



10-601 Introduction to Machine Learning

Machine Learning Department
School of Computer Science
Carnegie Mellon University

Reinforcement Learning

Matt Gormley
Lecture 26
April 13, 2018

Reminders

- **Homework 7: HMMs**
 - Out: Wed, Apr 04
 - Due: Mon, Apr 16 at 11:59pm
- **Schedule Changes**
 - Lecture on Fri, Apr 13
 - Recitation on Mon, Apr 23

VALUE ITERATION

Definitions for Value Iteration

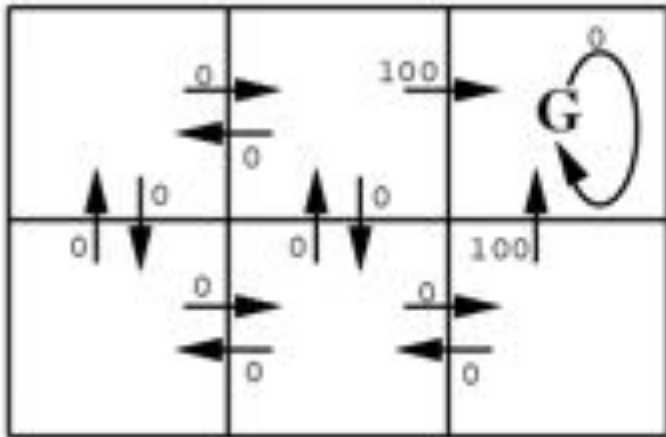
Whiteboard

- State trajectory
- Value function
- Bellman equations
- Optimal policy
- Optimal value function
- Computing the optimal policy
- Ex: Path Planning

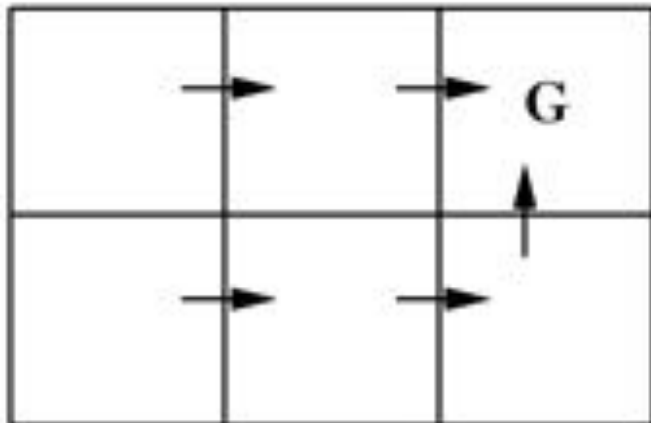
Example: Path Planning



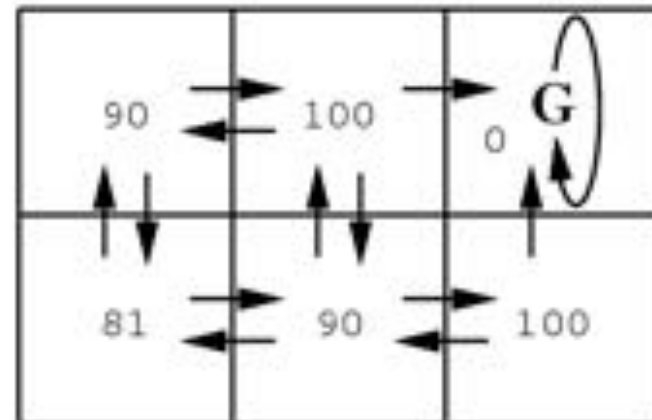
Example: Robot Localization



$r(s, a)$ (immediate reward) values



One optimal policy



$V^*(s)$ values

Value Iteration

Whiteboard

- Value Iteration Algorithm
- Synchronous vs. Asynchronous Updates
- Convergence Properties

Value Iteration

Algorithm 1 Value Iteration

```
1: procedure VALUEITERATION( $R(s, a)$  reward function,  $p(\cdot|s, a)$   
   transition probabilities)  
2:   Initialize value function  $V(s) = 0$  or randomly  
3:   while not converged do  
4:     for  $s \in \mathcal{S}$  do  
5:       for  $a \in \mathcal{A}$  do  
6:          $Q(s, a) = R(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'|s, a)V(s')$   
7:        $V(s) = \max_a Q(s, a)$   
8:   Let  $\pi(s) = \operatorname{argmax}_a Q(s, a)$ ,  $\forall s$   
9:   return  $\pi$ 
```

Policy Iteration

Whiteboard

- Policy Iteration Algorithm
- Solving the Bellman Equations for Fixed Policy
- Convergence Properties
- Value Iteration vs. Policy Iteration

Policy Iteration

Algorithm 1 Policy Iteration

- 1: **procedure** POLICYITERATION($R(s, a)$ reward function, $p(\cdot|s, a)$ transition probabilities)
- 2: Initialize policy π randomly
- 3: **while** not converged **do**
- 4: Solve Bellman equations for fixed policy π

$$V^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s' \in \mathcal{S}} p(s'|s, \pi(s)) V^\pi(s'), \quad \forall s$$

- 5: Improve policy π using new value function

$$\pi(s) = \operatorname{argmax}_a R(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'|s, a) V^\pi(s')$$

- 6: **return** π
-

Policy Iteration

Algorithm 1 Policy Iteration

1: **procedure** POLICYITERATION($R(s, a)$, γ , $p(\cdot|s, a)$
transition probabilities)

2: Initialize policy π randomly

3: **while** not converged **do**

4: Solve Bellman equations for fixed policy π

$$V^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s' \in \mathcal{S}} p(s'|s, \pi(s)) V^\pi(s'), \quad \forall s$$

5: Improve policy π using new value function

$$\pi(s) = \operatorname{argmax}_a R(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'|s, a) V^\pi(s')$$

6: **return** π

Compute value
function for fixed
policy is easy

System of $|\mathcal{S}|$
equations and $|\mathcal{S}|$
variables

Greedy policy
w.r.t. current
value function

Greedy policy might **remain the same** for a particular state if there is
no better action

Policy Iteration Convergence

In-Class Exercise:

How many policies are there for a finite sized state and action space?

In-Class Exercise:

Suppose policy iteration is shown to improve the policy at every iteration. Can you bound the number of iterations it will take to converge?

Value Iteration vs. Policy Iteration

- Value iteration requires $O(|A| |S|^2)$ computation per iteration
- Policy iteration requires $O(|A| |S|^2 + |S|^3)$ computation per iteration
- In practice, policy iteration converges in fewer iterations

Algorithm 1 Value Iteration

```
1: procedure VALUEITERATION( $R(s, a)$  reward function,  $p(\cdot|s, a)$ 
   transition probabilities)
2:   Initialize value function  $V(s) = 0$  or randomly
3:   while not converged do
4:     for  $s \in \mathcal{S}$  do
5:       for  $a \in \mathcal{A}$  do
6:          $Q(s, a) = R(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'|s, a) V(s')$ 
7:        $V(s) = \max_a Q(s, a)$ 
8:   Let  $\pi(s) = \operatorname{argmax}_a Q(s, a), \forall s$ 
9:   return  $\pi$ 
```

Algorithm 1 Policy Iteration

```
1: procedure POLICYITERATION( $R(s, a)$  reward function,  $p(\cdot|s, a)$ 
   transition probabilities)
2:   Initialize policy  $\pi$  randomly
3:   while not converged do
4:     Solve Bellman equations for fixed policy  $\pi$ 

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

5:     Improve policy  $\pi$  using new value function

$$\pi(s) = \operatorname{argmax}_a R(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'|s, a) V^\pi(s')$$

6:   return  $\pi$ 
```

Learning Objectives

Reinforcement Learning: Value and Policy Iteration

You should be able to...

1. Compare the reinforcement learning paradigm to other learning paradigms
2. Cast a real-world problem as a Markov Decision Process
3. Depict the exploration vs. exploitation tradeoff via MDP examples
4. Explain how to solve a system of equations using fixed point iteration
5. Define the Bellman Equations
6. Show how to compute the optimal policy in terms of the optimal value function
7. Explain the relationship between a value function mapping states to expected rewards and a value function mapping state-action pairs to expected rewards
8. Implement value iteration
9. Implement policy iteration
10. Contrast the computational complexity and empirical convergence of value iteration vs. policy iteration
11. Identify the conditions under which the value iteration algorithm will converge to the true value function
12. Describe properties of the policy iteration algorithm