# Policy Evaluation

# Announcements

P1 is delayed, will be released on Monday

HW0 due today

# Recap: Definitions

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

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

Policy $\pi : S \mapsto \Delta(A)$

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

# Recap: Bellman equation (consistency)

**Bellman Eq**

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

**Relationship between V and Q**

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

# 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)



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

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 $s \in S$

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

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

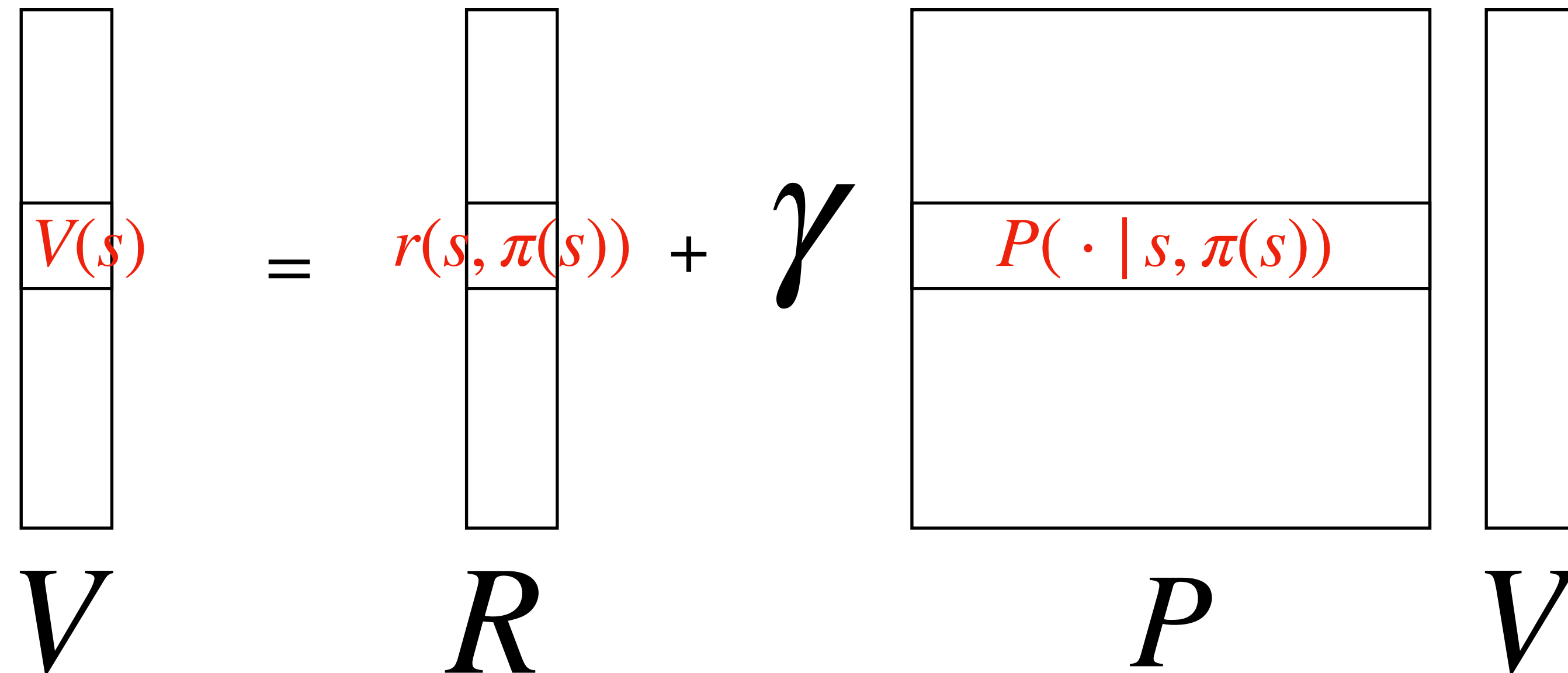we can **combine all $S$ many constraints together:**

$$V = R + \gamma PV$$

# 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:

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

# Exact Policy Evaluation

Since $V = r + \gamma P V$, 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:

$$V(s) = r(s, \pi(s)) + \gamma \, P(\,\cdot\,|\,s, \pi(s))$$

$$V \qquad\qquad R \qquad\qquad\qquad P \qquad V$$

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

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

# Outline:

✅ 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

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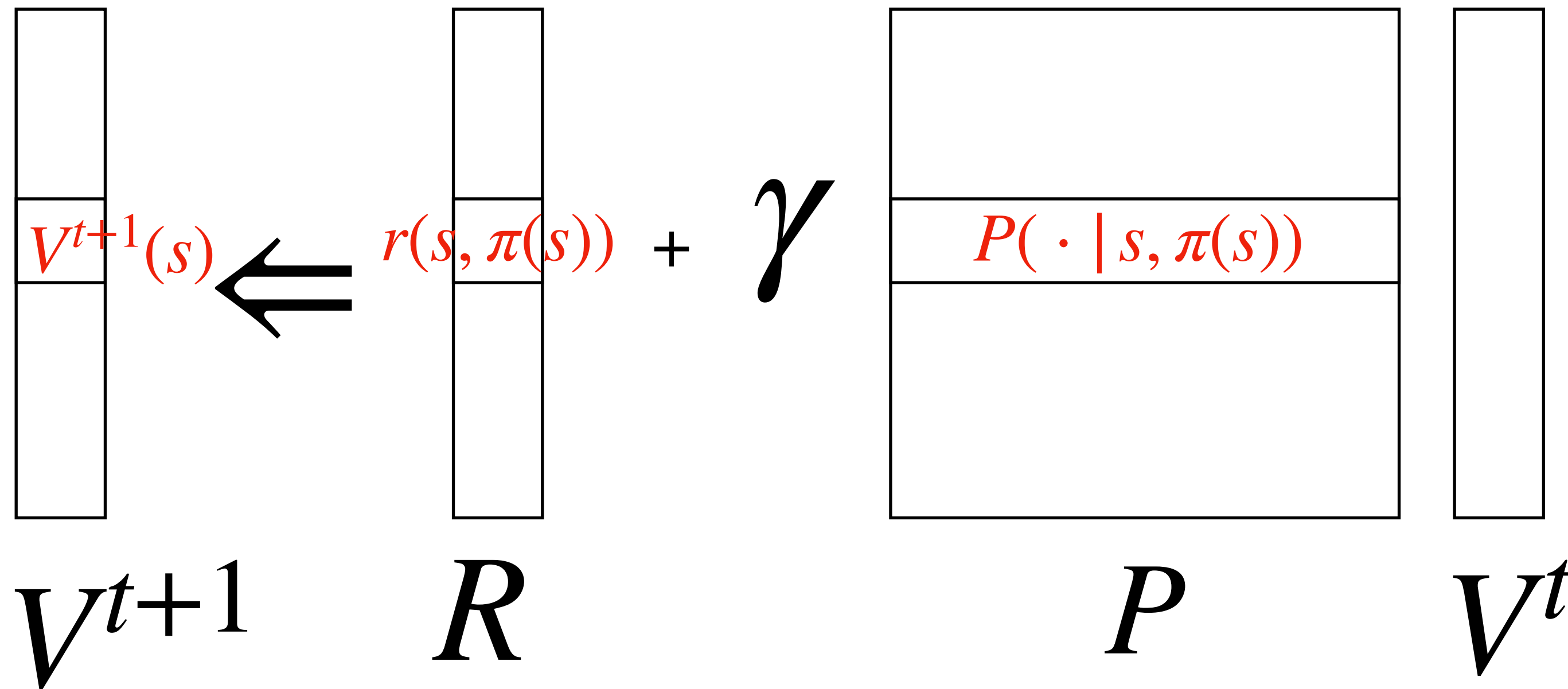
$$V^\pi = R + \gamma P V^\pi$$

# Iterative Policy Evaluation:

**Algorithm (Iterative PE):**
Start with some initialization $V^0 \in \mathbb{R}^{|S|}$, **repeat for** $t = 0\dots$:
$$V^{t+1} \Leftarrow R + \gamma P V^t$$

$$V^{t+1}(s) \Leftarrow r(s, \pi(s)) + \gamma \; P(\cdot \mid s, \pi(s)) \; V^t$$

$$V^{t+1} \qquad R \qquad\qquad P \qquad V^t$$

Q: What's computation complexity per iteration?
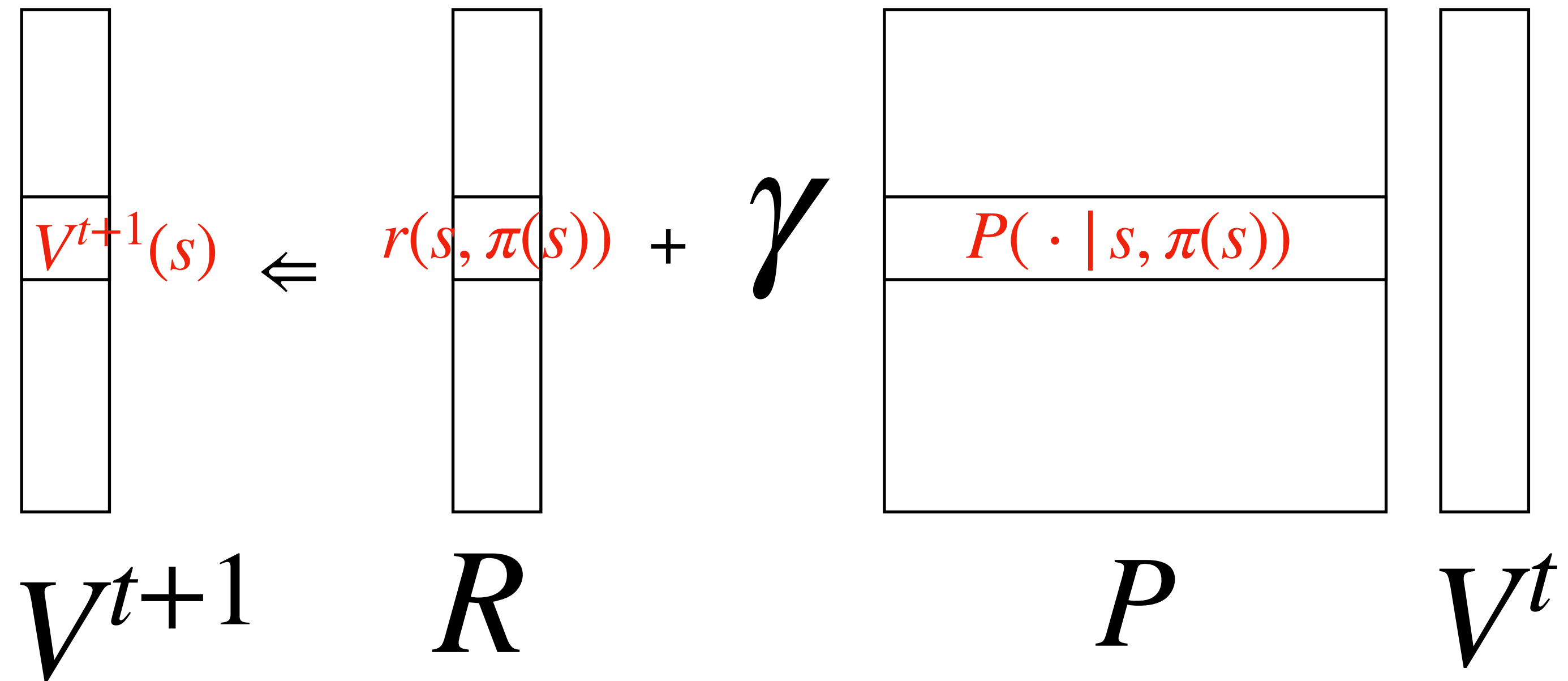
# Convergence of Iterative PE

**Theorem:**

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

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

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

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

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

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

# Summary so far:



$$V^{t+1}(s) \Leftarrow r(s, \pi(s)) + \gamma \quad P(\,\cdot\,|\,s, \pi(s))$$

$$V^{t+1} \qquad R \qquad\qquad P \qquad V^t$$

## Convergence:

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

# 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

# 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)