

# RL: Basics

## Policy Iteration

Marius Lindauer



Winter Term 2021

# Policy Search (PS)

- ▶ One option is searching to compute the best policy
- ▶ Number of deterministic policies is  $|A|^{|S|}$
- ▶ Policy iteration is generally more efficient than enumeration

# MDP Policy Iteration (PI)

- ▶ Set  $i = 0$
- ▶ Initialize  $\pi_0(s)$  randomly for all states  $s$
- ▶ While  $i \neq 0$  or  $\|\pi_i - \pi_{i-1}\|_1 > 0$  (L1-norm, measures if the policy changed for any state)
  - ▶  $V^{\pi_i} \leftarrow$  MDP V-function policy evaluation of  $\pi$
  - ▶  $\pi_{i+1} \leftarrow$  Policy improvement
  - ▶  $i \leftarrow i + 1$

# Definition: State-Action Value Q

- ▶ State-action value of a policy

$$Q^{\pi}(s, a) = R(s, a) + \gamma \sum_{s' \in S} P(s' \mid s, a) V^{\pi}(s')$$

~> Take action  $a$ , then follow the policy  $\pi$

# Policy Improvement

- ▶ Compute state-action value of a policy  $\pi_i$ 
  - ▶ For  $s$  in  $S$  and  $a$  in  $A$ :

$$Q^\pi(s, a) = R(s, a) + \gamma \sum_{s' \in S} P(s' \mid s, a) V^\pi(s')$$

# Policy Improvement

- ▶ Compute state-action value of a policy  $\pi_i$ 
  - ▶ For  $s$  in  $S$  and  $a$  in  $A$ :

$$Q^\pi(s, a) = R(s, a) + \gamma \sum_{s' \in S} P(s' \mid s, a) V^\pi(s')$$

- ▶ Compute new policy  $\pi_{i+1}$  for all  $s \in S$

$$\pi_{i+1}(s) \in \arg \max_{a \in A} Q^{\pi_i}(s, a). \forall s \in S$$

# MDP Policy Iteration (PI)

- ▶ Set  $i = 0$
- ▶ Initialize  $\pi_0(s)$  randomly for all states  $s$
- ▶ While  $i == 0$  or  $\|\pi_i - \pi_{i-1}\|_1 > 0$  (L1-norm, measures if the policy changed for any state)
  - ▶  $V^{\pi_i} \leftarrow$  MDP V function policy **evaluation** of  $\pi$   $\rightsquigarrow$  use  $Q$
  - ▶  $\pi_{i+1} \leftarrow$  Policy **improvement**
  - ▶  $i \leftarrow i + 1$

## Delving Deeper Into Policy Improvement Step

$$\begin{aligned}Q^\pi(s, a) &= R(s, a) + \gamma \sum_{s' \in S} P(s' \mid s, a) V^\pi(s') \\ \max_a Q^\pi(s, a) &\geq R(s, a) + \gamma \sum_{s' \in S} P(s' \mid s, a) V^\pi(s') \\ \pi_{i+1}(s) &\in \arg \max_{a \in A} Q^{\pi_i}(s, a)\end{aligned}$$

- ▶ Suppose we take  $\pi_{i+1}(s)$  for one action, then follow  $\pi_i$  forever
  - ▶ Our expected sum of rewards is at least as good as if we had always followed  $\pi_i$
- ▶ But new proposed policy is to always follow  $\pi_{i+1}$



# Monotonic Improvement in Policy

- ▶ Definition

$$V^{\pi_2} \geq V^{\pi_1} : V^{\pi_2}(s) \geq V^{\pi_1}(s). \forall s \in \mathcal{S}$$

- ▶ Proposition:  $V^{\pi_{i+1}} \geq V^{\pi_i}$

- ▶ where  $\pi_{i+1}$  is the new policy we get from policy improvement on  $\pi_i$
- ▶ with strict inequality if  $\pi_i$  is suboptimal

# MDP Policy Iteration (PI): Check your Understanding

- ▶ Set  $i = 0$
- ▶ Initialize  $\pi_0(s)$  randomly for all states  $s$
- ▶ While  $i \neq 0$  or  $\|\pi_i - \pi_{i-1}\|_1 > 0$  (L1-norm, measures if the policy changed for any state)
  - ▶  $V^{\pi_i} \leftarrow$  MDP V function policy **evaluation** of  $\pi$   $\rightsquigarrow$  use  $Q$
  - ▶  $\pi_{i+1} \leftarrow$  Policy **improvement**
  - ▶  $i \leftarrow i + 1$
- ▶ If policy doesn't change, can it ever change?
- ▶ Is there a maximum of iterations of policy iteration?