
 $\vec{x} \leftarrow f(\vec{x})$ does this work?
and "hope" it converges

here we can try:

start with some V then
 $\forall s$ do update

$$V(s) \leftarrow \max_a \left\{ r(s,a) + \gamma E_{s' \sim P(\cdot|s,a)} [V(s')] \right\}$$