

# **Policy Evaluation**

# Announcements

P1 is delayed, will be released on Monday

HW0 due today

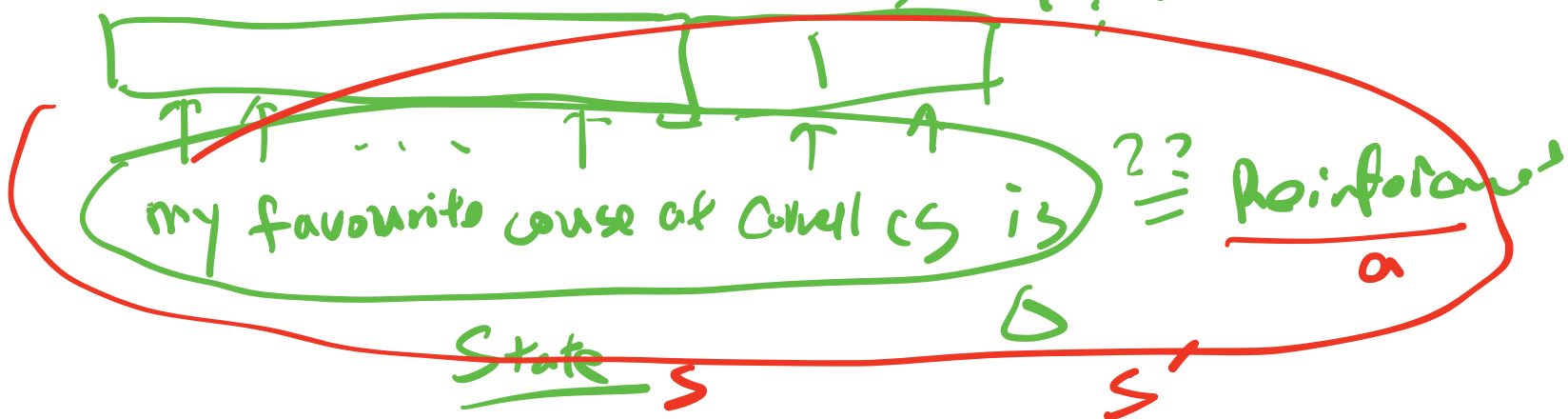
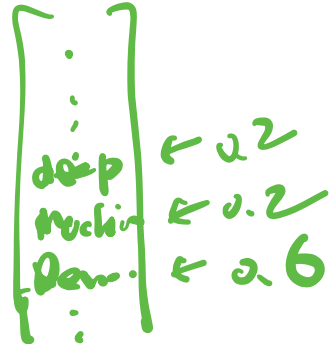
# Recap: Definitions

States → Actions

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

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

$$S' \sim P(\cdot | s, a)$$



# Recap: Definitions

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

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

$$\text{Policy } \pi : S \mapsto \Delta(A)$$

$$\underline{a \sim \pi(\cdot | s)}$$

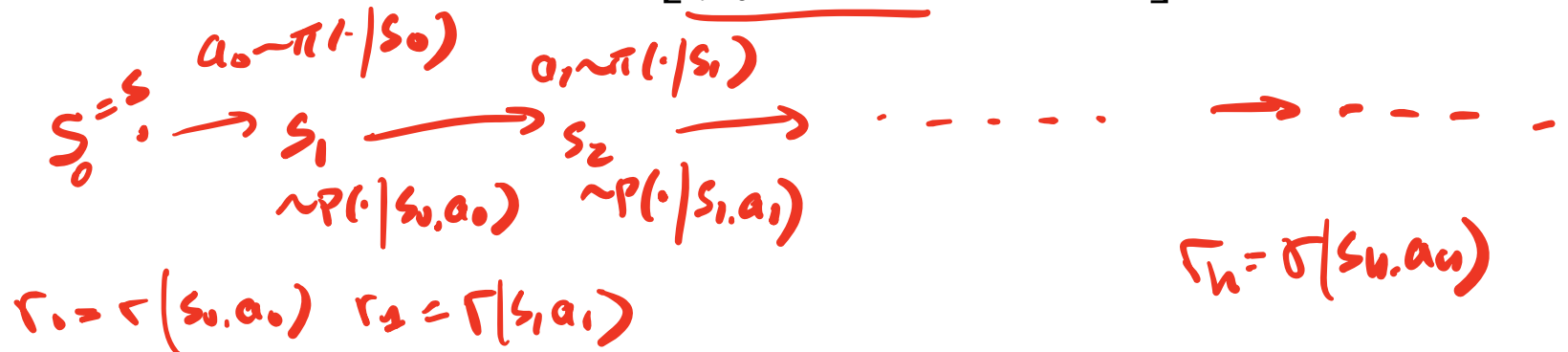
# Recap: Definitions

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

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

$$\text{Policy } \pi : S \mapsto \Delta(A)$$

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



# Recap: Definitions

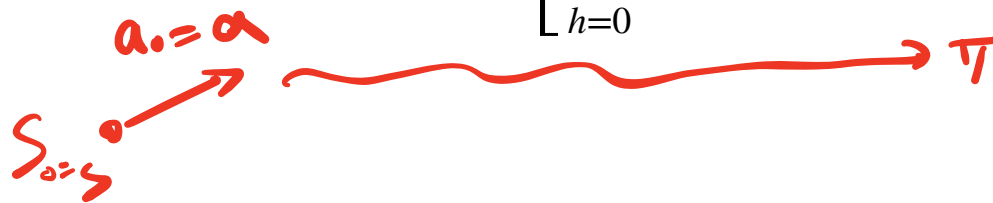
$$\mathcal{M} = \{S, A, P, r, \gamma\}$$

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

$$\text{Policy } \pi : S \mapsto \Delta(A)$$

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

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



# Recap: Bellman equation (consistency)

## Bellman Eq

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

$a \sim \pi(\cdot|s)$

$s \rightarrow s' \sim P(\cdot|s, a)$

$V^\pi(s')$

$\mathbb{E}_{s' \sim P(\cdot|s, a)}$

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

# Recap: Bellman equation (consistency)

## Bellman Eq

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

## Relationship between V and Q

Exercise: can you write  $V^\pi$  using  $Q^\pi$ , and then  $Q^\pi$  using  $V^\pi$

$$\begin{aligned} & \text{Handwritten: } V^\pi(s) \quad ?? \quad \mathbb{E}_{a \sim \pi(\cdot|s)} Q^\pi(s, a) = V^\pi(s) \\ & \text{Handwritten: } Q^\pi(s, a) \quad ?? \\ & \text{Handwritten: } = r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^\pi(s') \end{aligned}$$



# Today: Policy Evaluation

**Key Question:**

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

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

# Motivation for Policy Evaluation



We want to **evaluate** our strategy against some fixed opponent (AlphaGo constantly estimates the current probability of winning)

# Motivation for Policy Evaluation



We want to **evaluate** our strategy against some fixed opponent (AlphaGo constantly estimates the current probability of winning)



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

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

# Exact Policy Evaluation $\mathbb{E}_{a \sim \pi(\cdot|s)}$

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

# Exact Policy Evaluation

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

We know that for  $V^\pi$ , we have **Bellman equation:**

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

$$\underline{a \in \pi(s)}$$

# Exact Policy Evaluation

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

We know that for  $V^\pi$ , we have **Bellman equation**:

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

This gives us  $S$  many **linear** constraints



# Exact Policy Evaluation

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

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

# Exact Policy Evaluation

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

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

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

$$V = \begin{bmatrix} \vdots \\ V(s) \\ \vdots \end{bmatrix} \in \mathbb{R}^{|S|}$$

$$R = \begin{bmatrix} \vdots \\ r(s, \pi(s)) \\ \vdots \end{bmatrix} \in \mathbb{R}^{|S|}$$

$$P = \begin{bmatrix} \vdots \\ P_{s,s'} \\ \vdots \end{bmatrix} \in \mathbb{R}^{|S| \times |S|}$$

# Exact Policy Evaluation

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

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

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

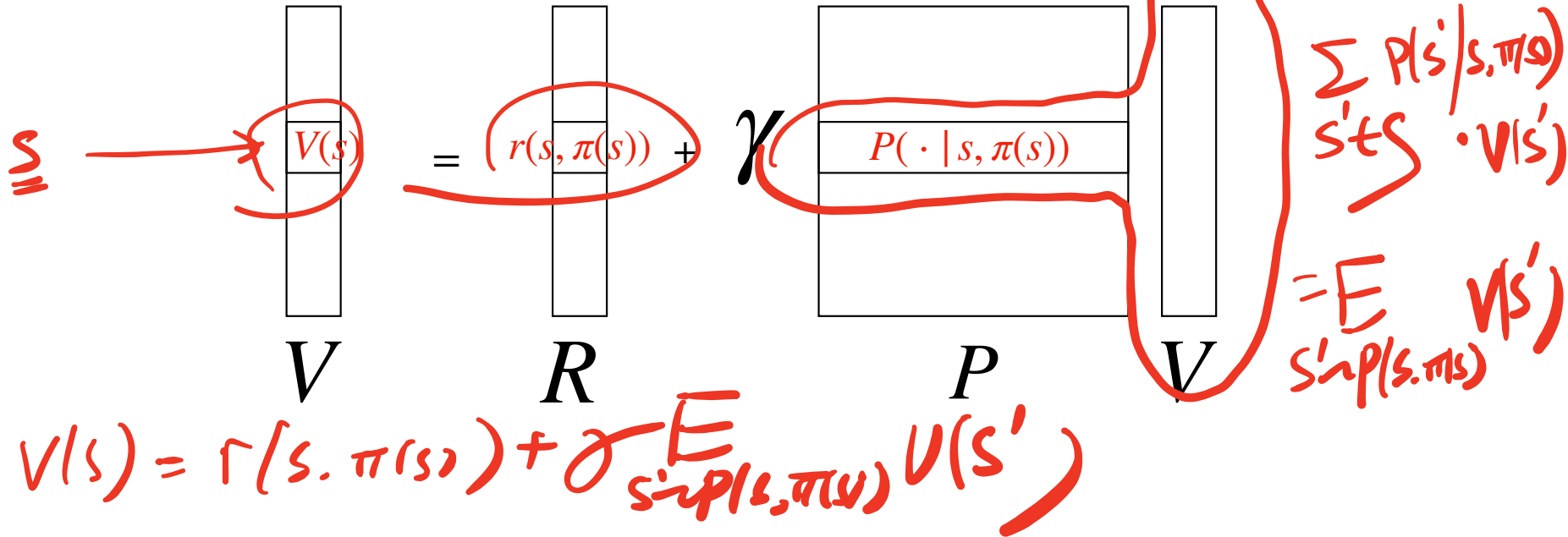
we can **combine all  $\mathcal{S}$  many constraints together:**

$$V = R + \gamma PV$$

# 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:

$$V = R + \gamma PV$$



## Exact Policy Evaluation

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

Since  $V = r + \gamma PV$ , we can obtain  $V$  as:

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

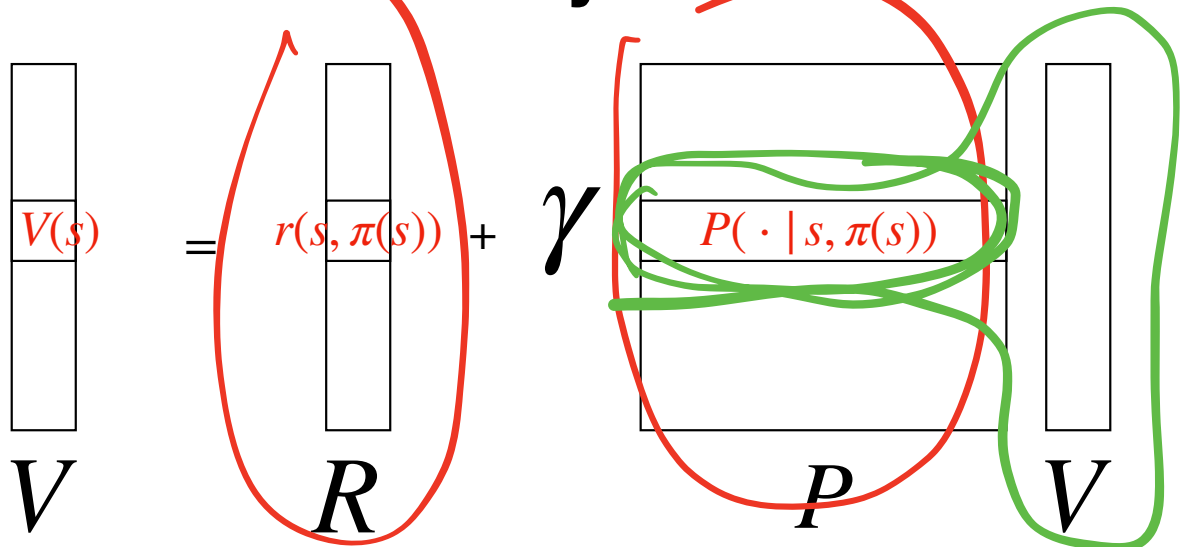
# Exact Policy Evaluation

Since  $V = r + \gamma PV$ , we can obtain  $V$  as:

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

See the assigned AJKS section for proving  $(I - \gamma P)$  is full rank (thus invertible)

# Summary so far:



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

$E_{s \sim P(\cdot | s, \pi(s))} \cdot V(s)$

## Summary so far:

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

$S \times S$

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

$S = 10,000$

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



# Outline:

$$V \approx V^\pi$$

 1. Exact Policy Evaluation

2. Approximate Policy Evaluation via an Iterative Algorithm

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

$V^\pi$  is a fix-point solution

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

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

$$\underline{R + \gamma P \cdot V^\pi} \Rightarrow V^\pi$$

$V^\pi$  is a fix-point solution

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

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

$$V^0 \quad V^1 = R + \gamma P V^0 \rightarrow V^2 = R + \gamma P V^1$$

$$\|V^{t+1} - V^t\|_2 \approx 0$$

# Iterative Policy Evaluation:

## Algorithm (Iterative PE):

Start with some initialization  $V^0 \in \mathbb{R}^{|S|}$ , repeat for  $t = 0 \dots$ :

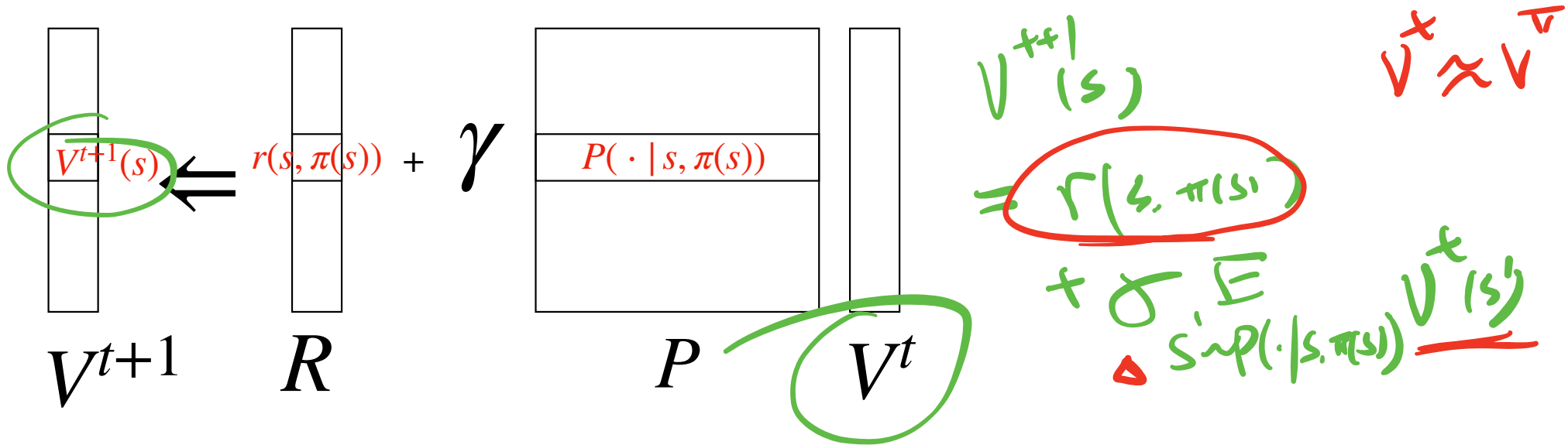
$$V^{t+1} \leftarrow R + \gamma P V^t$$

# Iterative Policy Evaluation:

## Algorithm (Iterative PE):

Start with some initialization  $V^0 \in \mathbb{R}^{|S|}$ , repeat for  $t = 0 \dots$ :

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

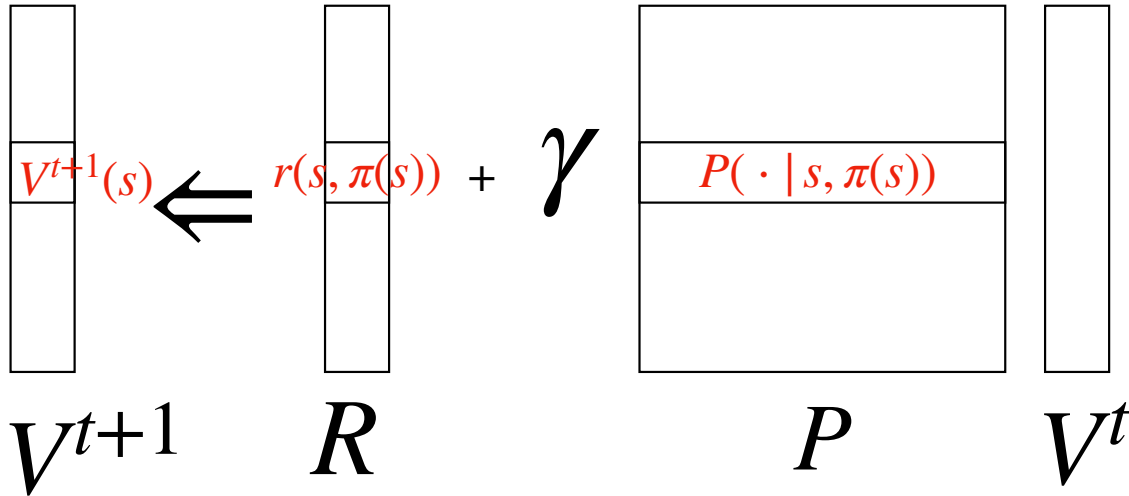


# Iterative Policy Evaluation:

## Algorithm (Iterative PE):

Start with some initialization  $V^0 \in \mathbb{R}^{|S|}$ , repeat for  $t = 0 \dots$ :

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



Q: What's computation complexity per iteration?

# Convergence of Iterative PE

## Theorem:

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

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

$$\max_s |V^0(s) - V^\pi(s)|$$

$$\underline{\gamma < 1}$$

$$\gamma = 0.99$$

$$(0.99)^{500} \approx 0$$

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

## 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:

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

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

*(Handwritten annotations: A red underline under "Aig", a red arrow pointing to the underlined term, and green circles around the expectation terms in both the current and target value functions.)*

# Convergence of Iterative PE

## Theorem:

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

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

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

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

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

# Convergence of Iterative PE

## Theorem:

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

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

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

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

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

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

# Convergence of Iterative PE

## Theorem:

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

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

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

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

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

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

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

# Convergence of Iterative PE

## Theorem:

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

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

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

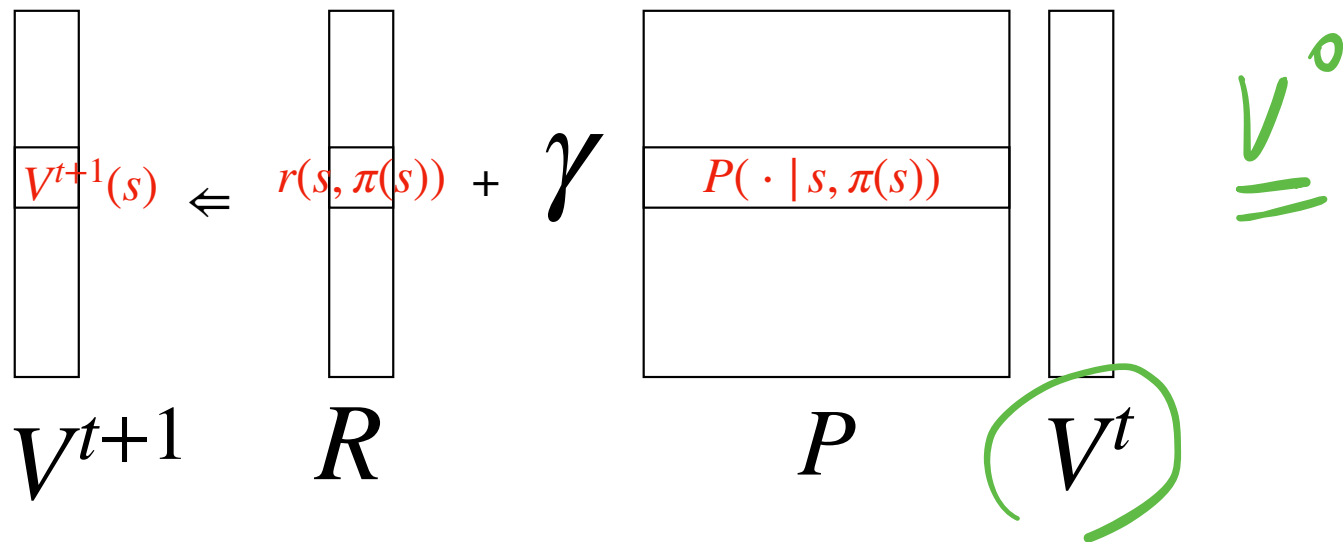
$$= \left| \underbrace{r(s, \pi(s))} + \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi(s))} V^t(s') - \left( \underbrace{r(s, \pi(s))} + \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi(s))} V^\pi(s') \right) \right|$$

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

$$\leq \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi(s))} \left| V^t(s') - V^\pi(s') \right| \leq \gamma \left\| V^t - V^\pi \right\|_\infty \Rightarrow \left\| V^{t+1} - V^\pi \right\|_\infty \leq \gamma \left\| V^t - V^\pi \right\|_\infty$$

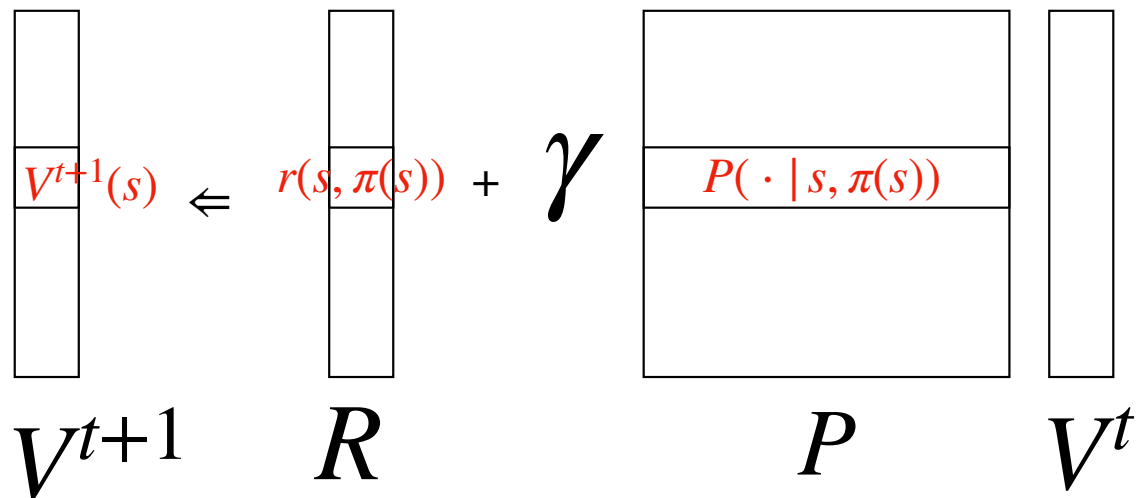
$$\leq \gamma^2 \left\| V^{t-1} - V^\pi \right\|_\infty$$

## Summary so far:



$$\|V^{t+1} - V^t\|_{\infty} \approx 0$$

## Summary so far:



Convergence:

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

# Extension to stochastic policy and Q functions

Your homework:

How to modify the two algorithms so that it can handle stochastic policy and learn Q functions

$$Q^T(s,a) = r(s,a) + \gamma \mathbb{E}_{s' \sim p(s'|s,a)} V^T(s')$$



# Summary

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

# Summary

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

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

# Summary

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

1. The **exact** algorithm  $V = (I - \gamma P)^{-1}R$  requires matrix inverse (computation complexity at least  $O(S^3)$ )
2. Iterative algorithm can quickly find an **approximate** solution (error shrinks in the rate of  $\gamma^t$ )

# Summary

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

1. The **exact** algorithm  $V = (I - \gamma P)^{-1}R$  requires matrix inverse (computation complexity at least  $O(S^3)$ )
2. Iterative algorithm can quickly find an **approximate** solution (error shrinks in the rate of  $\gamma^t$ )

(We will see many similar iterative algorithms later)