

# Policy Iteration

## Recap: Policy Evaluation

Given a policy  $\pi : S \mapsto A$ , compute  $V^\pi$ :

# Recap: Policy Evaluation

Given a policy  $\pi : S \mapsto A$ , compute  $V^\pi$ :

1. Solve a linear system:

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

# Recap: Policy Evaluation

Given a policy  $\pi : S \mapsto A$ , compute  $V^\pi$ :

## 1. Solve a linear system:

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

## 2. Fix-point iteration

$$\forall s : V^{t+1}(s) \leftarrow r(s, \pi(s)) + \gamma \sum_{s'} P(s'|s, \pi(s))V^t(s')$$

# Recap: Policy Evaluation

Given a policy  $\pi : S \mapsto A$ , compute  $V^\pi$ :

## 1. Solve a linear system:

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

## 2. Fix-point iteration

$$\forall s : V^{t+1}(s) \leftarrow r(s, \pi(s)) + \gamma \sum_{s'} P(s'|s, \pi(s))V^t(s')$$

Q: once we get  $V^\pi$ , how to get  $Q^\pi$ ?

$$Q^\pi(s, a) = r(s, a) + \gamma \sum_{s'} P(s'|s, \pi(s)) V^\pi(s')$$

# Recap: Value Iteration

1. VI

(a fix point iteration alg):

$$Q^{t+1} \leftarrow \mathcal{T} Q^t \quad \alpha^* = \mathcal{T} \alpha^*$$

$$\forall s, a \quad Q^{t+1}(s, a) \leftarrow r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a' \in A} Q^t(s', a')$$

## Recap: Value Iteration

1. VI

(a fix point iteration alg):

$$Q^{t+1} \leftarrow \mathcal{T} Q^t$$

Contraction



2. VI convergence: exponentially fast,

$$\text{i.e., } \|Q^t - Q^*\|_\infty \leq \gamma^t \|Q^0 - Q^*\|_\infty$$

## Recap: Value Iteration

1. VI

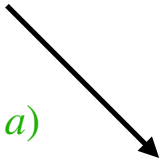
(a fix point iteration alg):

$$Q^{t+1} \leftarrow \mathcal{T} Q^t$$

Contraction



2. VI convergence: exponentially fast,  
i.e.,  $\|Q^t - Q^*\|_\infty \leq \gamma^t \|Q^0 - Q^*\|_\infty$



$$\pi^t(s) := \arg \max_a Q^t(s, a)$$



# Recap: Value Iteration

1. VI

(a fix point iteration alg):

$$Q^{t+1} \leftarrow \mathcal{T} Q^t$$

Contraction

$$\pi^t(s) := \arg \max_a Q^t(s, a)$$

3. Policy Performance:  $V^{\pi^t}(s) \geq V^*(s) - \frac{2\gamma^t}{1-\gamma} \|Q^0 - Q^*\|_\infty \forall s \in S$

2. VI convergence: exponentially fast,  
i.e.,  $\|Q^t - Q^*\|_\infty \leq \gamma^t \|Q^0 - Q^*\|_\infty$

$$Q^0 \in [0, \frac{1}{1-\gamma}]$$

$$Q^* \in [0, \frac{1}{1-\gamma}]$$

# Recap: Value Iteration

1. VI

(a fix point iteration alg):

$$Q^{t+1} \leftarrow \mathcal{T} Q^t$$

Contraction

2. VI convergence: exponentially fast,

$$\text{i.e., } \|Q^t - Q^*\|_\infty \leq \gamma^t \|Q^0 - Q^*\|_\infty$$

$$\pi^t(s) := \arg \max_a Q^t(s, a)$$

3. Policy Performance:  $V^{\pi^t}(s) \geq V^*(s) - \frac{2\gamma^t}{1-\gamma} \|Q^0 - Q^*\|_\infty \forall s \in S$

Note that  $Q^t \in \mathbb{R}^{|S||A|}$  is our estimator from VI,  
it does not correspond a  $Q^{\pi^t}$ !

$$Q^{\pi^t} \neq Q^t$$

## Question for Today:

~~(approximately)~~

Given an MDP  $\mathcal{M} = (S, A, P, r, \gamma)$ , How to find  $\pi^* : S \mapsto A$

# Outline:

1: An Iterative Algorithm: Policy Iteration

2: Convergence? How fast?

3: A new model: Finite horizon MDP

# Algorithm: Policy Iteration

↳  $\{\pi^0, \pi^1, \dots, \pi^T\}$

1. Initialization:  $\pi^0 : S \mapsto \Delta(A)$
2. For  $t = 0 \dots$ ,

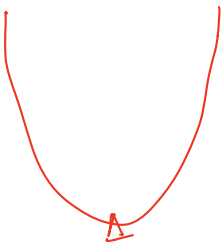
# Algorithm: Policy Iteration

1. Initialization:  $\pi^0 : S \mapsto \Delta(A)$

2. For  $t = 0 \dots$ ,

*linear program to compute  $V^{\pi^t} \Rightarrow Q^{\pi^t}$*

3. **Policy Evaluation:**  $Q^{\pi^t}(s, a), \forall s, a$



# Algorithm: Policy Iteration

$$Q^t \Rightarrow \pi^t = \arg \max_a Q^t(s, a)$$

$$\Downarrow$$

$$Q^{\pi^t}$$

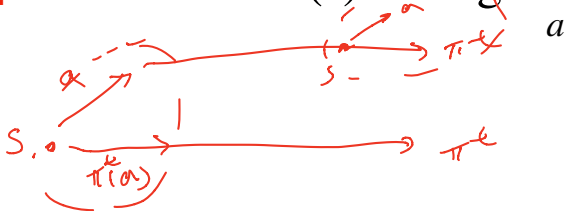
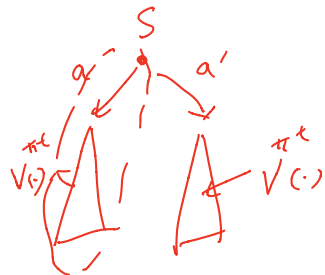
1. Initialization:  $\pi^0 : S \mapsto \Delta(A)$

2. For  $t = 0 \dots$ ,

3. **Policy Evaluation:**  $Q^{\pi^t}(s, a), \forall s, a$

4. **Policy Improvement**  $\pi^{t+1}(s) := \arg \max_a Q^{\pi^t}(s, a), \forall s$

Exact Alg



# Outline:

1: An Iterative Algorithm: Policy Iteration



2: Convergence? How fast?

$\epsilon$

3: A new model: Finite horizon MDP



# Key properties of Policy Iterations:

1. Monotonic improvement:

$$Q^{\pi^{t+1}}(s, a) \geq Q^{\pi^t}(s, a), \forall s, a$$

## Key properties of Policy Iterations:

1. Monotonic improvement:

$$Q^{\pi^{t+1}}(s, a) \geq Q^{\pi^t}(s, a), \forall s, a$$

2. Convergence:

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

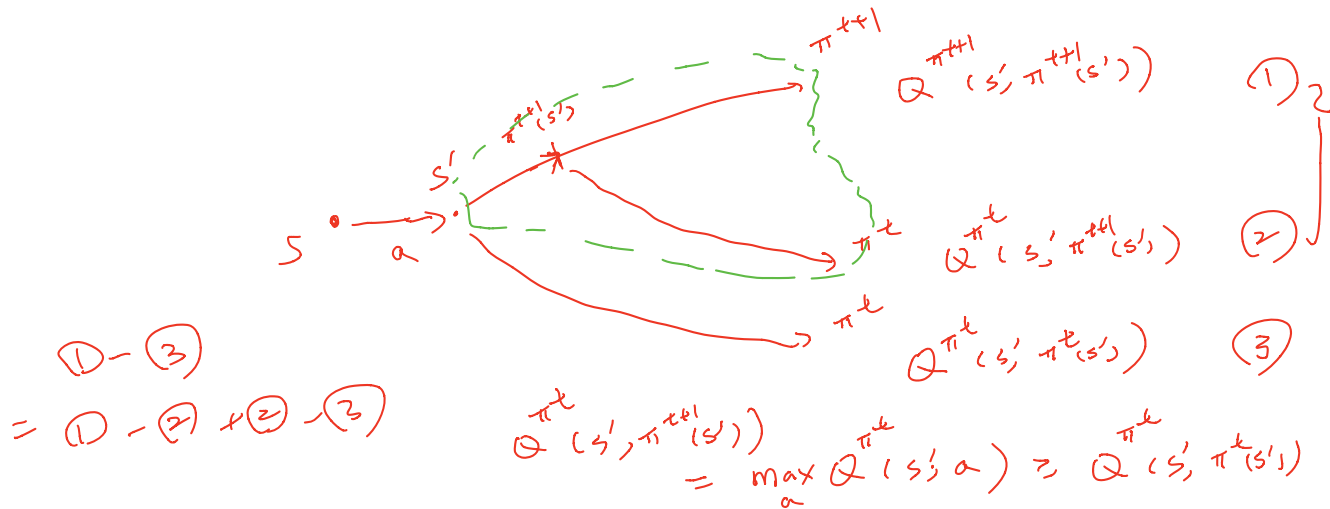
# Monotonic Improvement

Recall: Policy Improvement  $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \forall s$

# Monotonic Improvement

Recall: Policy Improvement  $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \forall s$

Let us compare  $Q^{\pi^{t+1}}(s, a)$  &  $Q^{\pi^t}(s, a)$ :



# Monotonic Improvement

Recall: Policy Improvement  $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \forall s$

Lemma: Monotonic improvement  $Q^{\pi^{t+1}}(s, a) \geq Q^{\pi^t}(s, a), \forall s, a$

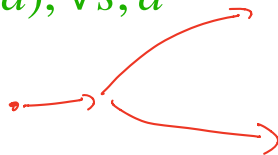
# Monotonic Improvement

Recall: Policy Improvement  $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \forall s$

Lemma: Monotonic improvement  $Q^{\pi^{t+1}}(s, a) \geq Q^{\pi^t}(s, a), \forall s, a$

*cancel  $r(s, a)$*

$$Q^{\pi^{t+1}}(s, a) - Q^{\pi^t}(s, a) = \gamma \mathbb{E}_{s' \sim P(s, a)} \left[ Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^t(s')) \right]$$

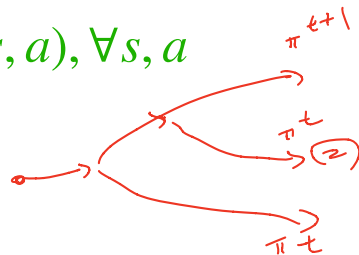


# Monotonic Improvement

Recall: Policy Improvement  $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \forall s$

Lemma: Monotonic improvement  $Q^{\pi^{t+1}}(s, a) \geq Q^{\pi^t}(s, a), \forall s, a$

$$Q^{\pi^{t+1}}(s, a) - Q^{\pi^t}(s, a) = \gamma \mathbb{E}_{s' \sim P(s, a)} \left[ Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^t(s')) \right]$$



$$= \gamma \mathbb{E}_{s' \sim P(s, a)} \left[ \underbrace{Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^{t+1}(s'))}_{\text{circled 2}} + \underbrace{Q^{\pi^t}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^t(s'))}_{\text{circled 2}} \right]$$

apply recursion

$\geq 0$

# Monotonic Improvement

$\gamma \in [0, 1]$

Recall: Policy Improvement  $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \forall s$

Lemma: Monotonic improvement  $Q^{\pi^{t+1}}(s, a) \geq Q^{\pi^t}(s, a), \forall s, a$

$$\begin{aligned} Q^{\pi^{t+1}}(s, a) - Q^{\pi^t}(s, a) &= \gamma \mathbb{E}_{s' \sim P(s, a)} \left[ Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^t(s')) \right] \\ &= \gamma \mathbb{E}_{s' \sim P(s, a)} \left[ Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^{t+1}(s')) + \underbrace{Q^{\pi^t}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^t(s'))}_{\geq 0} \right] \\ &\geq \gamma \mathbb{E}_{s' \sim P(s, a)} \left[ \underbrace{Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^{t+1}(s'))}_{\geq \frac{1}{1-\gamma}} \right] \end{aligned}$$

$\frac{1}{1-\gamma} \leq$  recursion  $\leq \frac{1}{1-\gamma}$

$Q^{\pi^k}(s', \pi^{t+1}(s'))$   
 $= \max_a Q^{\pi^k}(s', a)$   
 $\geq Q^{\pi^k}(s', \pi^k(s'))$



# Monotonic Improvement

$$Q^\pi(s,a) \in [0, \frac{1}{1-\gamma}]$$

Recall: Policy Improvement  $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \forall s$

Lemma: Monotonic improvement  $Q^{\pi^{t+1}}(s, a) \geq Q^{\pi^t}(s, a), \forall s, a$

$\forall s, a$

$$\begin{aligned}
 Q^{\pi^{t+1}}(s, a) - Q^{\pi^t}(s, a) &= \gamma \mathbb{E}_{s' \sim P(s, a)} \left[ Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^t(s')) \right] \\
 &= \gamma \mathbb{E}_{s' \sim P(s, a)} \left[ Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^{t+1}(s')) + Q^{\pi^t}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^t(s')) \right] \\
 &\geq \gamma \mathbb{E}_{s' \sim P(s, a)} \left[ Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^{t+1}(s')) \right] \geq \dots \geq -\gamma^\infty / (1 - \gamma) = 0
 \end{aligned}$$

$\Rightarrow \gamma \mathbb{E}_{s' \sim P_{s', \pi^{t+1}}(s')}$

$\Rightarrow \gamma^2 \left( \mathbb{E}_{s''} \mathbb{E}_{s'''} \left[ Q^{\pi^{t+1}}(s'', \pi^{t+1}(s'')) - Q^{\pi^t}(s'', \pi^{t+1}(s'')) \right] \right)$

# Monotonic Improvement

Recall: Policy Improvement  $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \forall s$

Lemma: Monotonic improvement  $Q^{\pi^{t+1}}(s, a) \geq Q^{\pi^t}(s, a), \forall s, a$

$$\begin{aligned}
 Q^{\pi^{t+1}}(s, a) - Q^{\pi^t}(s, a) &= \gamma \mathbb{E}_{s' \sim P(s, a)} \left[ Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^t(s')) \right] \\
 &= \gamma \mathbb{E}_{s' \sim P(s, a)} \left[ Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^{t+1}(s')) + Q^{\pi^t}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^t(s')) \right] \\
 &\geq \gamma \mathbb{E}_{s' \sim P(s, a)} \left[ Q^{\pi^{t+1}}(s', \pi^{t+1}(s')) - Q^{\pi^t}(s', \pi^{t+1}(s')) \right] \geq \dots \geq -\gamma^\infty / (1 - \gamma) = 0
 \end{aligned}$$

$V^{\pi^{t+1}}(s) \geq V^{\pi^t}(s), \forall s, ??$

$$V(s) = Q(s, \pi^{t+1}(s)) \geq Q(s, \pi^t(s)) \geq Q(s, \pi^t(s))$$
  

$$\uparrow \arg \max_a Q(s, a)$$

# Convergence analysis via Monotonic Improvement

Recall: Policy Improvement  $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \forall s$

Theorem: Convergence  $\|V^{\pi^{t+1}} - V^*\|_\infty \leq \gamma \|V^{\pi^t} - V^*\|_\infty \leq \dots \leq \gamma^{t+1} \|V^{\pi^0} - V^*\|_\infty$

# Convergence analysis via Monotonic Improvement

Recall: Policy Improvement  $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \forall s$

Theorem: Convergence  $\|V^{\pi^{t+1}} - V^{\star}\|_{\infty} \leq \gamma \|V^{\pi^t} - V^{\star}\|_{\infty}$

$$V^{\star}(s) - V^{\pi^{t+1}}(s) = \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^{\star}(s') \right] - \left[ \underline{r(s, \pi^{t+1}(s))} + \gamma \mathbb{E}_{s' \sim P(s, \pi^{t+1}(s))} \underline{V^{\pi^{t+1}}(s')} \right]$$

# Convergence analysis via Monotonic Improvement

Recall: Policy Improvement  $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \forall s$

Theorem: Convergence  $\|V^{\pi^{t+1}} - V^*\|_\infty \leq \gamma \|V^{\pi^t} - V^*\|_\infty$

*monotonic  
improvement  
on  $V^{\pi^{t+1}}$*

$$\begin{aligned} V^*(s) - V^{\pi^{t+1}}(s) &= \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right] - \left[ r(s, \pi^{t+1}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^{t+1}(s))} V^{\pi^{t+1}}(s') \right] \\ &\stackrel{\Delta}{\leq} \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right] - \left[ r(s, \pi^{t+1}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^{t+1}(s))} V^{\pi^t}(s') \right] \end{aligned}$$

# Convergence analysis via Monotonic Improvement

Recall: Policy Improvement  $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \forall s$

Theorem: Convergence  $\|V^{\pi^{t+1}} - V^*\|_\infty \leq \gamma \|V^{\pi^t} - V^*\|_\infty$

$$\begin{aligned} V^*(s) - V^{\pi^{t+1}}(s) &= \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right] - \left[ r(s, \pi^{t+1}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^{t+1}(s))} V^{\pi^{t+1}}(s') \right] \\ &\leq \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right] - \left[ \underbrace{r(s, \pi^{t+1}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^{t+1}(s))} V^{\pi^t}(s')}_{Q^{\pi^t}(s, \pi^{t+1}(s))} \right] \\ &= \max_a \left( r(s, a) + \mathbb{E}_{s' \sim P(s, a)} \gamma V^*(s') \right) - \max_a \left( r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^{\pi^t}(s') \right) \end{aligned}$$

# Convergence analysis via Monotonic Improvement

Recall: Policy Improvement  $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \forall s$

Theorem: Convergence  $\|V^{\pi^{t+1}} - V^{\star}\|_{\infty} \leq \gamma \|V^{\pi^t} - V^{\star}\|_{\infty}$

$$\begin{aligned} V^{\star}(s) - V^{\pi^{t+1}}(s) &= \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^{\star}(s') \right] - \left[ r(s, \pi^{t+1}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^{t+1}(s))} V^{\pi^{t+1}}(s') \right] \\ &\leq \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^{\star}(s') \right] - \left[ r(s, \pi^{t+1}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^{t+1}(s))} V^{\pi^t}(s') \right] \\ &= \max_a \left( r(s, a) + \mathbb{E}_{s' \sim P(s, a)} \gamma V^{\star}(s') \right) - \max_a \left( r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^{\pi^t}(s') \right) \\ &\leq \max_a \left( \cancel{r(s, a)} + \gamma \mathbb{E}_{s' \sim P(s, a)} V^{\star}(s') - \left( \cancel{r(s, a)} + \gamma \mathbb{E}_{s' \sim P(s, a)} V^{\pi^t}(s') \right) \right) \end{aligned}$$

# Convergence analysis via Monotonic Improvement

Recall: Policy Improvement  $\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \forall s$

Theorem: Convergence  $\|V^{\pi^{t+1}} - V^*\|_\infty \leq \gamma \|V^{\pi^t} - V^*\|_\infty$

*Monotonic improvement*

$\forall s$

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

$$\begin{aligned}
 V^*(s) - V^{\pi^{t+1}}(s) &= \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right] - \left[ r(s, \pi^{t+1}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^{t+1}(s))} V^{\pi^{t+1}}(s') \right] \\
 &\leq \max_a \left[ r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') \right] - \left[ r(s, \pi^{t+1}(s)) + \gamma \mathbb{E}_{s' \sim P(s, \pi^{t+1}(s))} V^{\pi^t}(s') \right] \\
 &= \max_a \left( r(s, a) + \mathbb{E}_{s' \sim P(s, a)} \gamma V^*(s') \right) - \max_a \left( r(s, a) + \gamma \mathbb{E}_{s' \sim P(s, a)} V^{\pi^t}(s') \right) \\
 &\leq \max_a \left( \cancel{r(s, a)} + \gamma \mathbb{E}_{s' \sim P(s, a)} V^*(s') - \left( \cancel{r(s, a)} + \gamma \mathbb{E}_{s' \sim P(s, a)} V^{\pi^t}(s') \right) \right) \\
 &\leq \gamma \|V^* - V^{\pi^t}\|_\infty
 \end{aligned}$$

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



# Summary of Policy Iteration

Iterate between Policy Evaluation and Policy Improvement:

$$\pi^{t+1}(s) := \arg \max_a Q^{\pi^t}(s, a), \forall s$$

↑ Exact  $Q^{\pi^t}$

# Summary of Policy Iteration

Iterate between Policy Evaluation and Policy Improvement:

$$\pi^{t+1}(s) := \arg \max_a Q^{\pi^t}(s, a), \forall s$$

*↪ policy evaluation is expensive*

Monotonic improvement + convergence:

$$Q^{\pi^{t+1}}(s, a) \geq Q^{\pi^t}(s, a), \forall s, a \rightarrow V^{\pi^{t+1}}(s) \geq V^{\pi^t}(s), \forall s$$

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

*Δ*

# Value Iteration vs Policy Iteration

How many iterations (computation complexity) need to find the ~~the~~ EXACT optimal policy?

$$Q^{\pi^{t+1}}(s,a) = Q^{\pi^t}(s,a) \quad \forall s,a$$

We will explore this problem in HW1

# Outline:

1: An Iterative Algorithm: Policy Iteration



2: Convergence? How fast?



3: A new model: Finite horizon MDP

# Finite horizon Markov Decision Process

$$\mathcal{M} = \{S, A, r, P, H, \mu_0\},$$
$$r : S \times A \mapsto [0,1], H \in \mathbb{N}^+, P : S \times A \mapsto \Delta(S), s_0 \sim \mu_0$$

# Finite horizon Markov Decision Process

$$\mathcal{M} = \{S, A, r, P, H, \mu_0\},$$
$$r : S \times A \mapsto [0,1], H \in \mathbb{N}^+, P : S \times A \mapsto \Delta(S), s_0 \sim \mu_0$$

i.e., the task always starts from  $s_0 \sim \mu_0$ , and lasts for  $H$  total steps

# Finite horizon Markov Decision Process

$$\mathcal{M} = \{S, A, r, P, H, \mu_0\},$$
$$r : S \times A \mapsto [0,1], H \in \mathbb{N}^+, P : S \times A \mapsto \Delta(S), s_0 \sim \mu_0$$

i.e., the task always starts from  $s_0 \sim \mu_0$ , and lasts for  $H$  total steps

Very common in control,  
e.g., keep tracking a pre-specified trajectory with fixed length and fixed initial state









# Finite horizon Markov Decision Process

$$\mathcal{M} = \{S, A, r, P, H, \mu_0\},$$
$$r : S \times A \mapsto [0,1], H \in \mathbb{N}^+, P : S \times A \mapsto \Delta(S), s_0 \sim \mu_0$$

# Finite horizon Markov Decision Process

$$\mathcal{M} = \{S, A, r, P, H, \mu_0\},$$
$$r : S \times A \mapsto [0,1], H \in \mathbb{N}^+, P : S \times A \mapsto \Delta(S), s_0 \sim \mu_0$$

Note that in finite horizon setting, we will consider **time-dependent policies**, i.e.,

$$\pi := \{\pi_0, \pi_1, \dots, \pi_{H-1}\}, \pi_h : S \mapsto A, \forall h$$

# Finite horizon Markov Decision Process

$$\mathcal{M} = \{S, A, r, P, H, \mu_0\},$$
$$r : S \times A \mapsto [0,1], H \in \mathbb{N}^+, P : S \times A \mapsto \Delta(S), s_0 \sim \mu_0$$

Note that in finite horizon setting, we will consider **time-dependent policies**, i.e.,

$$\pi := \{\pi_0, \pi_1, \dots, \pi_{H-1}\}, \pi_h : S \mapsto A, \forall h$$

Policy interacts with the MDP as follows:

$$\tau = \{s_0, a_0, s_1, a_1, \dots, \cancel{s_H, a_H}\}, s_0 \sim \mu_0, a_0 = \pi_0(s_0), s_1 \sim P(\cdot | s_0, a_0), a_1 = \pi_1(s_1), \dots$$

*Handwritten annotations:*  
A red line is drawn under the sequence  $s_1, a_1$ .  
A red triangle is drawn under  $a_0$ .  
A red triangle is drawn under  $a_1$ .

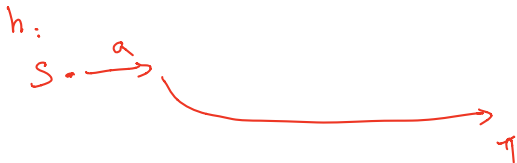
# V/Q functions in Finite horizon MDP

$$\pi = \{\pi_0, \pi_1, \dots, \pi_{H-1}\}$$

$$V_h^\pi(s) = \mathbb{E} \left[ \sum_{\tau=h}^{H-1} r(s_\tau, a_\tau) \mid s_h = s, a_\tau = \pi_\tau(s_\tau), s_{\tau+1} \sim P(\cdot | s_\tau, a_\tau) \right]$$

*no  $\delta$  any more*

$$Q_h^\pi(s, a) = \mathbb{E} \left[ \sum_{\tau=h}^{H-1} r(s_\tau, a_\tau) \mid (s_h, a_h) = (s, a), a_\tau = \pi_\tau(s_\tau), P \right]$$



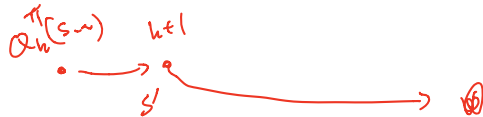
# V/Q functions in Finite horizon MDP

$$V_h^\pi(s) = \mathbb{E} \left[ \sum_{\tau=h}^{H-1} r(s_\tau, a_\tau) \mid s_h = s, a_\tau = \pi_\tau(s_\tau), s_{\tau+1} \sim P(\cdot | s_\tau, a_\tau) \right]$$

$$Q_h^\pi(s, a) = \mathbb{E} \left[ \sum_{\tau=h}^{H-1} r(s_\tau, a_\tau) \mid (s_h, a_h) = (s, a), a_\tau = \pi_\tau(s_\tau), P \right]$$

Bellman Equation:

$$Q_h^\pi(s, a) = r(s, a) + \mathbb{E}_{s' \sim P(\cdot | s, a)} [V_{h+1}^\pi(s')]$$



$$Q_h^\pi(s, a) = r(s, a) + \mathbb{E} \left[ V_{h+1}^\pi(s') \right]$$

# Compute Optimal Policy via DP

$$\pi^* = \{\pi_0^*, \pi_1^*, \dots, \pi_{H-1}^*\}$$



# Compute Optimal Policy via DP

$$\pi^* = \{\pi_0^*, \pi_1^*, \dots, \pi_{H-1}^*\}$$

We use Dynamic Programming, and do DP backward in time; start at  $H - 1$

# Compute Optimal Policy via DP

$$\pi^* = \{\pi_0^*, \pi_1^*, \dots, \pi_{H-1}^*\}$$

We use Dynamic Programming, and do DP backward in time; start at  $H - 1$

$$Q_{H-1}^*(s, a) = r(s, a)$$

# Compute Optimal Policy via DP

$$\pi^* = \{\pi_0^*, \pi_1^*, \dots, \pi_{H-1}^*\}$$

We use Dynamic Programming, and do DP backward in time; start at  $H - 1$

$$Q_{H-1}^*(s, a) = r(s, a) \quad \pi_{H-1}^*(s) = \arg \max_a Q_{H-1}^*(s, a)$$

# Compute Optimal Policy via DP

$$\pi^* = \{\pi_0^*, \pi_1^*, \dots, \pi_{H-1}^*\}$$

We use Dynamic Programming, and do DP backward in time; start at  $H - 1$

$$Q_{H-1}^*(s, a) = r(s, a) \quad \pi_{H-1}^*(s) = \arg \max_a Q_{H-1}^*(s, a)$$

$$V_{H-1}^* = \max_a Q_{H-1}^*(s, a) = Q_{H-1}^*(s, \pi_{H-1}^*(s))$$

# Compute Optimal Policy via DP

$$\pi^* = \{\pi_0^*, \pi_1^*, \dots, \pi_{H-1}^*\}$$

We use Dynamic Programming, and do DP backward in time; start at  $H - 1$

$$Q_{H-1}^*(s, a) = r(s, a) \quad \pi_{H-1}^*(s) = \arg \max_a Q_{H-1}^*(s, a)$$

$$V_{H-1}^* = \max_a Q_{H-1}^*(s, a) = Q_{H-1}^*(s, \pi_{H-1}^*(s))$$

Now assume that we have already computed  $V_{h+1}^*$ ,  $h \leq H - 2$   
(i.e., we know how to perform optimally starting at  $h + 1$ )

# Compute Optimal Policy via DP

$$\pi^* = \{\pi_0^*, \pi_1^*, \dots, \pi_{H-1}^*\}$$

We use Dynamic Programming, and do DP backward in time; start at  $H - 1$

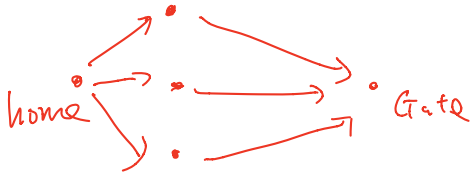
$$Q_{H-1}^*(s, a) = r(s, a) \quad \pi_{H-1}^*(s) = \arg \max_a Q_{H-1}^*(s, a)$$

$$V_{H-1}^* = \max_a Q_{H-1}^*(s, a) = Q_{H-1}^*(s, \pi_{H-1}^*(s))$$

Now assume that we have already computed  $V_{h+1}^*$ ,  $h \leq H - 2$   
(i.e., we know how to perform optimally starting at  $h + 1$ )

$$Q_h^*(s, a) = r(s, a) + \mathbb{E}_{s' \sim P(\cdot | s, a)} V_{h+1}^*(s')$$

# Compute Optimal Policy via DP



$$\pi^* = \{\pi_0^*, \pi_1^*, \dots, \pi_{H-1}^*\}$$

$$Q_{H-1}^*(s, a) = r(s, a) + \mathbb{E}_{s' \sim P_{sa}} V_{H-1}^*(s')$$

$$Q_{H-1}^*(s, a') = r(s, a') + \mathbb{E}_{s' \sim P_{sa'}} V_{H-1}^*(s')$$

We use Dynamic Programming, and do DP backward in time; start at  $H-1$

$$Q_{H-1}^*(s, a) = r(s, a) \quad \pi_{H-1}^*(s) = \arg \max_a Q_{H-1}^*(s, a)$$

$$V_{H-1}^* = \max_a Q_{H-1}^*(s, a) = Q_{H-1}^*(s, \pi_{H-1}^*(s))$$

Now assume that we have already computed  $V_{h+1}^*$ ,  $h \leq H-2$   
(i.e., we know how to perform optimally starting at  $h+1$ )



$$Q_h^*(s, a) = r(s, a) + \mathbb{E}_{s' \sim P(\cdot|s, a)} V_{h+1}^*(s') \quad \Rightarrow V_h^*(s)$$

$$\pi_h^*(s) = \arg \max_a Q_h^*(s, a)$$

# Summary on Finite horizon MDP

$$\mathcal{M} = \{S, A, r, P, H, \mu_0\},$$
$$r : S \times A \mapsto [0,1], H \in \mathbb{N}^+, P : S \times A \mapsto \Delta(S), s_0 \sim \mu_0$$

## Comparing to the infinite horizon MDP:

1. Policy will be time dependent
2. Value Iteration takes H steps (DP from  $H - 1 \rightarrow 0$ )
3. Compute exact  $\pi^*$  —no need to use  $\gamma^t$  argument
4. No more discount factor