

Markov Decision Process II

CSE 4617: Artificial Intelligence



Isham Tashdeed
Lecturer, CSE



$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k(s')]$$

Value Iteration



V_2

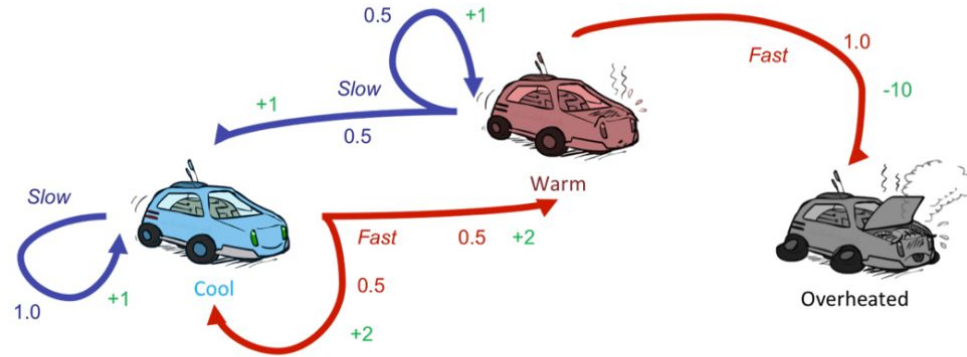
| | | |
|-----|-----|---|
| 3.5 | 2.5 | 0 |
|-----|-----|---|

V_1

| | | |
|---|---|---|
| 2 | 1 | 0 |
|---|---|---|

V_0

| | | |
|---|---|---|
| 0 | 0 | 0 |
|---|---|---|

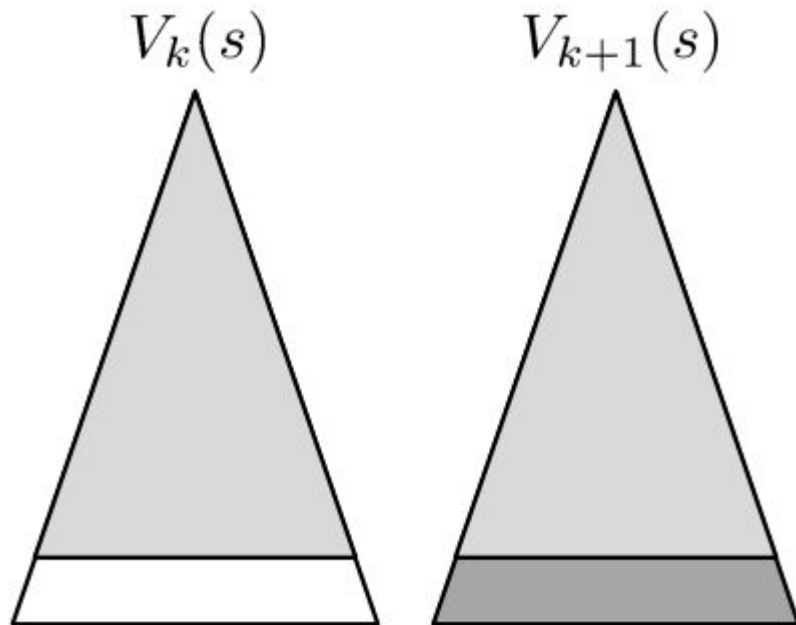


Assume no discount!

Will It Converge?

How do we know that the V_k vectors are going to converge?

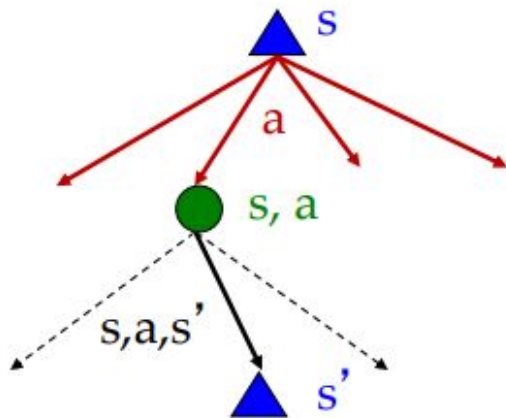
- Case 1 \rightarrow Time limited MDPs (Always terminates after m steps)
- Case 2 \rightarrow Value of $0 < \gamma < 1$
 - The last layer is at best R_{max} or at worst R_{min}
 - The difference between the two trees is at most $\gamma^k \max|R|$



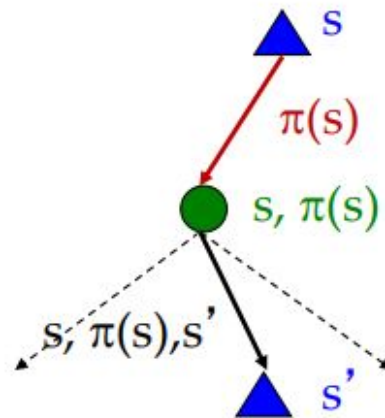
Policy Evaluation

Fixed Policies

Do the optimal action

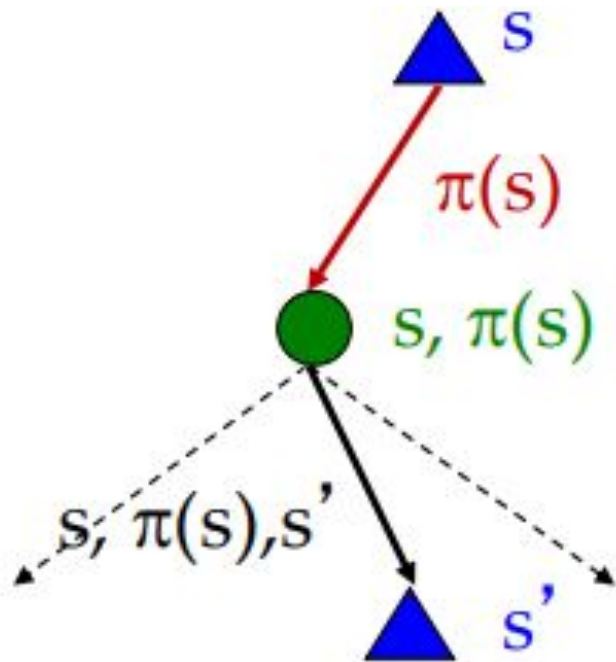


Do what π says to do



- If we have some fixed policy $\pi(s) \rightarrow$ Only one action per state, No need to max over all actions
- For regular value iteration $\rightarrow O(S^2A)$ per iteration

Utilities of a Fixed Policy



- Compute the utility of a state s under a fixed policy $\pi \rightarrow$ usually a non-optimal policy
- $V^\pi(s) \rightarrow$ Total expected utility when starting from state s and following π
- $$V^\pi(s) = \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V^\pi(s')]$$
- How to calculate $V^\pi(s)$ or $V^\pi(s')$?
- $V_0^\pi(s) = 0$
- $$V_{k+1}^\pi(s) \leftarrow \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V_k^\pi(s')]$$
- Efficiency $\rightarrow O(S^2)$ per iteration
- Or try solving the linear equation!

Policy Evaluation

Always Go Right



Always Go Forward



Policy Extraction

Policy Extraction

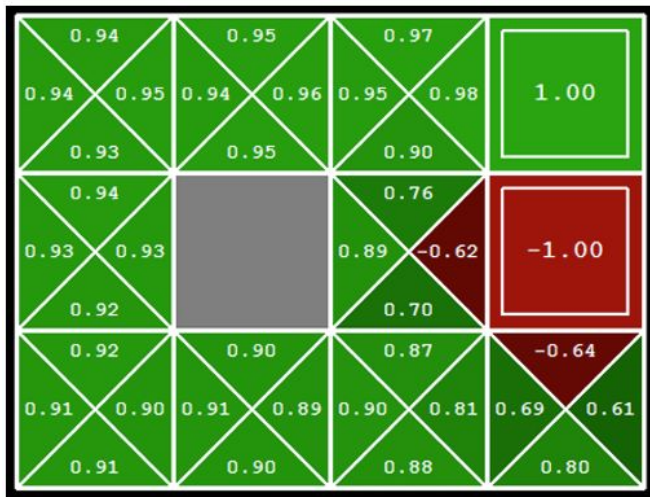
| | | | |
|------|------|------|-------|
| 0.95 | 0.96 | 0.98 | 1.00 |
| 0.94 | | 0.89 | -1.00 |
| 0.92 | 0.91 | 0.90 | 0.80 |

- Given the optimal values $V^*(s)$, how to extract the policy?
- We need to run a one step expectimax to find the policy

$$\pi^*(s) = \operatorname{argmax}_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

- This is policy extraction

Policy Extraction



- Given the optimal q-values $Q^*(s,a)$, how to extract the policy?
- We need to run a one step expectimax to find the policy

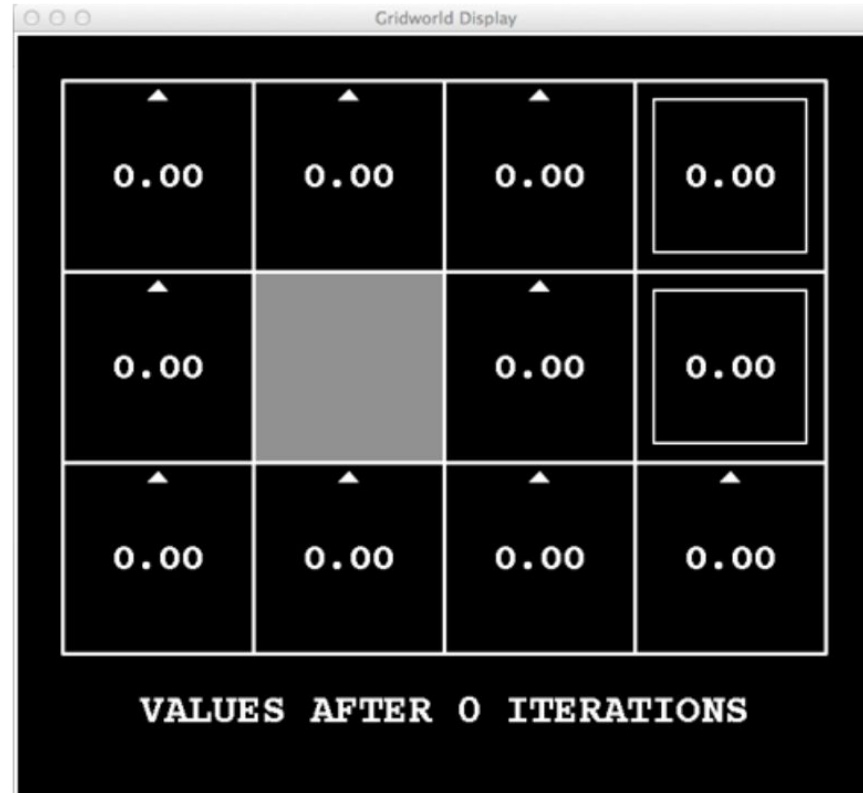
$$\pi^*(s) = \operatorname{argmax}_a Q^*(s, a)$$

- Actions are much easier to select from q-values, than values
- This is the very basic of reinforcement learning

Policy Iteration

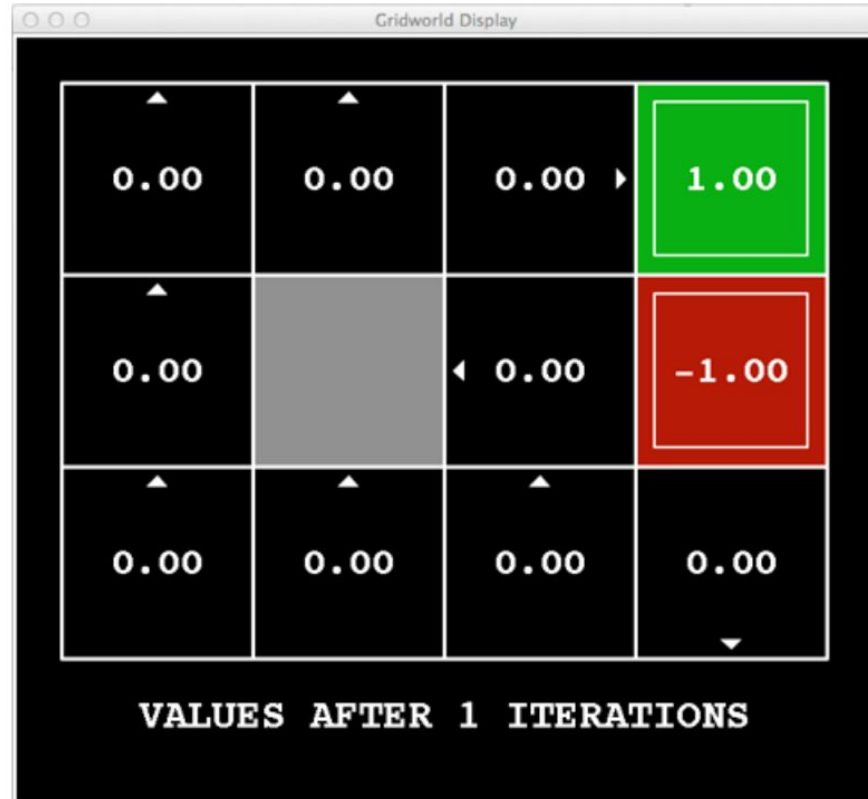
Gridworld Value Iteration

$k = 0$



Gridworld Value Iteration

$k = 1$



Gridworld Value Iteration

$k = 2$



Gridworld Value Iteration

$k = 3$



Gridworld Value Iteration

$k = 4$



Gridworld Value Iteration

$k = 5$



Gridworld Value Iteration

$k = 6$



Gridworld Value Iteration

$k = 7$



Gridworld Value Iteration

$k = 8$



Gridworld Value Iteration

$k = 9$

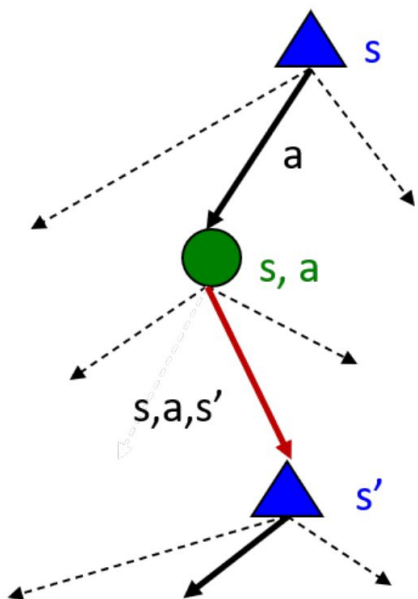


Gridworld Value Iteration

$k = 10$



Problems with Value Iteration



s is a
state

(s, a) is a
 q -state

(s, a, s') is a
transition

- Repeats the Bellman updates at every state each iteration
- It is slow $\rightarrow O(S^2A)$
- The direction at each state rarely changes
- Policy converges long before the values converge

Policy Iteration

- Step 1: Policy evaluation \rightarrow Calculate utilities for some fixed policy until convergence (Not optimal policy)
- Step 2: Policy improvement \rightarrow Update policy using one-step look-ahead with resulting converged (but not optimal) utilities as future values
- Repeat until policy converges



Policy Iteration



- Evaluation: for fixed current policy π , find the values with policy evaluation:
 - Iterate until convergence

$$V_{k+1}^{\pi_i}(s) \leftarrow \sum_{s'} T(s, \pi_i(s), s') [R(s, \pi_i(s), s') + \gamma V_k^{\pi_i}(s')]$$

Policy Iteration



- Evaluation: For fixed current policy π , find the values with policy evaluation:
 - Iterate until convergence

$$V_{k+1}^{\pi_i}(s) \leftarrow \sum_{s'} T(s, \pi_i(s), s') [R(s, \pi_i(s), s') + \gamma V_k^{\pi_i}(s')]$$

- Improvement: For fixed values, get a better policy using policy extraction:

$$\pi_{i+1}(s) = \operatorname{argmax}_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^{\pi_i}(s')]$$

- Repeat until convergence

Policy Iteration

Both value iteration and policy iteration compute the same thing

In value iteration:

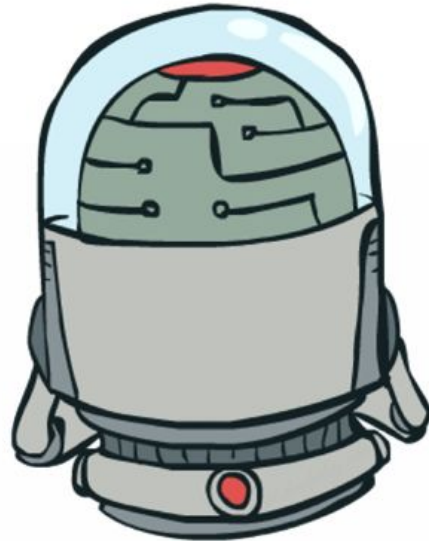
- Both utility values and policy is updated at each iteration
- Slow

In policy iteration:

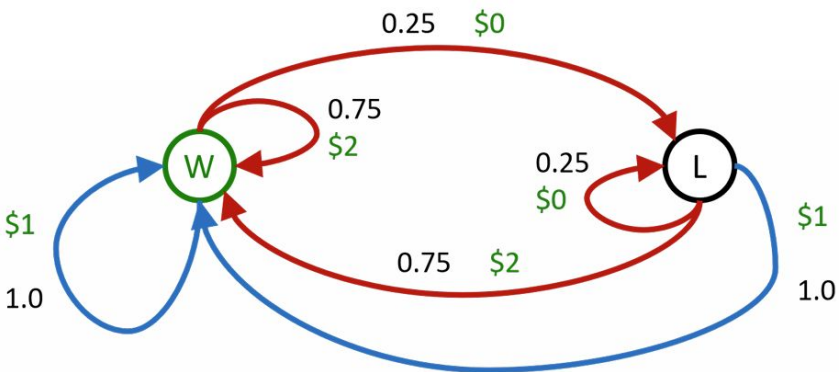
- First we update the utility values for a fixed policy \rightarrow faster as we only consider one action
- After that, new policy is chosen \rightarrow Slow, similar to value iteration

Reinforcement Learning

Reinforcement Learning



Reinforcement Learning

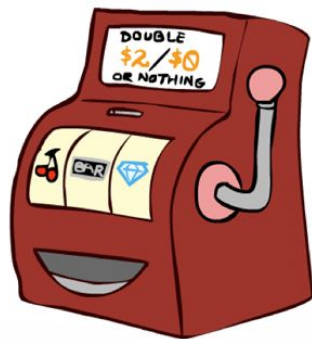
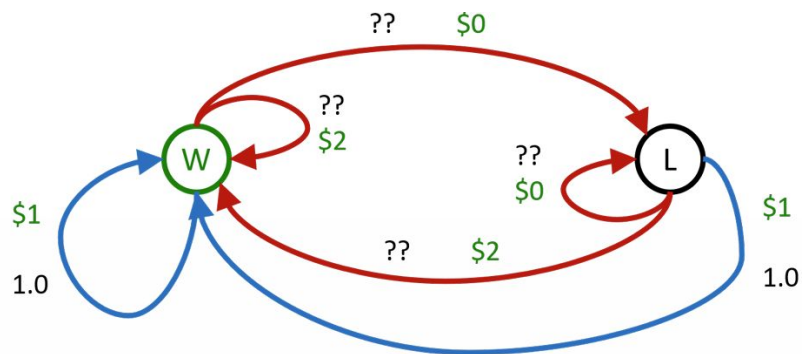


- States \rightarrow Win, Lose
- Actions \rightarrow Blue, Red
- $\gamma = 1$
- 100 time steps

Solving MDPs is offline planning

- You know the details of the MDP
- You find all the quantities through computation
- You don't actually play the game!

Reinforcement Learning



| | | | | |
|-----|-----|-----|-----|-----|
| \$0 | \$0 | \$0 | \$2 | \$0 |
| \$2 | \$0 | \$0 | \$0 | \$0 |

- States \rightarrow Win, Lose
- Actions \rightarrow Blue, Red
- $\gamma = 1$
- 100 time steps

Reinforcement Learning is online planning

- You **don't know** the details of the MDP
- You actually play the game!

Reinforcement Learning

In reinforcement learning

- Exploration: Try unknown actions to get information.
- Exploitation: Eventually, you have to use what you know
- Regret: Even if you learn intelligently, you make mistakes
- Sampling: You have to try things repeatedly
- Difficulty: Learning can be much harder than solving a known MDP



Thank you