# CS 473: Artificial Intelligence

## MDP Planning: Value Iteration and Policy Iteration



Travis Mandel *(subbing for Dan Weld)*
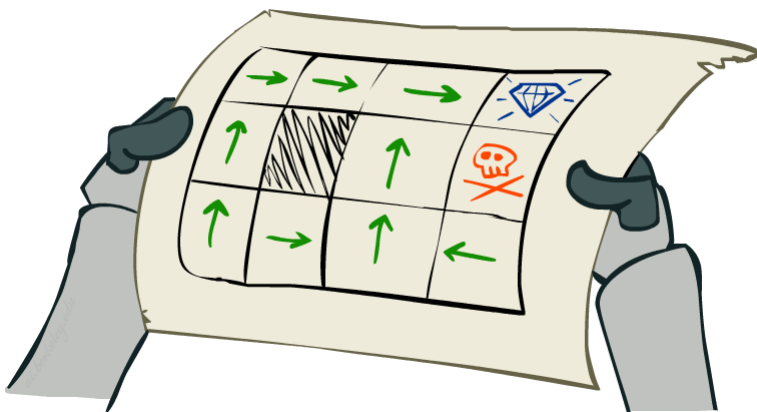
University of Washington

---

# Reminder: Midterm Monday!!

- Will cover everything from Search to Value Iteration
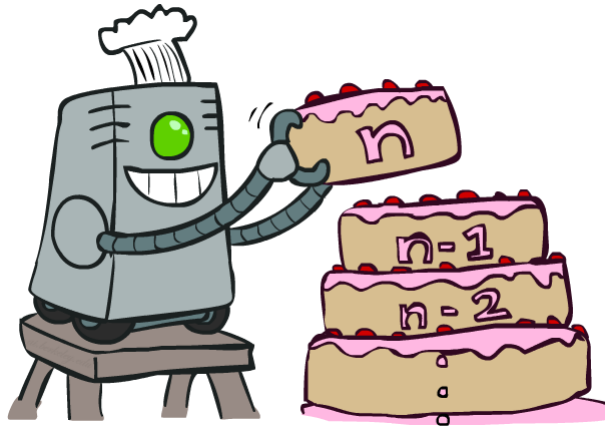  - One page (double-sided, 8.5 x 11) notes allowed

# Reminder: MDP Planning

- Given an MDP, find optimal policy $\pi^*$: S$\rightarrow$A that maximizes expected discounted reward
  - Sometimes called "Solving" the MDP
- Being so long-term complicates things
  - Simplifies things if we know long-term value of state

# MDP Planning



- Value Iteration
  - Prioritized Sweeping

- Policy Iteration

# Value Iteration



# Value Iteration

- **Forall s, initialize $V_0(s) = 0$** *no time steps left means an expected reward of zero*
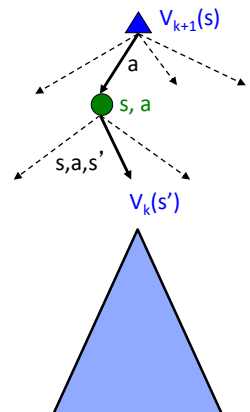
- **Repeat**
  K += 1
  $Q_{k+1}(s, a) = \Sigma_{s'} T(s, a, s') [ R(s, a, s') + \gamma V_k(s')]$

  $V_{k+1}(s) = \text{Max}_a Q_{k+1}(s, a)$

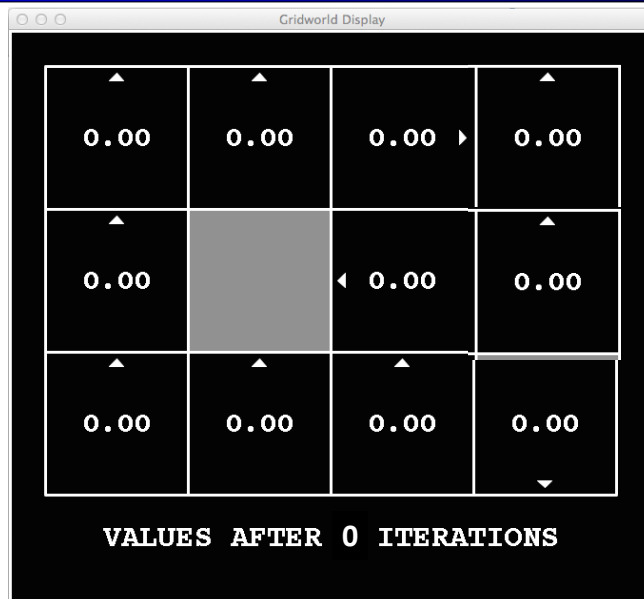  } do $\forall s, a$

- **Repeat until $|V_{k+1}(s) - V_k(s)| < \epsilon$, forall s** *"convergence"*

Successive approximation; dynamic programming

$V_{k+1}(s)$
a
s, a
s,a,s'
$V_k(s')$

3

# k=0



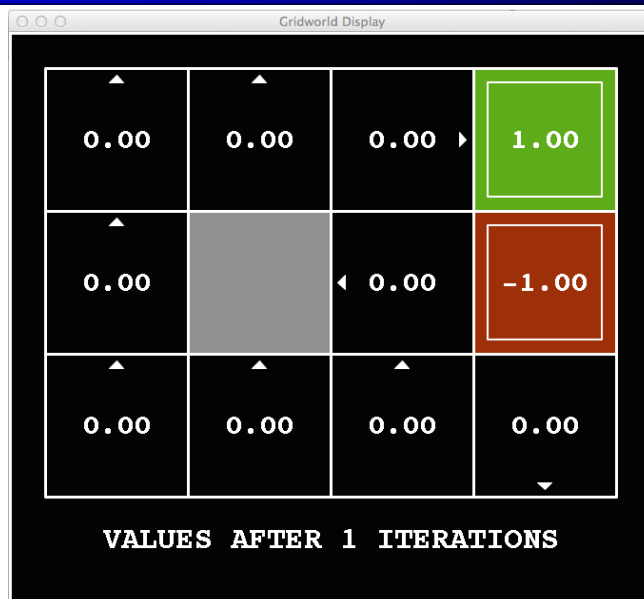VALUES AFTER 0 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=1

If agent is in 4,3, it only has one legal action: get jewel. It gets a reward and the game is over.
If agent is in the pit, it has only one legal action, die. It gets a penalty and the game is over.

Agent does NOT get a reward for moving INTO 4,3.



VALUES AFTER 1 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=2



VALUES AFTER 2 ITERATIONS

0.8 (0 + 0.9*1)
+ 0.1  (0 + 0.9*0)
+ 0.1  (0 + 0.9*0)

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=3



VALUES AFTER 3 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=4



VALUES AFTER 4 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=5



VALUES AFTER 5 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=6



VALUES AFTER 6 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=7



VALUES AFTER 7 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=8



|  |  |  |  |
|---|---|---|---|
| 0.63 ▶ | 0.74 ▶ | 0.85 ▶ | 1.00 |
| ▲ 0.53 |  | ▲ 0.57 | −1.00 |
| ▲ 0.42 | 0.39 ▶ | ▲ 0.46 | ◀ 0.26 |

VALUES AFTER 8 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=9



|  |  |  |  |
|---|---|---|---|
| 0.64 ▶ | 0.74 ▶ | 0.85 ▶ | 1.00 |
| ▲ 0.55 |  | ▲ 0.57 | −1.00 |
| ▲ 0.46 | 0.40 ▶ | ▲ 0.47 | ◀ 0.27 |

VALUES AFTER 9 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=10



VALUES AFTER 10 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=11



VALUES AFTER 11 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=12



VALUES AFTER 12 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=100



VALUES AFTER 100 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# VI: Policy Extraction



# Computing Actions from Values

- Let's imagine we have the optimal values V*(s)

- How should we act?
  - In general, it's not obvious!

- We need to do a mini-expectimax (one step)

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

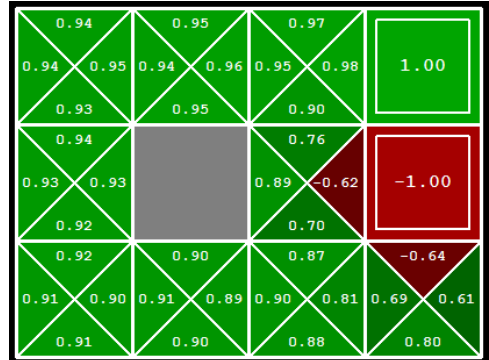- This is called policy extraction, since it gets the policy implied by the values

# Computing Actions from Q-Values

- Let's imagine we have the optimal q-values:
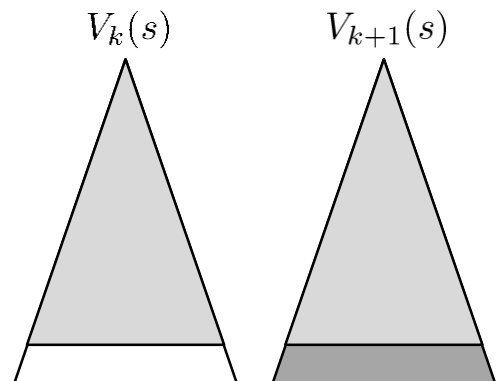
- How should we act?
  - Completely trivial to decide!

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



- Important lesson: actions are easier to select from q-values than values!

# Convergence*

- How do we know the $V_k$ vectors will converge?

- Case 1: If the tree has maximum depth M, then $V_M$ holds the actual untruncated values

- Case 2: If the discount is less than 1
  - Sketch: For any state $V_k$ and $V_{k+1}$ can be viewed as depth k+1 expectimax results in nearly identical search trees
  - The max difference happens if big reward at k+1 level
  - That last layer is at best all $R_{MAX}$
  - But everything is discounted by $\gamma^k$ that far out
  - So $V_k$ and $V_{k+1}$ are at most $\gamma^k$ max|R| different
  - So as k increases, the values converge

$$V_k(s) \qquad V_{k+1}(s)$$

# Value Iteration - Recap

- **Forall s, Initialize $V_0(s) = 0$**    *no time steps left means an expected reward of zero*

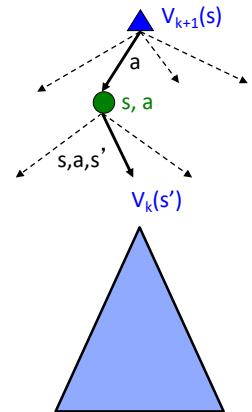- **Repeat**            *do Bellman backups*
    K += 1
    Repeat for all states, s, and all actions, a:

    $$Q_{k+1}(s, a) = \Sigma_{s'}\ T(s, a, s')\ [\ R(s, a, s') + \gamma\ V_k(s')\ ]$$

    $$V_{k+1}(s) = \text{Max}_a\ Q_{k+1}(s, a)$$

    $V_{k+1}(s)$
    a
    s, a
    s,a,s'
    $V_k(s')$

- **Until $|V_{k+1}(s) - V_k(s)| < \varepsilon$,**     **forall s**   *"convergence"*

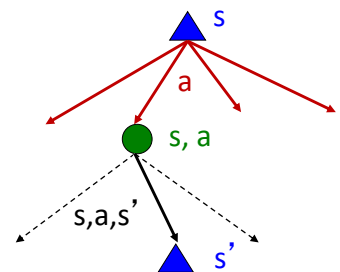- **Theorem: will converge to unique optimal values**

---

# Problems with Value Iteration

- Value iteration repeats the Bellman updates:

    $$Q_{k+1}(s, a) = \Sigma_{s'}\ T(s, a, s')\ [\ R(s, a, s') + \gamma\ V_k(s')\ ]$$

    $$V_{k+1}(s) = \text{Max}_a\ Q_{k+1}(s, a)$$

    s
    a
    s, a
    s,a,s'
    s'

- Problem 1: It's slow – $O(S^2A)$ per iteration

- Problem 2: The "max" at each state rarely changes

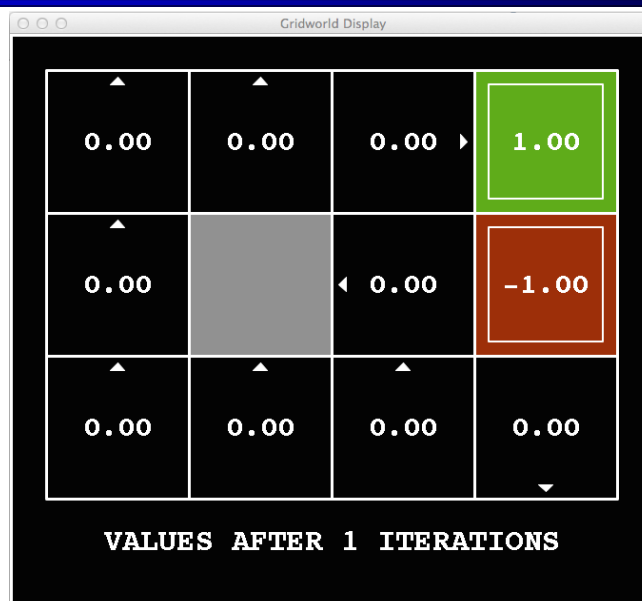- Problem 3: The policy often converges long before the values

[Demo: value iteration (L9D2)]

# VI → Asynchronous VI

- Is it essential to back up *all* states in each iteration?
  - No!

- States may be backed up
  - many times or not at all
  - in any order

- As long as no state gets starved…
  - convergence properties still hold!!

# k=1



Noise = 0.2
Discount = 0.9
Living reward = 0

# k=2



VALUES AFTER 2 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=3



VALUES AFTER 3 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=8



VALUES AFTER 8 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=9



VALUES AFTER 9 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=10



VALUES AFTER 10 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=11



VALUES AFTER 11 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=12



VALUES AFTER 12 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=100



VALUES AFTER 100 ITERATIONS

Noise = 0.2
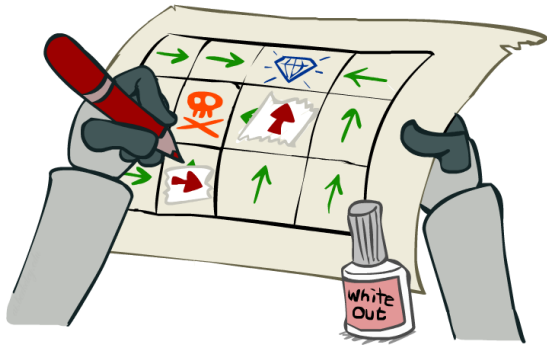Discount = 0.9
Living reward = 0

# Asynch VI: Prioritized Sweeping

- Why backup a state if values of successors *unchanged*?
- Prefer backing a state
  - whose successors had *most* change
- Priority Queue of (state, expected change in value)
- Backup in the order of priority
- After backing up state s', update priority queue
  - for all predecessors s (ie all states where an action can reach s')
  - Priority(s) ← $T(s,a,s') * |V^{k+1}(s') - V^k(s')|$
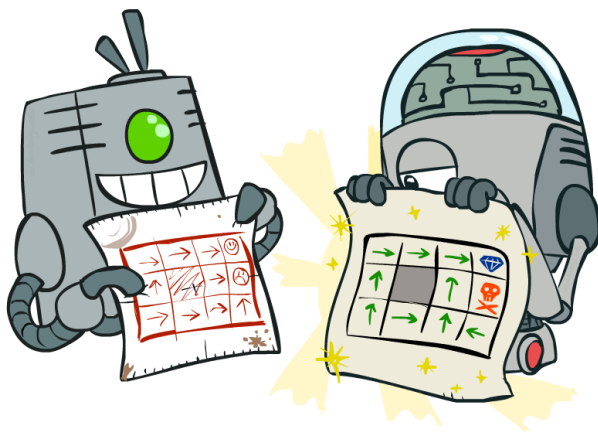
# Prioritized Sweeping

- Pros?

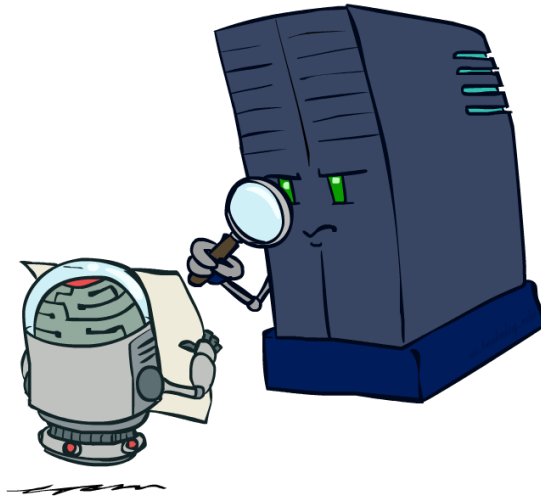- Cons?

# MDP Planning

- Value Iteration
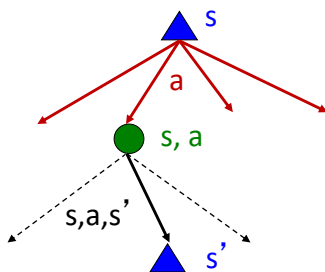  - Prioritized Sweeping

- Policy Iteration

# Policy Methods

Policy Iteration =
1. Policy Evaluation
2. Policy Improvement
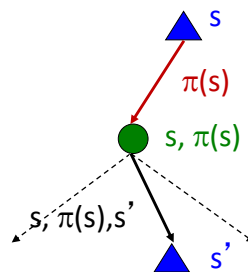
# Part 1 - Policy Evaluation



# Fixed Policies

Do the optimal action

Do what $\pi$ says to do



$s$

$a$

$s, a$

$s,a,s'$

$s'$

$s$

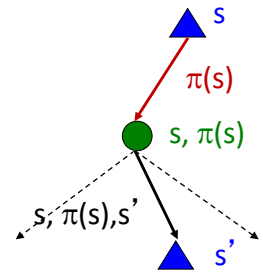$\pi(s)$

$s, \pi(s)$

$s, \pi(s),s'$

$s'$

- Expectimax trees max over all actions to compute the optimal values

- If we fixed some policy $\pi(s)$, then the tree would be simpler – only one action per state
    - … though the tree's value would depend on which policy we fixed
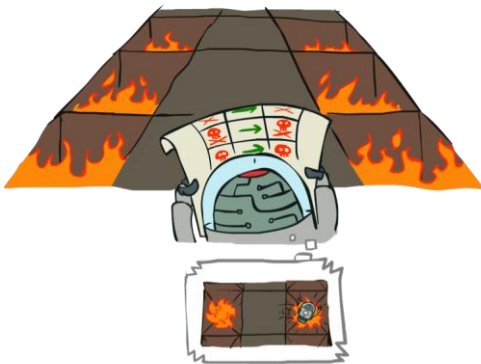
# Computing Utilities for a Fixed Policy

- **A new basic operation:** compute the utility of a state s under a fixed (generally non-optimal) policy

- Define the utility of a state s, under a fixed policy $\pi$:

  $V^\pi(s)$ = expected total discounted rewards starting in s and following $\pi$

- Recursive relation (variation of Bellman equation):

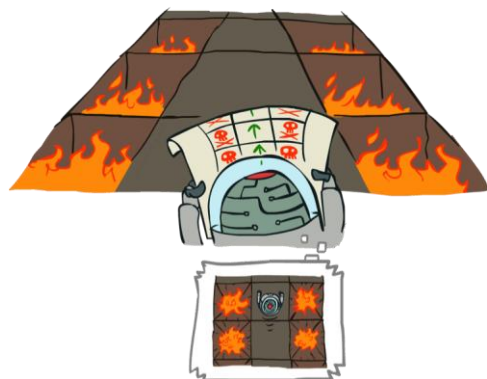$$V^\pi(s) = \sum_{s'} T(s, \pi(s), s')[R(s, \pi(s), s') + \gamma V^\pi(s')]$$

s

$\pi$(s)

s, $\pi$(s)

s, $\pi$(s),s'

s'

---

# Example: Policy Evaluation

Always Go Right                    Always Go Forward

# Example: Policy Evaluation

Always Go Right

| | | |
|---|---|---|
| -10.00 | 100.00 | -10.00 |
| -10.00 | 1.09 ▸ | -10.00 |
| -10.00 | -7.88 ▸ | -10.00 |
| -10.00 | -8.69 ▸ | -10.00 |

Always Go Forward

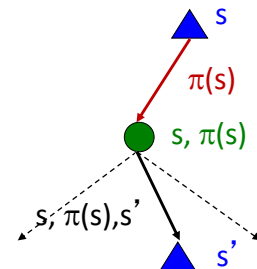| | | |
|---|---|---|
| -10.00 | 100.00 | -10.00 |
| -10.00 | 70.20 | -10.00 |
| -10.00 | 48.74 | -10.00 |
| -10.00 | 33.30 | -10.00 |

---

# Iterative Policy Evaluation Algorithm

- How do we calculate the V's for a fixed policy $\pi$?

- Idea 1: Turn recursive Bellman equations into updates
  (like value iteration)

$$V_0^\pi(s) = 0$$

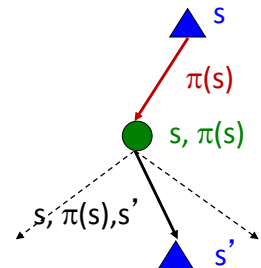$$V_{k+1}^\pi(s) \leftarrow \sum_{s'} T(s, \pi(s), s')[R(s, \pi(s), s') + \gamma V_k^\pi(s')]$$

s

$\pi$(s)

s, $\pi$(s)

s, $\pi$(s),s'

s'

- Efficiency: $O(S^2)$ per iteration
  - Often converges in much smaller number of iterations compared to VI

# Linear Policy Evaluation Algorithm

- How do we calculate the V's for a fixed policy $\pi$?

- Idea 2: Without the maxes, the Bellman equations are just a linear system of equations

$$V^\pi(s) = \sum_{s'} T\big(s, \pi(s), s'\big)[R\big(s, \pi(s), s'\big) + \gamma V^\pi(s')]$$

- Solve with Matlab (or your favorite linear system solver)
  - S equations, S unknowns = $O(S^3)$ and EXACT!
  - In large spaces, still too expensive

s

$\pi(s)$

s, $\pi(s)$

s, $\pi(s)$,s'

s'

---

# Part 2 - Policy Iteration

# Policy Iteration

- Initialize π(s) to random actions
- Repeat
  - Step 1: Policy evaluation: calculate utilities of π at each s using a nested loop
  - Step 2: Policy improvement: update policy using one-step look-ahead
    "For each s, what's the best action I could execute, assuming I then follow π?
    Let π'(s) = this best action.
    π = π'
- Until policy doesn't change

# Policy Iteration Details

- Let i =0
- Initialize $\pi_i(s)$ to random actions
- Repeat
  - Step 1: Policy evaluation:
    - Initialize k=0;   Forall s, $V_0^\pi(s) = 0$
    - Repeat until $V^\pi$ converges
      - For each state s, $V_{k+1}^{\pi_i}(s) \leftarrow \sum_{s'} T(s, \pi_i(s), s') \left[ R(s, \pi_i(s), s') + \gamma V_k^{\pi_i}(s') \right]$
      - Let k += 1
  - Step 2: Policy improvement:
    - For each state, s, $\pi_{i+1}(s) = \arg\max_a \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma V^{\pi_i}(s') \right]$
    - If $\pi_i == \pi_{i+1}$ then it's optimal; return it.
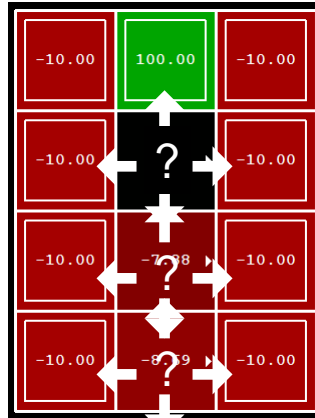    - Else let i += 1

# Example

Initialize $\pi_0$ to "always go right"

Perform policy evaluation

Perform policy improvement
   Iterate through states

Has policy changed?

Yes!  i += 1

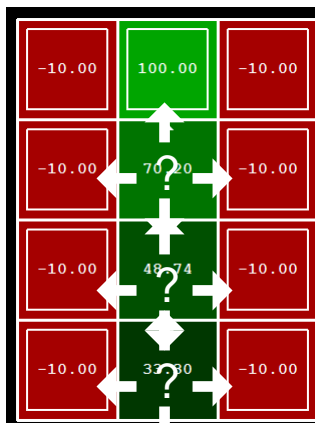| -10.00 | 100.00 | -10.00 |
|--------|--------|--------|
| -10.00 | ? | -10.00 |
| -10.00 | -7.?8 ? | -10.00 |
| -10.00 | -8.?9 ? | -10.00 |

---

# Example

$\pi_1$ says "always go up"

Perform policy evaluation

Perform policy improvement
   Iterate through states

Has policy changed?

No!  We have the optimal policy

| -10.00 | 100.00 | -10.00 |
|--------|--------|--------|
| -10.00 | 70.?0 ? | -10.00 |
| -10.00 | 48.?4 ? | -10.00 |
| -10.00 | 33.?0 ? | -10.00 |

# Example: Policy Evaluation

Always Go Right

| | | |
|---|---|---|
| -10.00 | 100.00 | -10.00 |
| -10.00 | 1.09 ▸ | -10.00 |
| -10.00 | -7.88 ▸ | -10.00 |
| -10.00 | -8.69 ▸ | -10.00 |

Always Go Forward

| | | |
|---|---|---|
| -10.00 | 100.00 | -10.00 |
| -10.00 | 70.20 | -10.00 |
| -10.00 | 48.74 | -10.00 |
| -10.00 | 33.30 | -10.00 |

# Policy Iteration Properties

- Policy iteration finds the optimal policy, guaranteed (assuming exact evaluation)!
- Often converges (much) faster

# Comparison

- Both value iteration and policy iteration compute the same thing (all optimal values)

- In value iteration:
  - Every iteration updates both the values and (implicitly) the policy
  - We don't track the policy, but taking the max over actions implicitly recomputes it
  - What is the space being searched?

- In policy iteration:
  - We do fewer iterations
  - Each one is slower (must update all $V^\pi$ and then choose new best $\pi$)
  - What is the space being searched?

- Both are dynamic programs for planning in MDPs