

Policy Evaluation

Announcements

HW0 is out today (due March 2rd, 11:59 pm ET)

(Gradescope entry code: BP3B3N)

Office hours start this week:

Wen: Tuesday and Thursday, 10:55am - 11:30am

Wen-Ding: Friday 3pm-4pm

Hadi: Wednesday 2:30-3:30pm

Recap: Definitions

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

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

Recap: Definitions

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

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

Policy $\pi : S \mapsto A$

A^S

Recap: Definitions

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

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

Policy $\pi : S \mapsto A$

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: Definitions

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

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

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

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

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

Recap: Optimal Policy

For Discounted infinite horizon MDP, \exists a deterministic policy $\pi^* : S \mapsto A$:

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

Recap: Optimal Policy

For Discounted infinite horizon MDP, \exists a deterministic policy $\pi^* : S \mapsto A$:

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

Bellman Optimality (DP):

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

Recap: Optimal Policy

For Discounted infinite horizon MDP, \exists a deterministic policy $\pi^* : S \mapsto A$:

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

Bellman Optimality (DP):

1. For V^* , we have $V^*(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right], \forall s$
2. For V that satisfies $V(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V(s') \right], \forall s$,
we have $V(s) = V^*(s), \forall s$

Recap: Optimal Policy

For Discounted infinite horizon MDP, \exists a deterministic policy $\pi^* : S \mapsto A$:

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

Bellman Optimality (DP):

Bellman Eq

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

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

2. For V that satisfies $V(s) = \max_a \left[r(s, a) + \gamma \mathbb{E}_{s' \sim P(s,a)} V(s') \right], \forall s$,
we have $V(s) = V^*(s), \forall s$

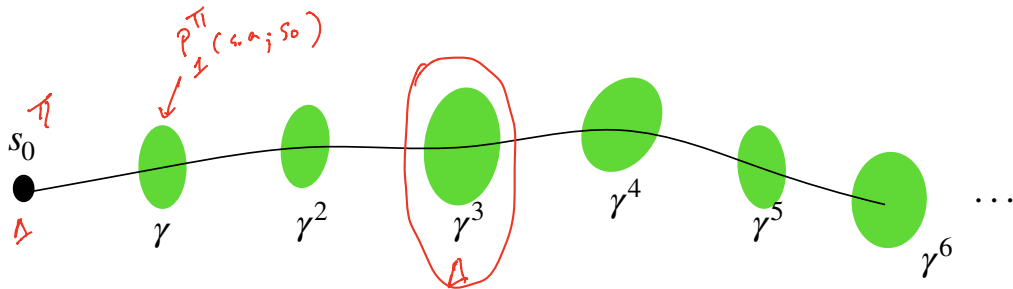
In HW0, we will study Bellman Optimality for Q^*/Q

Recap: State-action distribution

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

Recap: State-action distribution

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



Recap: State-action distribution

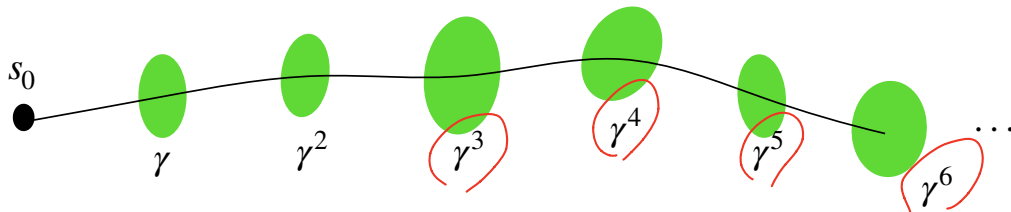
$$\frac{1-\gamma}{1-\gamma} \Rightarrow 1 + \gamma + \gamma^2 + \gamma^3 + \dots$$

$\star = \frac{1}{1-\gamma}$

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

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

$\nearrow V^\pi(s_0)$ using $d_{s_0}^\pi(s, a)$



Today: Policy Evaluation

Key Question:

**Given MDP $\mathcal{M} = (S, A, r, P, \gamma)$ & a $\pi : S \mapsto A$,
how good is π ?**

i.e., how to compute $V^\pi(s), \forall s$?

Motivation for Policy Evaluation



We want to **evaluate** our strategy against some opponent (we can abstract our strategy as policy π)

Motivation for Policy Evaluation



We want to **evaluate** our strategy against some opponent (we can abstract our strategy as policy π)



We want to **evaluate** our recommendation strategy before we release it to users

A more fundamental motivation...

Recall that we have A^S many policies.
To select the optimal policy, we need to do evaluation

Outline:

1. **Exact** Policy Evaluation

$$V = V^\pi$$

2. **Approximate** Policy Evaluation via an Iterative Algorithm

$$V \approx V^\pi$$

↑
BE

Exact Policy Evaluation

Setup: we have MDP $\mathcal{M} = (S, A, P, \gamma, r)$, and policy π , we want to compute V^π

Exact Policy Evaluation

Setup: we have MDP $\mathcal{M} = (S, A, P, \gamma, r)$, and policy π , we want to compute V^π

We know that for V^π , we have **Bellman equation:**

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

Exact Policy Evaluation

Setup: we have MDP $\mathcal{M} = (S, A, P, \gamma, r)$, and policy π , we want to compute V^π

We know that for V^π , we have **Bellman equation:**

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

This gives us S many linear constraints

Exact Policy Evaluation

Let's form linear constraints. Denote $V(s)$ as our estimator for $s \in \mathcal{S}$

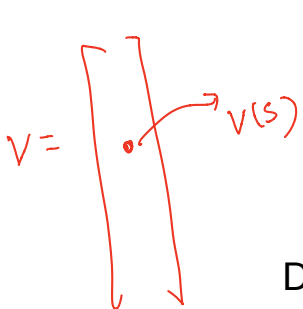
$$\forall s, V(s) = r(s, \pi(s)) + \sum_{s' \in \mathcal{S}} P(s' | s, \pi(s)) V(s')$$

$$V(s_1) = r(s_1, \pi(s_1)) + \sum_{s' \in \mathcal{S}} P(s' | s_1, \pi(s_1)) V(s')$$

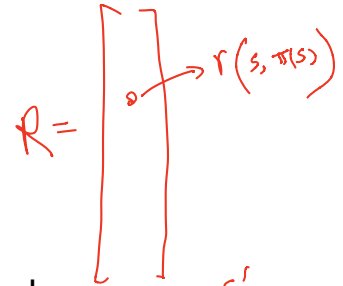
$$V(s_2) = r(s_2, \pi(s_2)) + \sum_{s' \in \mathcal{S}} P(s' | s_2, \pi(s_2)) V(s')$$

Exact Policy Evaluation

Let's form linear constraints. Denote $V(s)$ as our estimator for $s \in S$



$$\forall s, V(s) = \underbrace{r(s, \pi(s))}_{\text{circled in red}} + \gamma \sum_{s' \in S} P(s' | s, \pi(s)) V(s')$$



Denote $V \in \mathbb{R}^{|S|}$, $R \in \mathbb{R}^{|S|}$, where $R_s = r(s, \pi(s))$, and $P \in \mathbb{R}^{|S| \times |S|}$, where $P_{s',s} = P(s' | s, \pi(s))$,

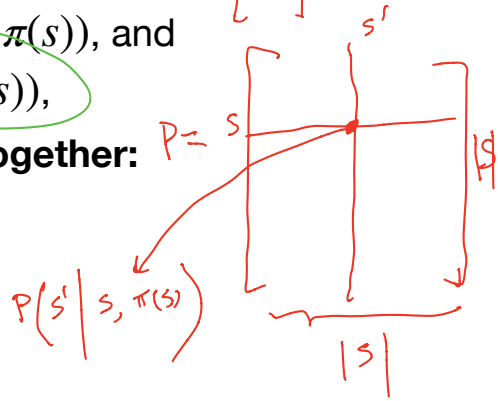
we can **combine all S many constraints together:**

★ Correction:

$$V = R + \gamma PV$$

$$P_{s,s'} = P(s' | s, \pi(s))$$

i.e., row is indexed by s , column is indexed by s'



$$\pi: S \rightarrow A$$

Exact Policy Evaluation

$V \in \mathbb{R}^{|S|}$, $R \in \mathbb{R}^{|S|}$, where $R_s = r(s, \pi(s))$, and $P \in \mathbb{R}^{|S| \times |S|}$, where $P_{s',s} = P(s' | s, \pi(s))$, we can combine all constraints together:

$$V = R + \gamma PV$$

$\sum_{s'} P(s' | s, \pi(s)) \cdot V(s')$
 $= E_{s \sim P(\cdot | s, \pi(s))} [V(s')]$

Exact Policy Evaluation

Since $V = \overset{R}{x} + \gamma PV$, we can obtain V as:

$$V = \underbrace{(I - \gamma P)^{-1}} R$$

Exact Policy Evaluation

$$V - \gamma P V = R \Rightarrow (I - \gamma P) V = R$$

Since $V = \overset{I}{\underset{R}{r}} + \gamma P V$, we can obtain V as:

$$V = (I - \gamma P)^{-1} R \checkmark$$

In HW0, we will show that $(I - \gamma P)$ is full rank (thus invertible)

A is full rank:

$$\forall x, \text{ s.t. } x \neq \vec{0} \quad Ax \neq \vec{0}$$

Summary so far:

$$V = R + \gamma P V$$

$$V = (I - \gamma P)^{-1} R$$

Summary so far:

$$\begin{array}{c} \boxed{} \\ \boxed{V(s)} \\ \boxed{} \\ \hline V \end{array} = \begin{array}{c} \boxed{} \\ \boxed{r(s, \pi(s))} \\ \boxed{} \\ \hline R \end{array} + \gamma \begin{array}{c} \boxed{} \\ \boxed{P(\cdot | s, \pi(s))} \\ \boxed{} \\ \hline P \end{array} \begin{array}{c} \boxed{} \\ \boxed{} \\ \boxed{} \\ \hline V \end{array}$$

$$\underline{\underline{Q^\pi \in \mathbb{R}^{|S|}}}$$

$$V = (I - \gamma P)^{-1} R \quad \rightarrow \quad |S \times S|$$

Downside: computation expensive: matrix inverse is $O(S^3)$



Outline:

✓ 1. Exact Policy Evaluation $\leftarrow O(s^3)$

2. Approximate Policy Evaluation via an Iterative Algorithm

(An approximation solution could be enough, i.e., trade accuracy for computation)

Detour: fix-point solution

Consider $x^\star = f(x^\star)$, $f: [a, b] \mapsto [a, b]$

Detour: fix-point solution

Consider $x^* = f(x^*)$, $f: [a, b] \mapsto [a, b]$

Common approach to find x^* :

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

$f f f f \dots (x^0)$

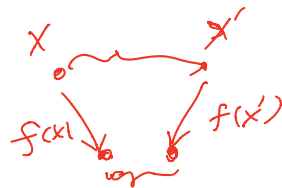
x^*

Detour: fix-point solution

Consider $x^* = f(x^*)$, $f: [a, b] \mapsto [a, b]$

Common approach to find x^* :

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



If f is a contraction mapping,

i.e., $\forall x, x', |f(x) - f(x')| \leq \gamma |x - x'|$, for some $\gamma \in [0, 1)$, then:

$x^t \rightarrow x^*$, as $t \rightarrow \infty$

$$\left| \underset{\Delta}{x^{t+1}} - x^* \right| = \left| \underset{\Delta}{f(x^t)} - f(x^*) \right| \leq \gamma \left| \underset{\Delta}{x^t} - x^* \right| \leq \gamma^2 \left| x^{t+1} - x^* \right|$$

V^π is a fix-point solution:

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

V^π is a fix-point solution:

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

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

$$V^\pi = (R + \gamma P V^\pi)$$

$f(V^\pi)$

Iterative Policy Evaluation:

$$r(s,a) \in [0,1]$$
$$1 + \gamma + \gamma^2 + \gamma^3 + \dots = \frac{1}{1-\gamma}$$

Algorithm (Iterative PE):

Start with some initialization $V^0 \in [0, 1/(1-\gamma)]^{|S|}$, repeat for $t = 0 \dots$:

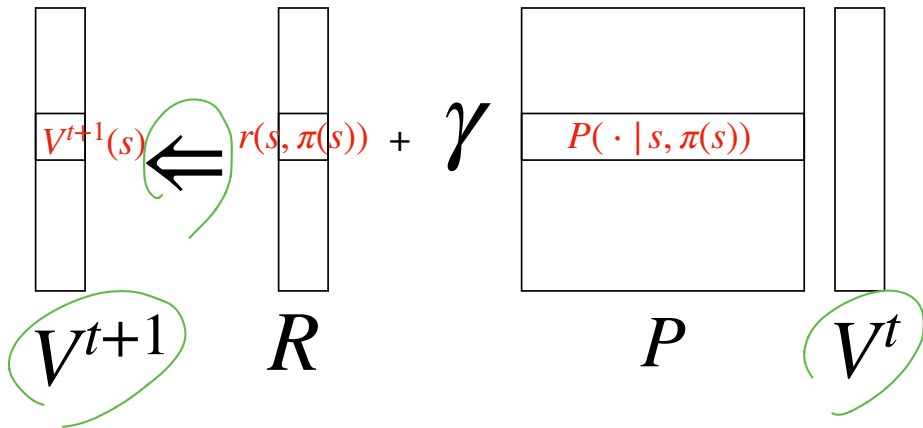
$$V^{t+1} \leftarrow R + \gamma P V^t$$

Iterative Policy Evaluation:

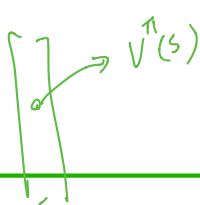
Algorithm (Iterative PE):

Start with some initialization $V^0 \in [0, 1/(1 - \gamma)]^{|S|}$, repeat for $t = 0 \dots$:

$$V^{t+1} \leftarrow R + \gamma P V^t$$



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



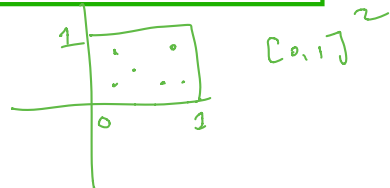
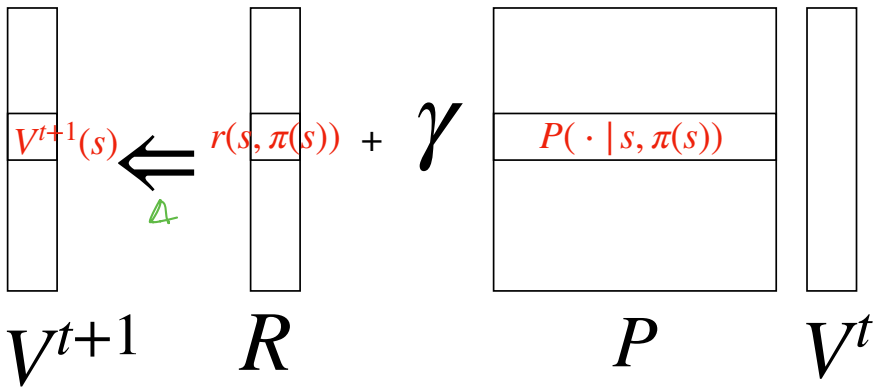
Iterative Policy Evaluation:

$$x \in [a, b]^d$$

$$x_i \in [a, b], \forall i \in [d]$$

Algorithm (Iterative PE):
 Start with some initialization $V^0 \in [0, 1/(1-\gamma)]^{|S|}$, repeat for $t = 0 \dots$:

$$V^{t+1} \leftarrow R + \gamma P V^t$$



Q: What's computation complexity per iteration?

$$O(S^2)$$

per-iteration

Iterative Policy Evaluation:

$$V^{t+1} \leftarrow R + \gamma P V^t$$

The diagram illustrates the convergence of the Bellman optimality equation for value iteration. It shows two equations:

$$V^{t+1}(s) = r(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi(s))} V^t(s')$$
$$V^\pi(s) = r(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi(s))} V^\pi(s')$$

Hand-drawn green annotations include:

- A large circle around the entire diagram.
- A circle around the $r(s, \pi(s))$ term in both equations.
- A circle around the $\gamma \mathbb{E}_{s' \sim P(s, \pi(s))} V^t(s')$ term in the top equation.
- A circle around the $\gamma \mathbb{E}_{s' \sim P(s, \pi(s))} V^\pi(s')$ term in the bottom equation.
- An arrow pointing from the $\gamma \mathbb{E}_{s' \sim P(s, \pi(s))} V^t(s')$ term to the $\gamma \mathbb{E}_{s' \sim P(s, \pi(s))} V^\pi(s')$ term, labeled $\gamma \cdot \epsilon$.
- An equals sign $= 0$ between the two equations.
- A label ϵ at the bottom right, indicating the error term.
- A label $\gamma \epsilon$ at the bottom center, indicating the error term.

Convergence of Iterative PE

Theorem:

Recall $\gamma \in [0,1)$. After t iterations, we have:

$$\forall s, \left| V^t(s) - V^\pi(s) \right| \leq \gamma^t \left\| V^0 - V^\pi \right\|_\infty$$

$\|v\|_\infty = \max_{i \in \mathcal{S}} |V_i|$

Convergence of Iterative PE

Theorem:

Recall $\gamma \in [0,1)$. After t iterations, we have:

$$\forall s, \left| V^t(s) - V^\pi(s) \right| \leq \gamma^t \left\| V^0 - V^\pi \right\|_\infty$$

$$\forall s, \left| V^{t+1}(s) - V^\pi(s) \right|$$

Convergence of Iterative PE

Theorem:

Recall $\gamma \in [0, 1)$. After t iterations, we have:

$$\forall s, \left| V^t(s) - V^\pi(s) \right| \leq \gamma^t \left\| V^0 - V^\pi \right\|_\infty$$

$$\begin{aligned} & \forall s, \left| V^{t+1}(s) - V^\pi(s) \right| \\ &= \left| \underbrace{r(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi(s))} V^t(s')}_{\substack{\text{By Alg} \\ \swarrow}} - \left(\underbrace{r(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi(s))} V^\pi(s')}_{\swarrow \text{BE on } V^\pi} \right) \right| \end{aligned}$$

Convergence of Iterative PE

Theorem:

Recall $\gamma \in [0, 1)$. After t iterations, we have:

$$\forall s, \left| V^t(s) - V^\pi(s) \right| \leq \gamma^t \left\| V^0 - V^\pi \right\|_\infty$$

$$\begin{aligned} & \forall s, \left| V^{t+1}(s) - V^\pi(s) \right| \\ &= \left| r(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi(s))} V^t(s') - \left(r(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi(s))} V^\pi(s') \right) \right| \\ &= \gamma \left| \mathbb{E}_{s' \sim P(\cdot | s, \pi(s))} V^t(s') - \mathbb{E}_{s' \sim P(\cdot | s, \pi(s))} V^\pi(s') \right| \end{aligned}$$

Convergence of Iterative PE

Theorem:

Recall $\gamma \in [0, 1)$. After t iterations, we have:

$$\forall s, \left| V^t(s) - V^\pi(s) \right| \leq \gamma^t \left\| V^0 - V^\pi \right\|_\infty$$

$$\begin{aligned} & \forall s, \left| V^{t+1}(s) - V^\pi(s) \right| \\ &= \left| r(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi(s))} V^t(s') - \left(r(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi(s))} V^\pi(s') \right) \right| \\ &= \gamma \left| \underbrace{\mathbb{E}_{s' \sim P(\cdot | s, \pi(s))} V^t(s')} - \underbrace{\mathbb{E}_{s' \sim P(\cdot | s, \pi(s))} V^\pi(s')} \right| \\ &\leq \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi(s))} \left| V^t(s') - V^\pi(s') \right| \end{aligned}$$

Convergence of Iterative PE

Theorem:

Recall $\gamma \in [0, 1)$. After t iterations, we have:

$$\forall s, \left| V^t(s) - V^\pi(s) \right| \leq \gamma^t \left\| V^0 - V^\pi \right\|_\infty$$

$$\forall s, \left| V^{t+1}(s) - V^\pi(s) \right|$$

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

$$= \gamma \left| \mathbb{E}_{s' \sim P(\cdot | s, \pi(s))} V^t(s') - \mathbb{E}_{s' \sim P(\cdot | s, \pi(s))} V^\pi(s') \right|$$

$$\leq \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi(s))} \left| V^t(s') - V^\pi(s') \right|$$

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

$$\leq \gamma \left\| V^t - V^\pi \right\|_\infty$$

Convergence of Iterative PE

Theorem:

Recall $\gamma \in [0, 1)$. After t iterations, we have:

$$\forall s, \left| V^t(s) - V^\pi(s) \right| \leq \gamma^t \left\| V^0 - V^\pi \right\|_\infty$$

$$\forall s, \left| V^{t+1}(s) - V^\pi(s) \right|$$

$$\forall s \quad \left| V^{t+1}(s) - V^\pi(s) \right| \leq \gamma \left\| V^t - V^\pi \right\|_\infty$$
$$\Rightarrow \left\| V^{t+1} - V^\pi \right\|_\infty \leq \gamma \left\| V^t - V^\pi \right\|_\infty$$

$$\stackrel{\textcircled{=}}{=} \left| r(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi(s))} V^t(s') - \left(r(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi(s))} V^\pi(s') \right) \right|$$

$$= \gamma \left| \mathbb{E}_{s' \sim P(\cdot | s, \pi(s))} V^t(s') - \mathbb{E}_{s' \sim P(\cdot | s, \pi(s))} V^\pi(s') \right|$$

$$\leq \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi(s))} \left| V^t(s') - V^\pi(s') \right|$$

$$\leq \gamma \left\| V^t - V^\pi \right\|_\infty \Rightarrow \left\| V^{t+1} - V^\pi \right\|_\infty \leq \gamma \left\| V^t - V^\pi \right\|_\infty \leq \dots \leq \gamma^{t+1} \left\| V^0 - V^\pi \right\|_\infty$$

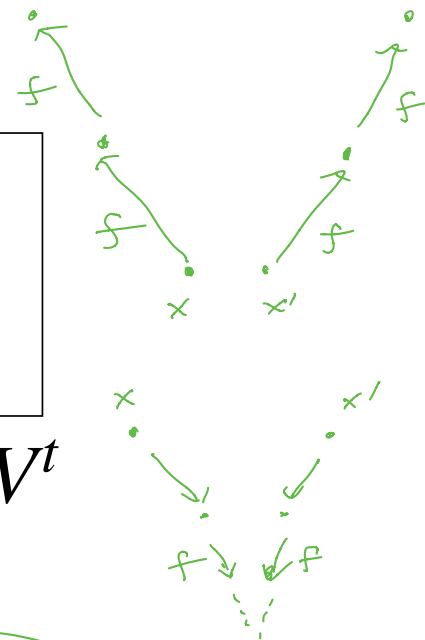
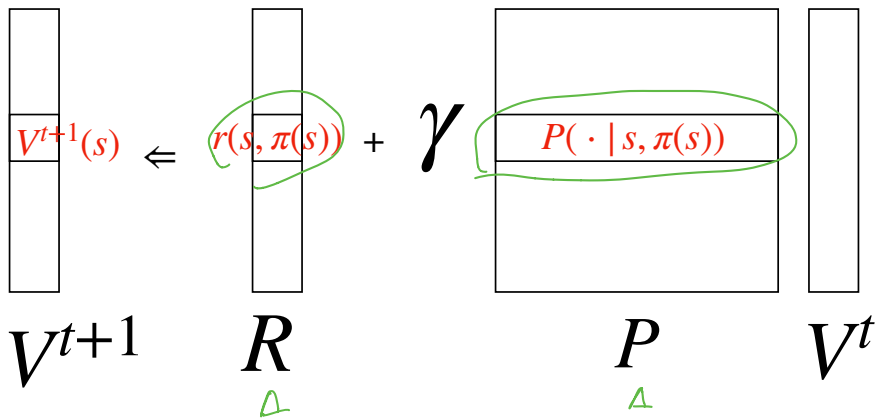
$\gamma < 1$

Summary so far:

Handwritten notes:

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

$$f(v) := R + \gamma P \cdot v$$




Convergence:

$$\| V^{t+1} - V^\pi \|_\infty \leq \gamma \| V^t - V^\pi \|_\infty \leq \gamma^{t+1} \| V^0 - V^\pi \|_\infty$$

The entire equation is circled in green.

Outline:

 1. Exact Policy Evaluation

 2. Approximate Policy Evaluation via an Iterative Algorithm

Summary

↙ Deterministic

Key Question today: Given MDP \mathcal{M} , and a policy π , How to compute $V^\pi(s), \forall s$?

Summary

Key Question today: Given MDP \mathcal{M} , and a policy π , How to compute $V^\pi(s), \forall s$?

1. The exact algorithm $V = (I - \gamma P)^{-1}R$ requires matrix inverse $O(S^3)$

Summary

Key Question today: Given MDP \mathcal{M} , and a policy π , How to compute $V^\pi(s), \forall s$?

1. The exact algorithm $V = (I - \gamma P)^{-1}R$ requires matrix inverse $O(S^3)$

1. For iterative PE algorithm, to find a ϵ accurate value function, we need # of iterations:

Summary

Key Question today: Given MDP \mathcal{M} , and a policy π , How to compute $V^\pi(s), \forall s$?

1. The exact algorithm $V = (I - \gamma P)^{-1}R$ requires matrix inverse $O(S^3)$

1. For iterative PE algorithm, to find a ϵ accurate value function, we need # of iterations:

$$\ln \left(\frac{\|V^0 - V^*\|_\infty}{\epsilon} \right) / \ln(1/\gamma)$$

$$\begin{aligned} \|V^t - V^\pi\|_\infty &\leq \gamma^t \|V_0 - V^\pi\|_\infty \\ &\leq \epsilon \\ \Rightarrow \gamma^t &\leq \frac{\|V_0 - V^\pi\|_\infty}{\epsilon} \end{aligned}$$

Summary

Key Question today: Given MDP \mathcal{M} , and a policy π , How to compute $V^\pi(s), \forall s$?

1. The exact algorithm $V = (I - \gamma P)^{-1}R$ requires matrix inverse $O(S^3)$

1. For iterative PE algorithm, to find a ϵ accurate value function, we need # of iterations:

$$\ln \left(\frac{\|V^0 - V^*\|_\infty}{\epsilon} \right) / \ln(1/\gamma)$$

Computation wise, we need $O \left(S^2 \ln \left(\frac{1}{\epsilon} \right) \right)$

versus S^3

$\epsilon = 0.01$

Summary

Key Question today: Given MDP \mathcal{M} , and a policy π , How to compute $V^\pi(s), \forall s$?

Bellman Equation



A fix-point equation:

$$V^\pi = R + \gamma P V^\pi$$

Summary

Key Question today: Given MDP \mathcal{M} , and a policy π , How to compute $V^\pi(s), \forall s$?

Bellman Equation



A fix-point equation:

$$V^\pi = R + \gamma P V^\pi$$

**Fix-point iteration
framework**

Summary

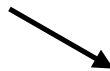
Key Question today: Given MDP \mathcal{M} , and a policy π , How to compute $V^\pi(s), \forall s$?

Bellman Equation



A fix-point equation:

$$V^\pi = R + \gamma P V^\pi$$



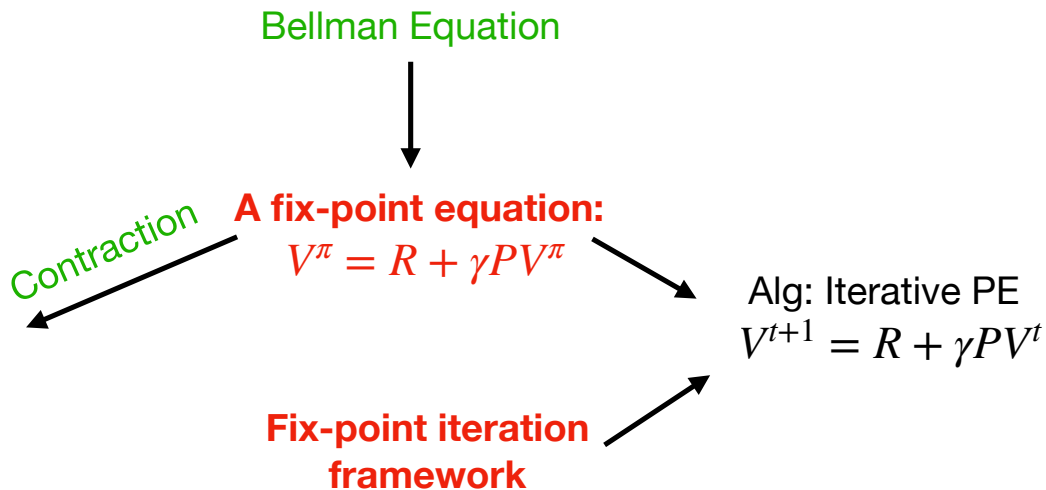
Alg: Iterative PE
 $V^{t+1} = R + \gamma P V^t$

**Fix-point iteration
framework**



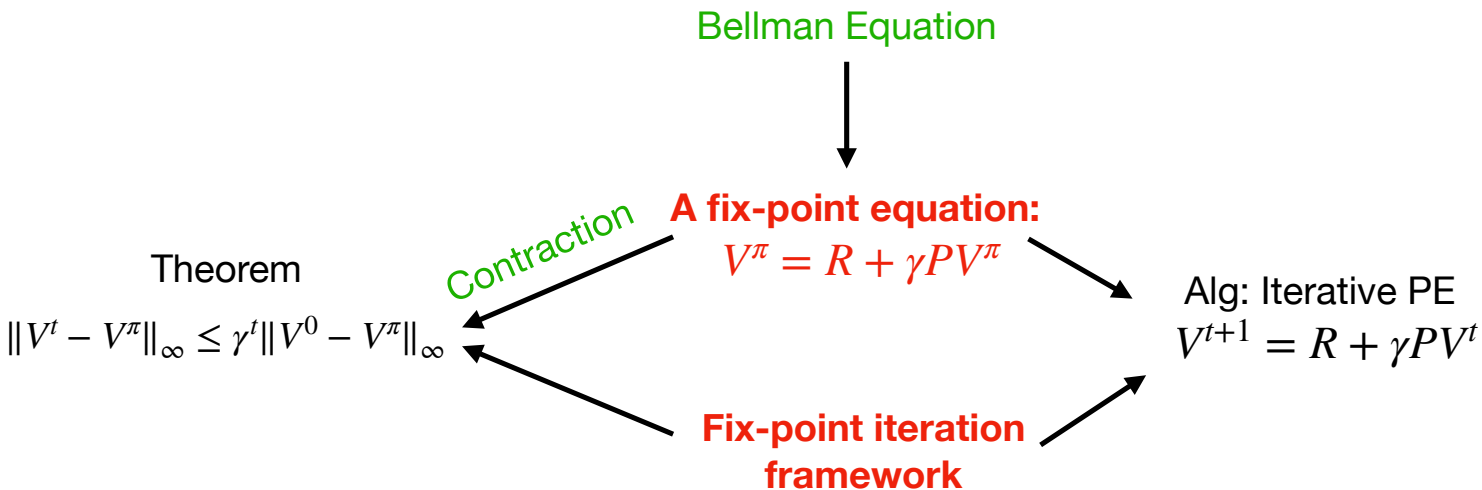
Summary

Key Question today: Given MDP \mathcal{M} , and a policy π , How to compute $V^\pi(s), \forall s$?



Summary

Key Question today: Given MDP \mathcal{M} , and a policy π , How to compute $V^\pi(s), \forall s$?



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

$$f: |f(x') - f(x)| \leq \gamma |x' - x|$$

$$V^{\pi} = \|R + \gamma P V - (R + \gamma P V')\|_{\infty} \leq \gamma \|V - V'\|_{\infty}$$

Next two lectures:

Given MDP \mathcal{M} , how to compute the optimal policy π^* , and V^*

$$\pi^* \Leftarrow \arg \max_{\pi} \left(r(s, \pi) + \gamma E_{s' \sim P_{\pi}(\cdot | s)} V^{\pi}(s') \right)$$

$V^{\pi} \rightarrow V^*$