

# Graduate Artificial Intelligence

## CS 640

# Value and Policy Iteration

Jeffrey Considine  
jconsidi@bu.edu

# Plan for Today

- Policies
- Value iteration methods
- Policy iteration methods

# Last Time

Re: specifying optimal values for Markov decision processes,

- What's the **best possible expected value** for this state?
- Specified **recursively** in terms of best possible expected values for other states.

Still open: how do we do this maximization?

# What is a Policy?

A policy is a function mapping states to actions.

- A deterministic policy returns a single action.
- A probabilistic policy returns a probability distribution of actions.

Usually denoted as  $\pi$  or with subscripts for context...

# Probabilistic vs Deterministic Policies

## Deterministic policy advantages

- Usually sufficient for optimal performance.
  - Exceptions with simultaneous adversarial choices.
- Much smaller to represent.
  - actions for states

## Probabilistic policy advantages

- Can always represent deterministic policies.
  - All probability on one action.
- Often convenient for numerical optimizations.

# Representing Policies for Finite MDPs

- Represent a deterministic policy as a table of actions.
- Represent a probabilistic policy as an matrix for states and actions.

# Evaluating Policies for MDPs (take 1)

Previously,

For a specific policy (optimal or not),

# Evaluating Policies for MDPs (take 2)

Rewrite

to



# Evaluating Policies for MDPs (take 3a)

Rewrite

with

# Evaluating Policies for MDPs (take 3b)

Rewrite

to

Rewriting with  $\pi$  and  $\gamma$  yields a Markov reward process.

And we already know how to solve for their values.

# What is an Optimal Policy?

So, what is ?

A policy is optimal if and only if

Optimal policies always exist. Is that surprising?

# Any Questions?

???

# Numerical Problems of Analytical Solution

Previously for Markov reward processes, we derived this solution.

- The interpretation of  $\mathbf{v}$  is the **expected number of -discounted visits** to each state.
- If  $\rho = 1$ , then **at least one state has an infinite number of visits**, and the matrix inversion will fail.
- This **breaks the calculation** even if those states have zero rewards.

# Computational Cost of Analytical Solution

If , then what is the computational cost of the analytical solution?

???

# Value Iteration

1. Set  $V$ .
2. For  $i = 0, 1, 2, \dots$ 
  - For all states  $s$ ,
  - Set  $V(s) = \max_a \sum_{s'} P(s'|s,a) V(s')$

What is  $V$ ?

How do we implement that expectation?

# Value Iteration as Fixed Horizon Evaluations

Prove: For integer  $n$ ,  $V_n$  is the best possible value achieved within  $n$  steps.

- Base case:  $V_0$  is the best possible value achieved within 0 steps.
  - 0 steps so 0 value.
- Inductive hypothesis: Assume that  $V_{n-1}$  is the best possible value achieved within  $n-1$  steps
- Induction step: Then  $V_n$  is the best possible value achieved within  $n$  steps.
  - This is true by construction.



# Bounding Distance Values

How close is a given  $s$  to  $s^*$ ?

How much discounted reward is remaining after  $t$  steps?

# How to Act given (Near) Optimal Values?

For a given  $V$ , pick the value maximizing action.

# Any Questions?

???

# Policy Iteration

- Policy iteration methods iterate on policies instead of value function.
- Similar intuitions but often converging in fewer iterations.
- Basic idea:
  - Start with any policy .
  - Calculate .
  - Calculate the best assuming will be used for one step and followed by using .

# Naïve Implementations of Policy Iteration

1. Initialize to any policy.
2. For  $i = 0, 1, 2, \dots$ 
  - a) Compute using value iteration.
  - b) Set

Ignoring cost of 2a, why might this be better than value iteration?

# Incremental Value Updates

1.  $J$  is always a lower bound on  $J^*$  since  $J^*$  is achievable.
2.  $J^*$  is achieved through  $\pi^*$  specifically.
3.  $J$  is no worse than  $J^*$ .
4.  $J$  is no worse than  $J^*$ .
5. We can speedup computation of  $J$  by initializing with  $J^*$ .

# Policy Iteration with Warm Start Values

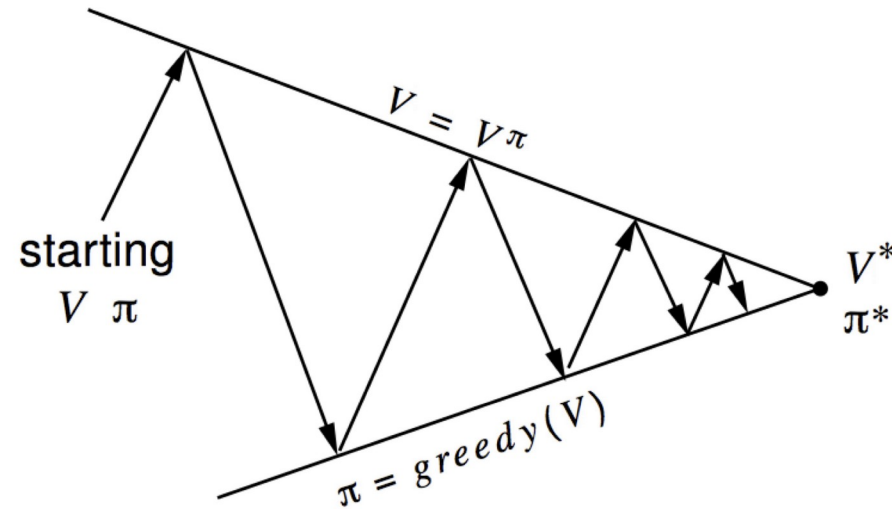
1. Initialize  $\pi$  to any policy.
2. Set  $V$  using value iteration.
3. For  $i = 0, 1, 2, \dots$ 
  - a) Set  $\pi_i$
  - b) Compute  $V_i$  starting from  $V$ .

# Interleaved Policy and Value Updates

1. Initialize  $\pi$  to any policy.
2. Set  $V$  using value iteration.
3. For  $i = 0, 1, 2, \dots$ 
  - a) Set  $\pi$
  - b) Set  $V$

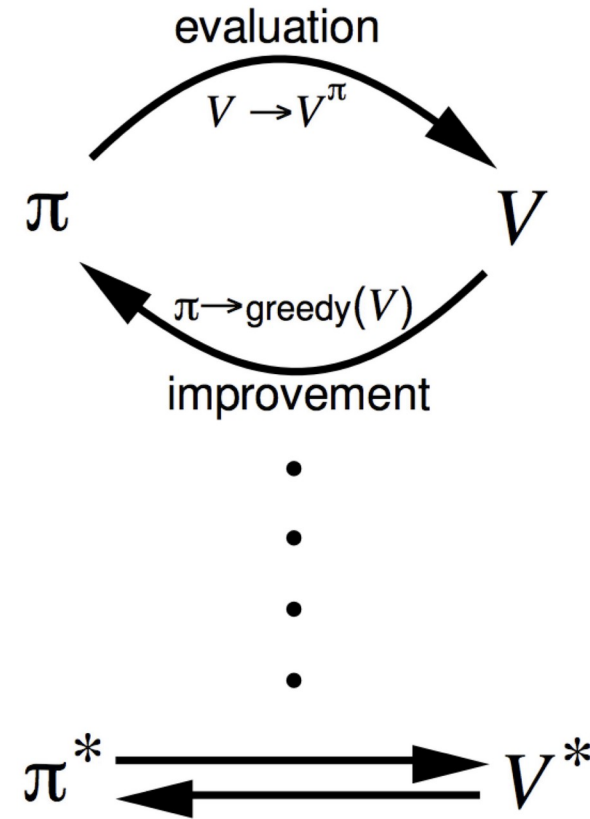


# Convergence of Policy Iteration



**Policy evaluation** Estimate  $v_\pi$   
Iterative policy evaluation

**Policy improvement** Generate  $\pi' \geq \pi$   
Greedy policy improvement



# Any Questions?

???