

# Announcements

## Assignments:

- P3: Optimization; Due today, 10 pm
- HW7 (online) 10/22 Tue, 10 pm

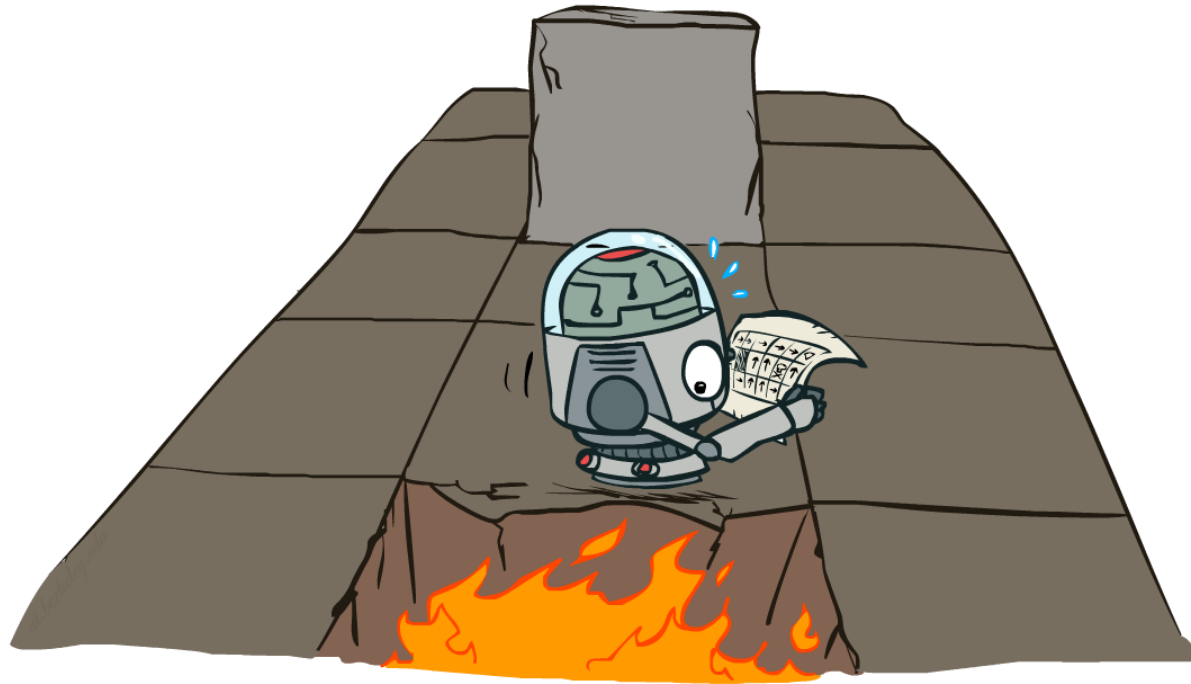
Recitations canceled on: October 18 (Mid-Semester Break) and October 25 (Day for Community Engagement)

We will provide recitation worksheet (reference for midterm/final)

Piazza post for In-class Questions

# AI: Representation and Problem Solving

## Markov Decision Processes II



Instructors: Fei Fang & Pat Virtue

Slide credits: CMU AI and <http://ai.berkeley.edu>

# Learning Objectives

- Write **Bellman Equation** for state-value and Q-value for optimal policy and a given policy
- Describe and implement **value iteration algorithm** (through **Bellman update**) for **solving MDPs**
- Describe and implement **policy iteration algorithm** (through **policy evaluation** and **policy improvement**) for **solving MDPs**
- Understand **convergence** for value iteration and policy iteration
- Understand concept of **exploration, exploitation, regret**

# MDP Notation

Standard expectimax:  $V(s) = \max_a \sum_{s'} P(s'|s, a) V(s')$

Bellman equations:  $V(s) = \max_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V(s')]$

Value iteration:  $V_{k+1}(s) = \max_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V_k(s')], \quad \forall s$

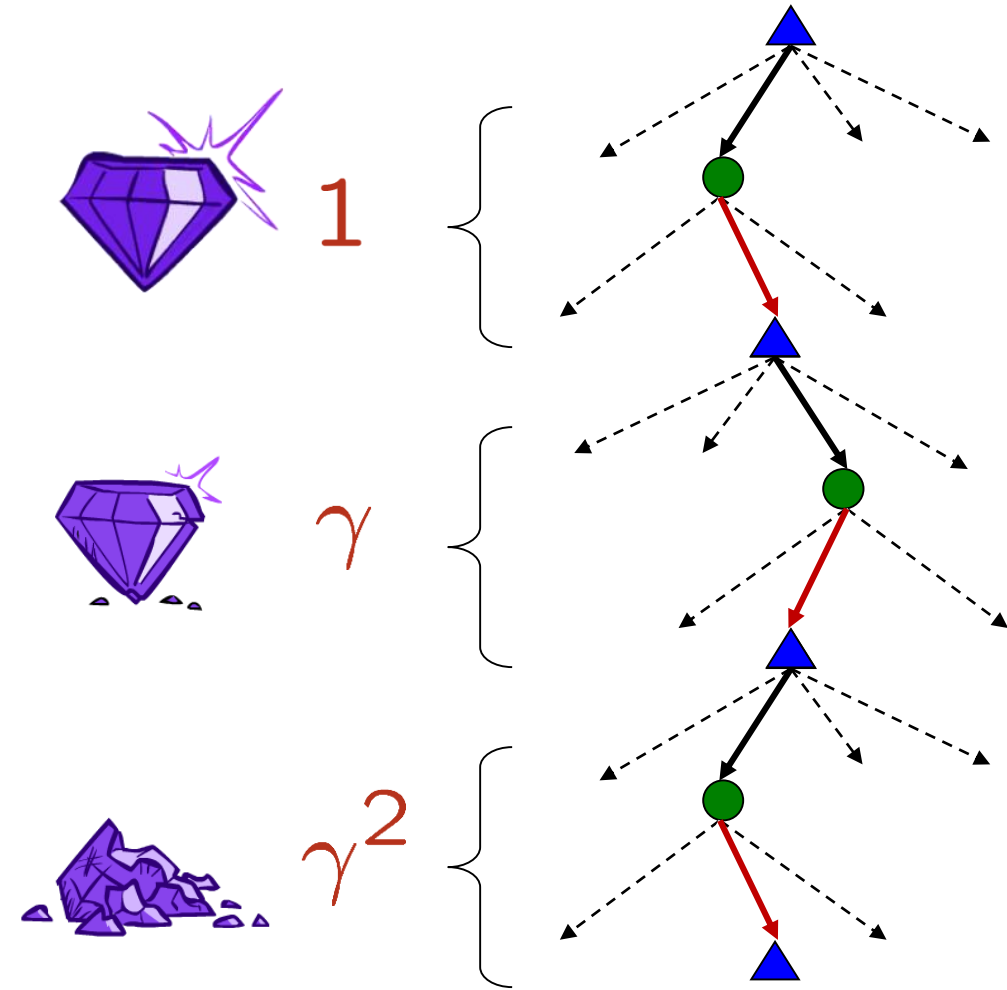
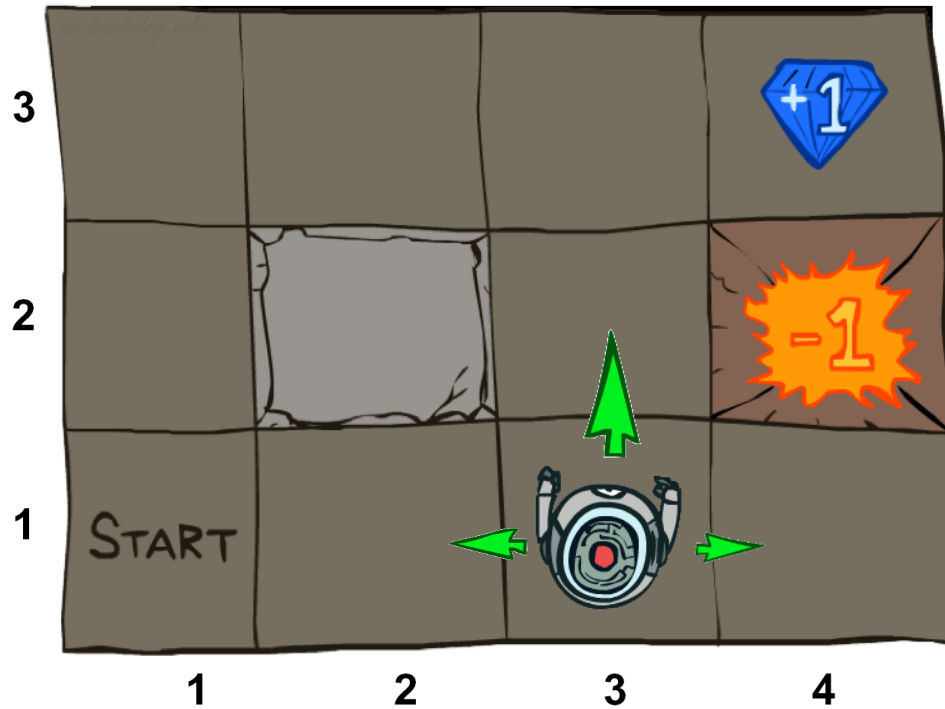
Q-iteration:  $Q_{k+1}(s, a) = \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma \max_{a'} Q_k(s', a')], \quad \forall s, a$

Policy extraction:  $\pi_V(s) = \operatorname{argmax}_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V(s')], \quad \forall s$

Policy evaluation:  $V_{k+1}^\pi(s) = \sum_{s'} P(s'|s, \pi(s)) [R(s, \pi(s), s') + \gamma V_k^\pi(s')], \quad \forall s$

Policy improvement:  $\pi_{new}(s) = \operatorname{argmax}_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V^{\pi_{old}}(s')], \quad \forall s$

# Example: Grid World



- Goal: maximize sum of (discounted) rewards

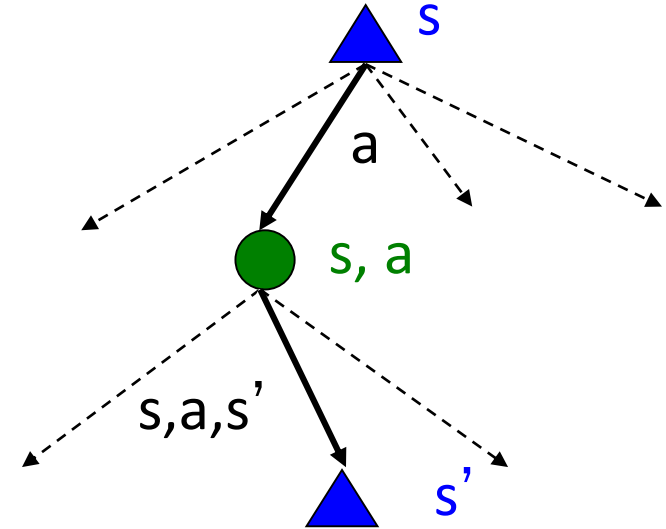
# MDP Quantities

## Markov decision processes:

- States  $S$
- Actions  $A$
- Transitions  $P(s' | s, a)$  (or  $T(s, a, s')$ )
- Rewards  $R(s, a, s')$  (and discount  $\gamma$ )
- Start state  $s_0$

## MDP quantities:

- Policy = map of states to actions
- Utility = sum of (discounted) rewards
- (State) Value = expected utility starting from a state (max node)
- Q-Value = expected utility starting from a state-action pair, i.e., q-state (chance node)



# MDP Optimal Quantities

- The optimal policy:

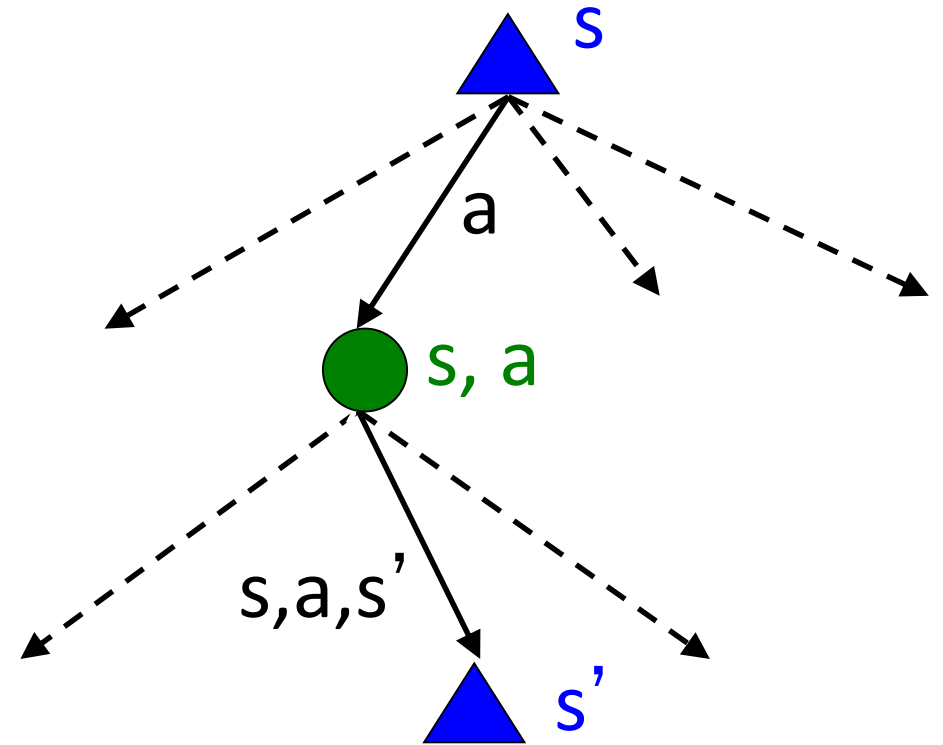
$\pi^*(s)$  = optimal action from state  $s$

- The (true) value (or utility) of a state  $s$ :

$V^*(s)$  = expected utility starting in  $s$  and acting optimally

- The (true) value (or utility) of a q-state  $(s,a)$ :

$Q^*(s,a)$  = expected utility starting out having taken action  $a$  from state  $s$  and (thereafter) acting optimally



$$V^*(s) < +\infty \text{ if } \gamma < 1 \text{ and } R(s, a, s') < \infty$$
$$U([r_0, \dots, r_\infty]) = \sum_{t=0}^{\infty} \gamma^t r_t \leq R_{\max} / (1 - \gamma)$$

Solve MDP: Find  $\pi^*$ ,  $V^*$  and/or  $Q^*$

[Demo: gridworld values (L9D1)]

# Piazza Poll 1

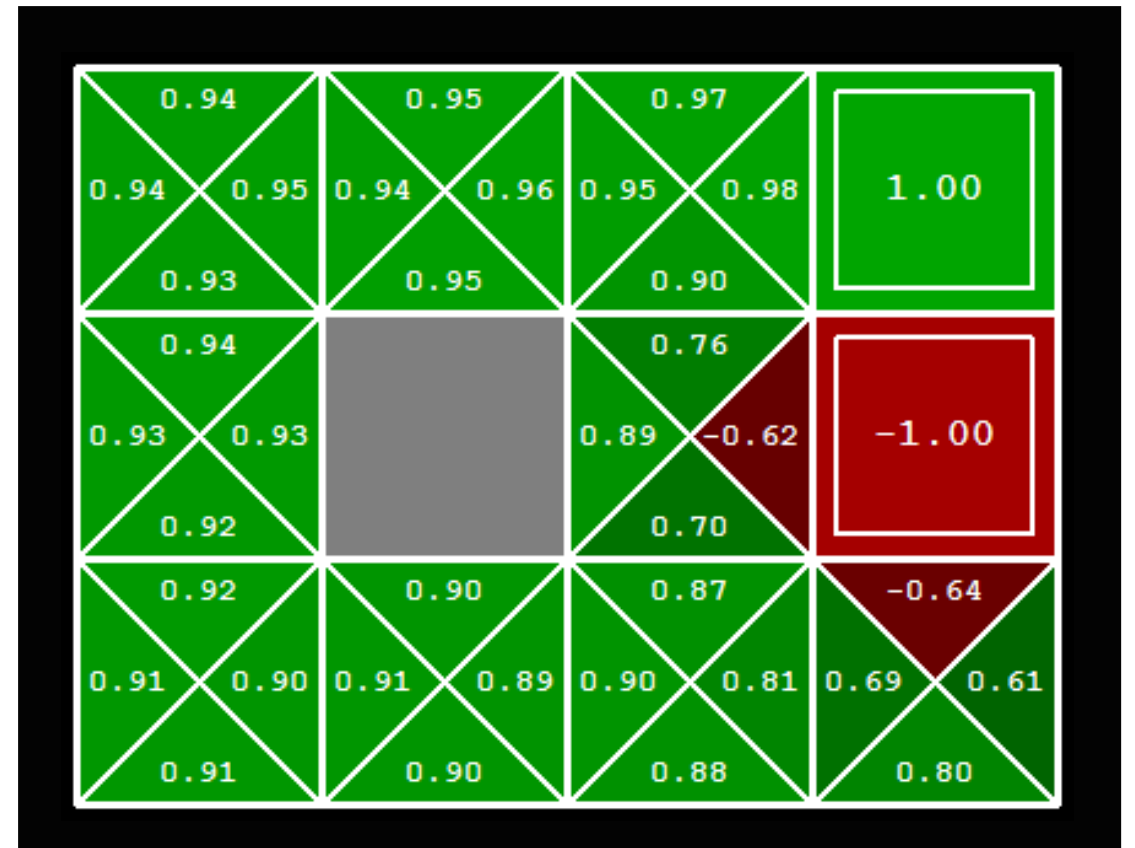
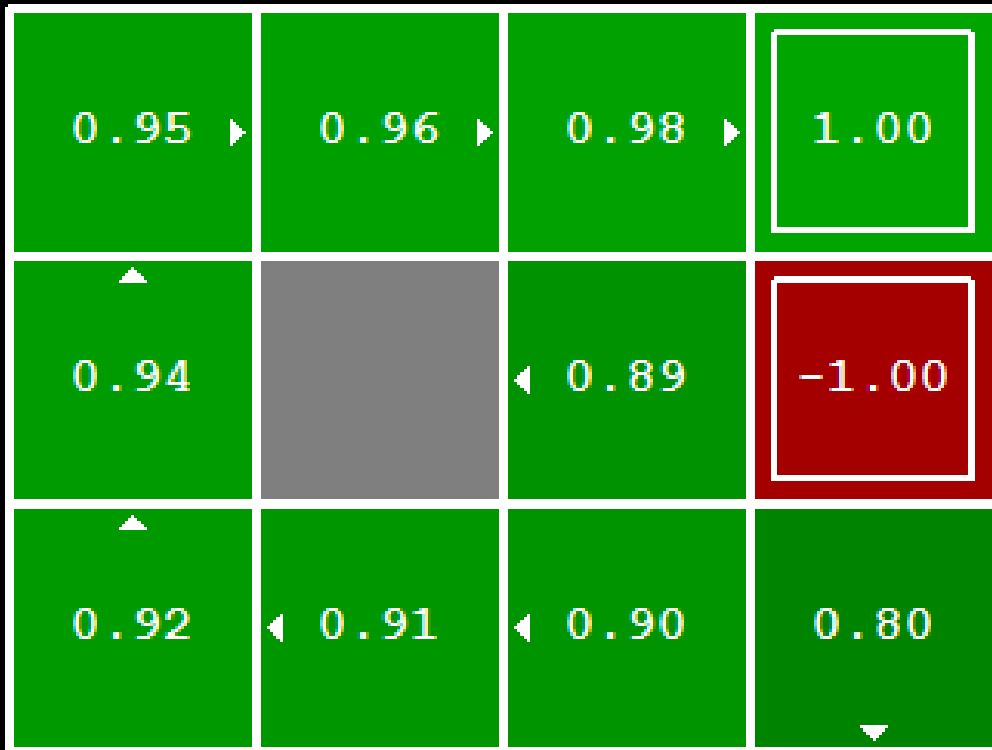
Which ones are true about optimal policy  $\pi^*(s)$ , true values  $V^*(s)$  and true Q-Values  $Q^*(s, a)$ ?

A:  $\pi^*(s) = \operatorname{argmax}_a V^*(s')$

where  $s' = \operatorname{argmax}_{s''} V^*(s, a, s'')$

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

C:  $V^*(s) = \max_a Q^*(s, a)$



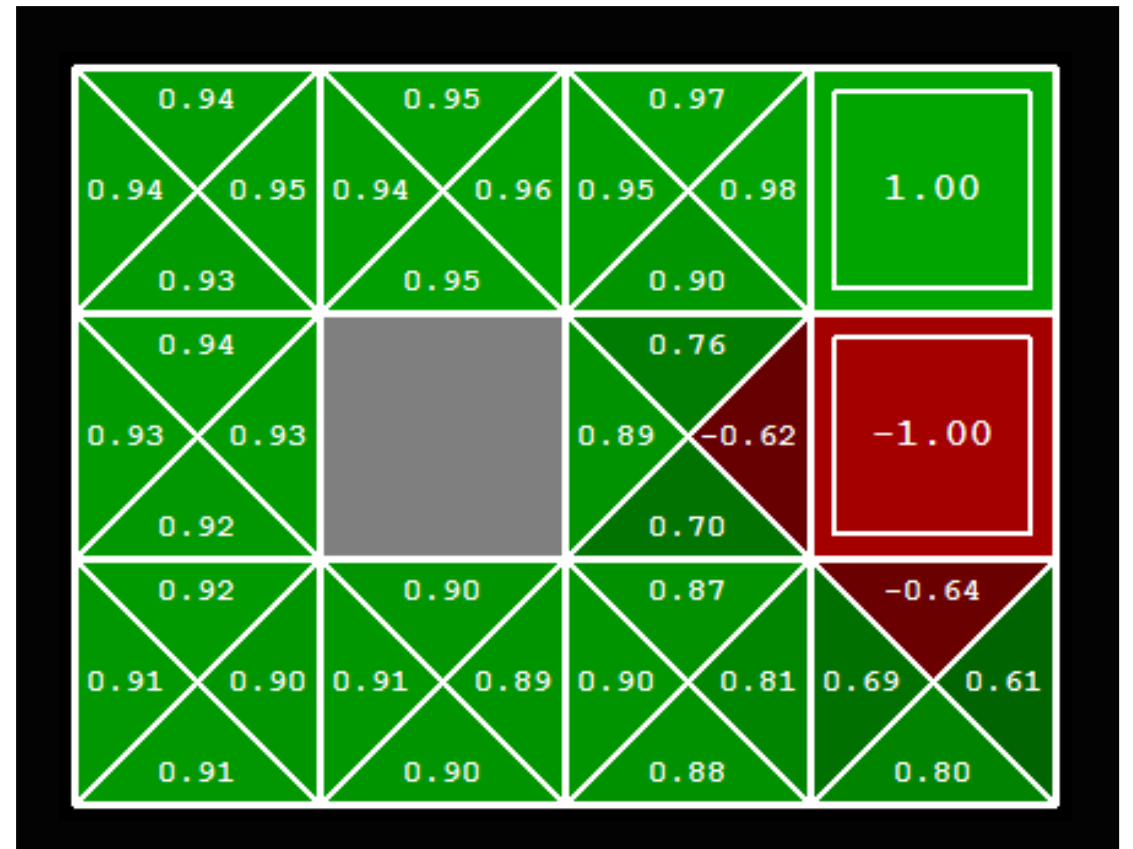
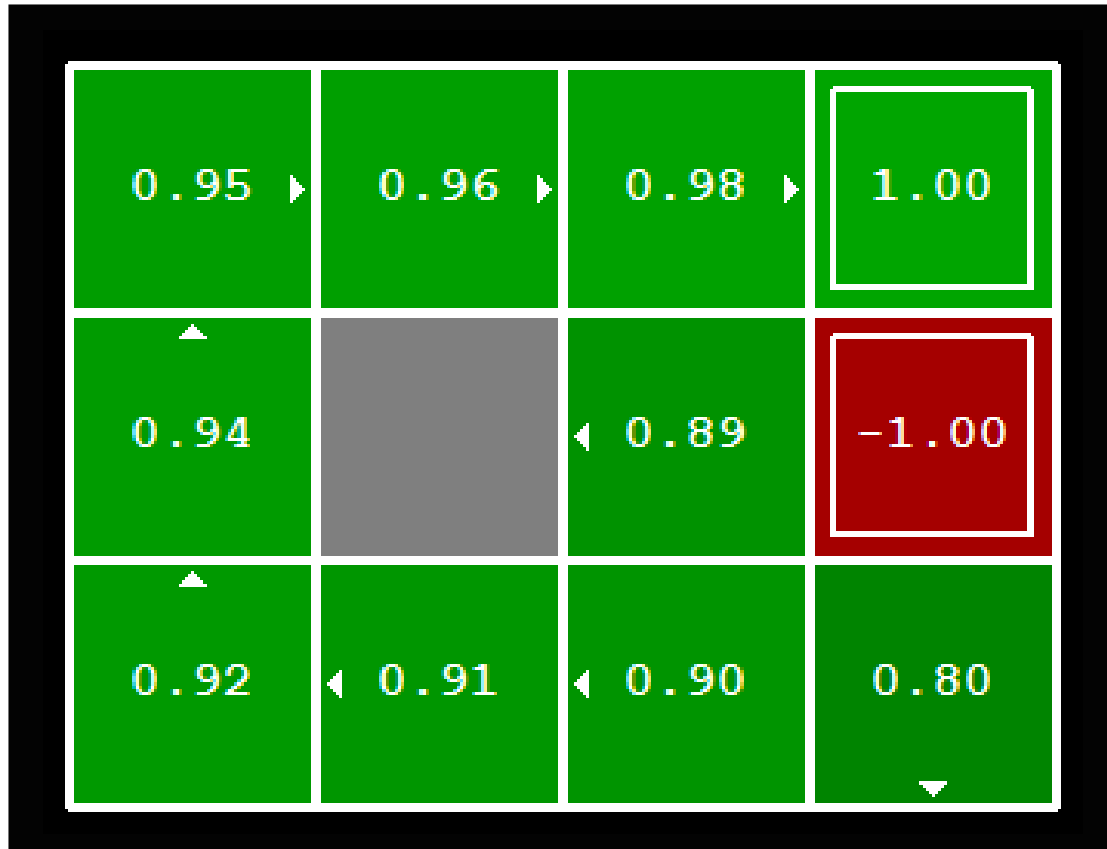


# Piazza Poll 1

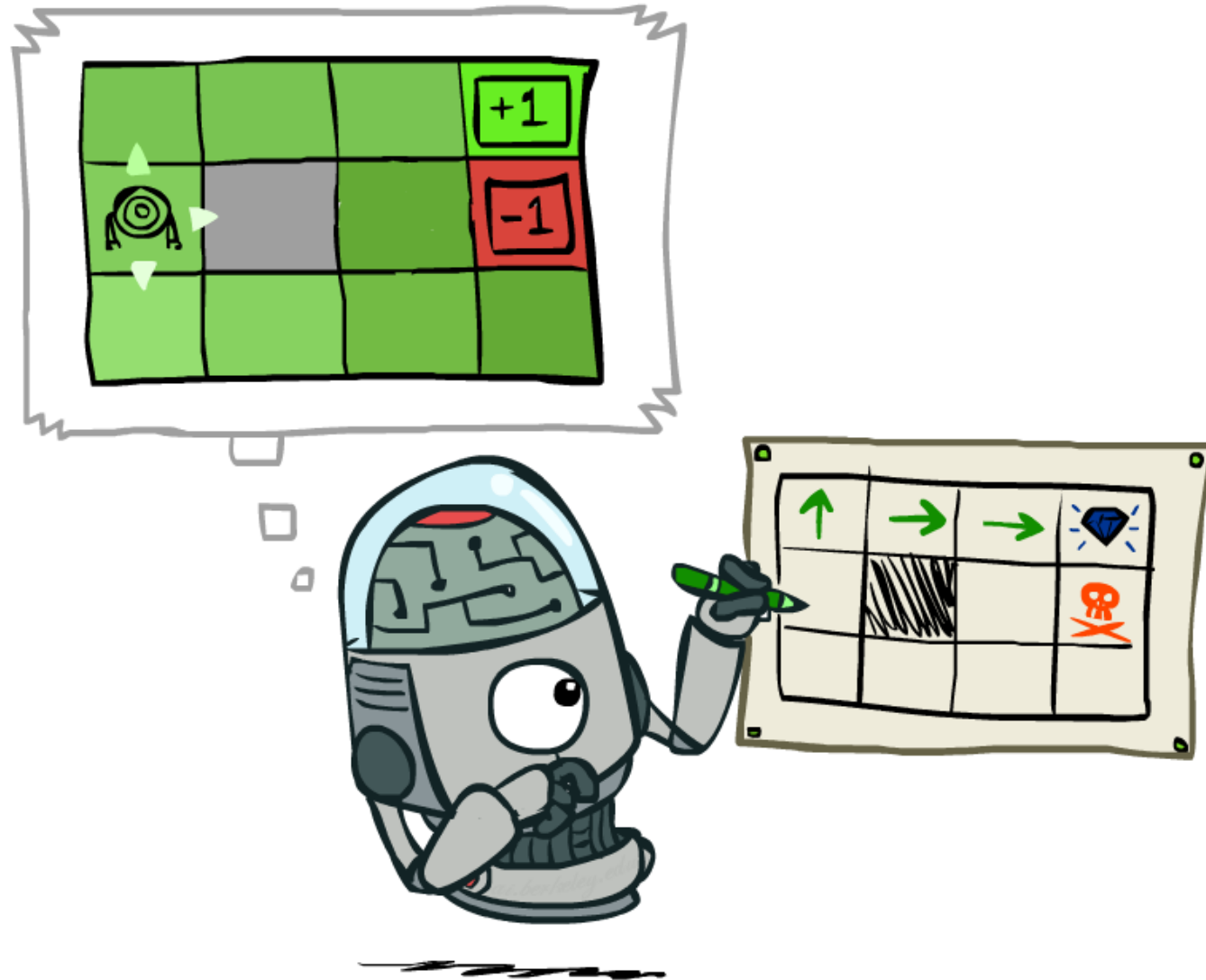
$V^*(s, a, s'')$ : Represent  $V^*(s'')$  where  $s''$  is reachable through  $(s, a)$ . Not a standard notation.

$$\pi^*(s) = \operatorname{argmax}_a \sum_{s'} P(s'|s, a) * (R(s, a, s') + \gamma V^*(s')) \neq \operatorname{argmax}_{s'} V^*(s')$$

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



# Computing Optimal Policy from Values



# Computing Optimal Policy from Values

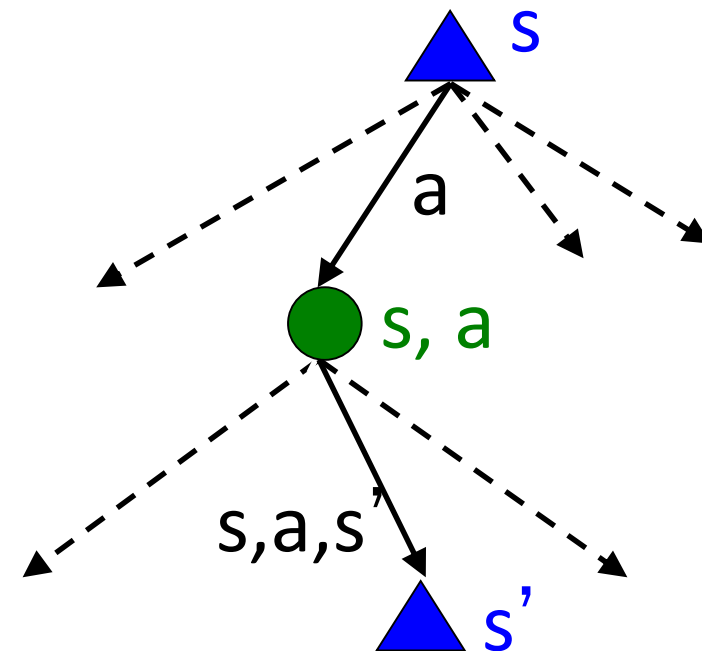
Let's imagine we have the optimal values  $V^*(s)$

How should we act?

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

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

Sometimes this is called policy extraction, since it gets the policy implied by the values



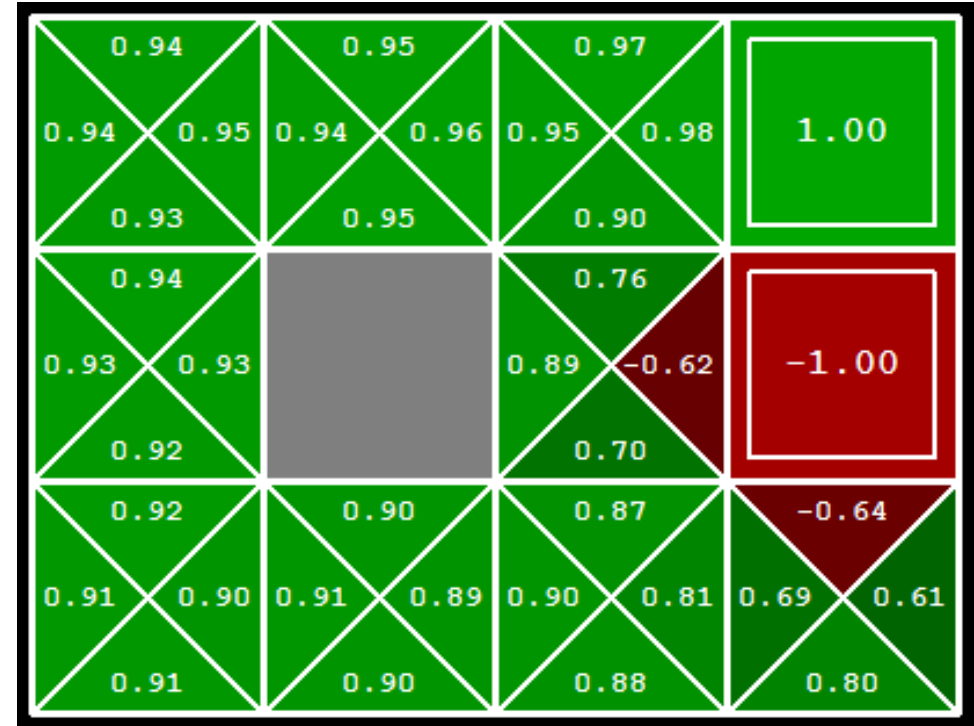
# Computing Optimal Policy from Q-Values

Let's imagine we have the optimal q-values:

How should we act?

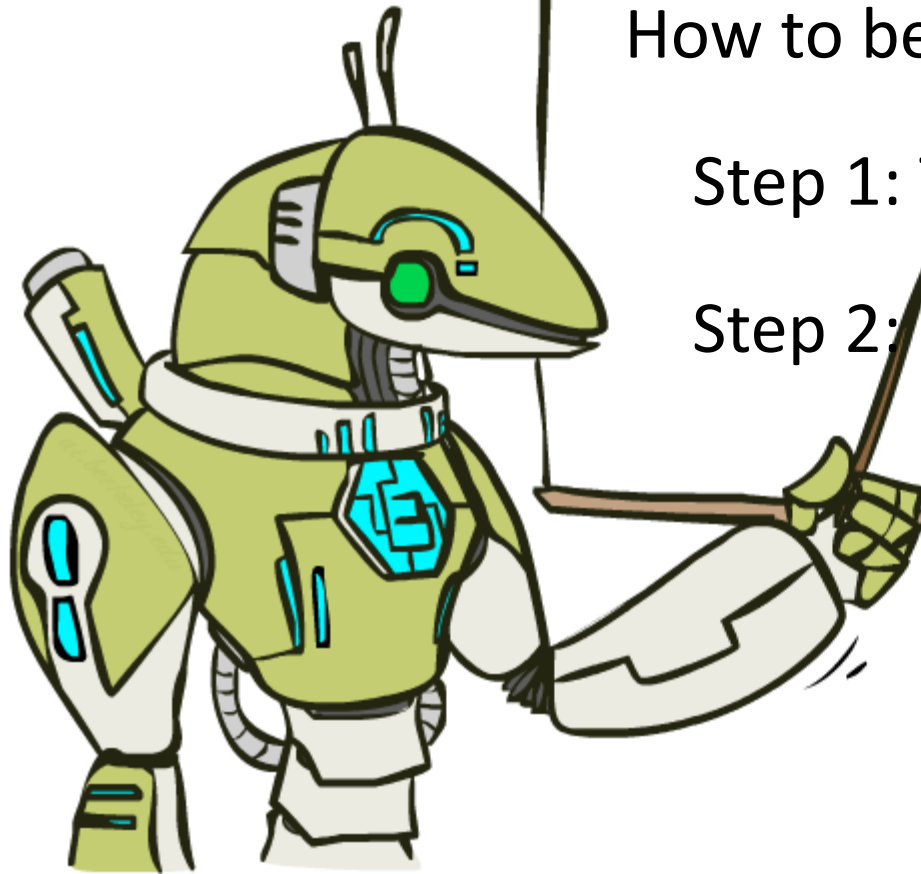
- Completely trivial to decide!

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



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

# The Bellman Equations



How to be optimal:

Step 1: Take correct first action

Step 2: Keep being optimal

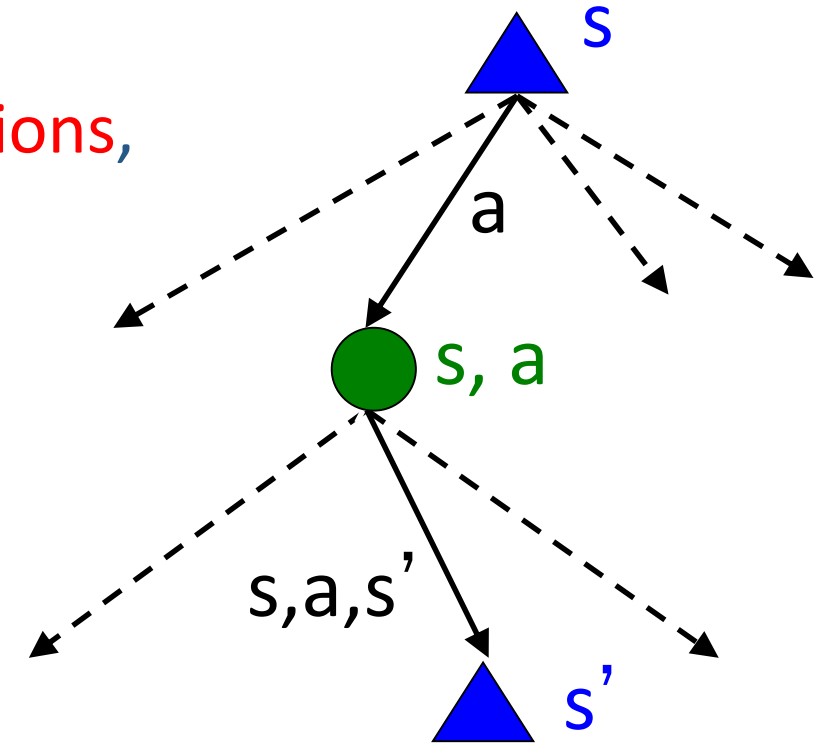
# The Bellman Equations

Definition of “optimal utility” leads to **Bellman Equations**, which **characterize** the **relationship** amongst optimal utility values

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

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

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



Expectimax-like computation

**Necessary and sufficient** conditions for optimality

Solution is **unique**

# The Bellman Equations

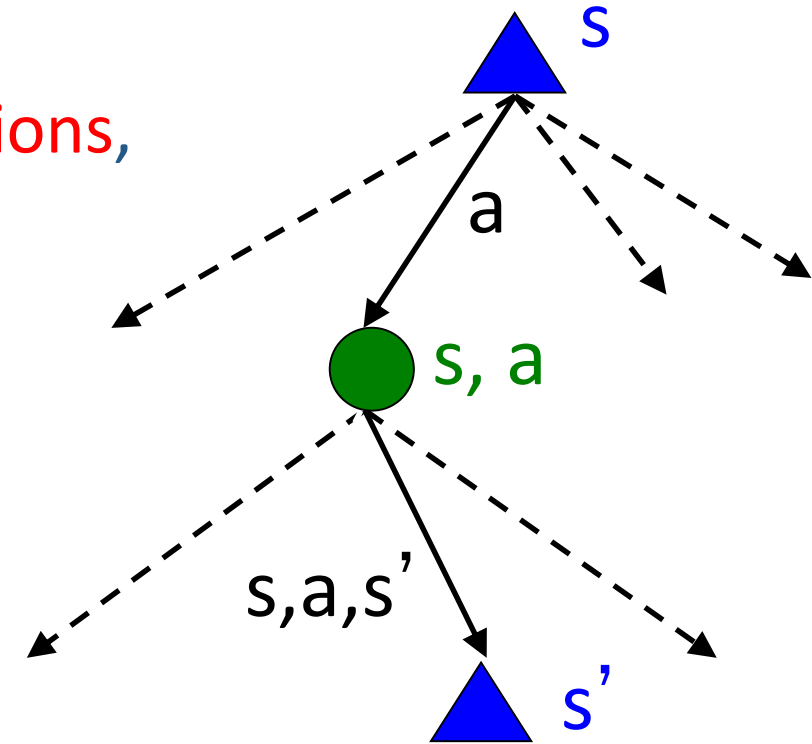
Definition of “optimal utility” leads to **Bellman Equations**, which **characterize** the **relationship** amongst optimal utility values

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

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

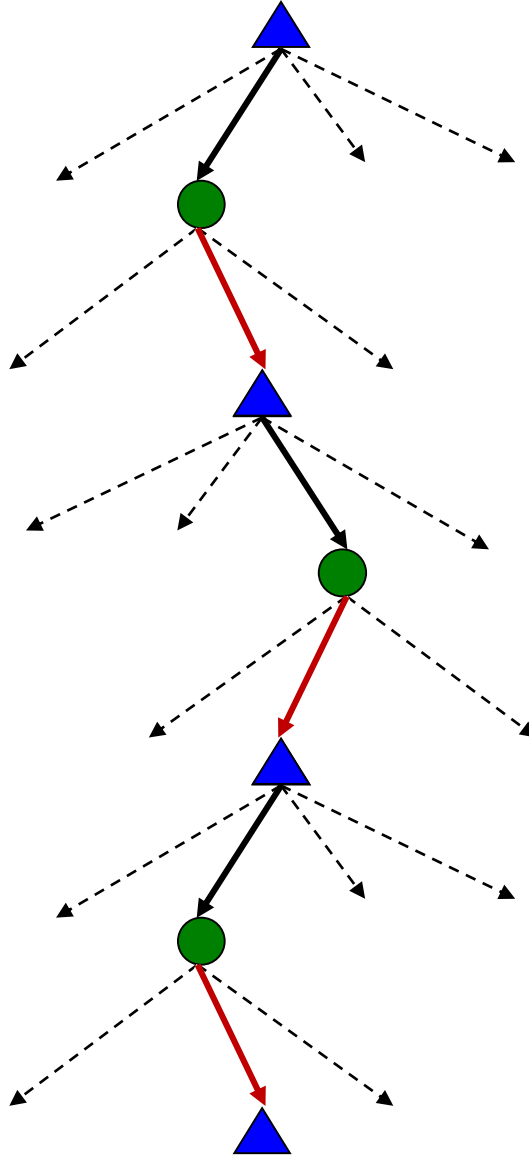
$$V^*(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

**Necessary and sufficient** conditions for optimality  
Solution is **unique**



Expectimax-like computation with one-step lookahead and a “perfect” heuristic at leaf nodes

# Solving Expectimax



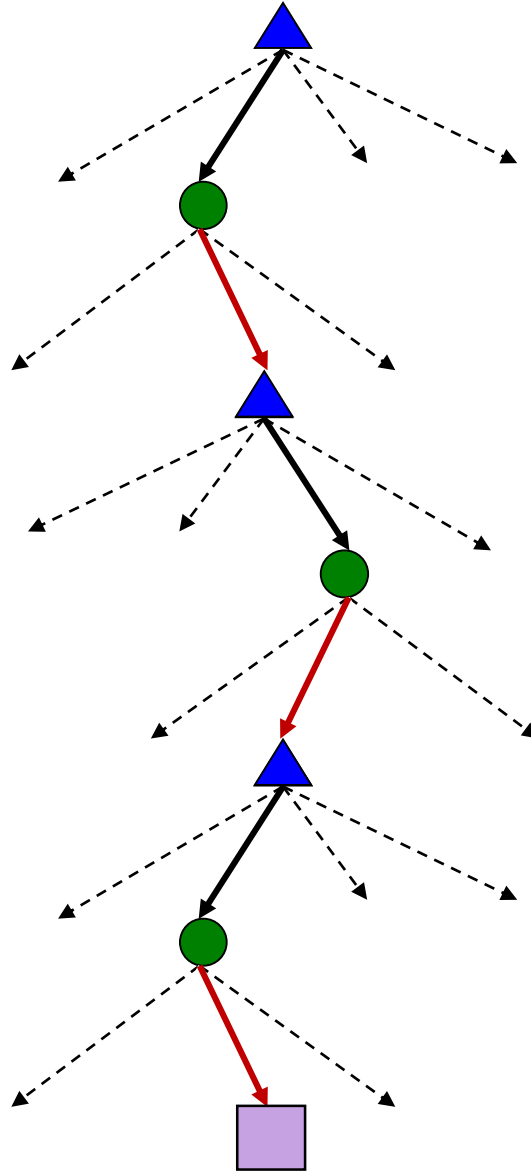
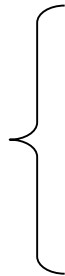


# Solving MDP

Limited Lookahead



$\gamma^2$

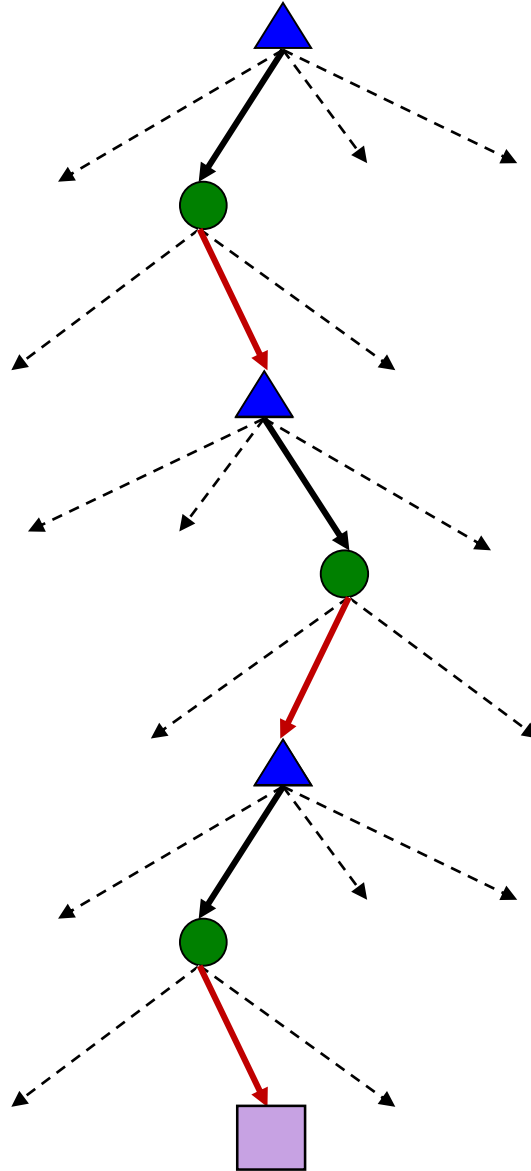
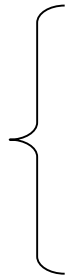


# Solving MDP

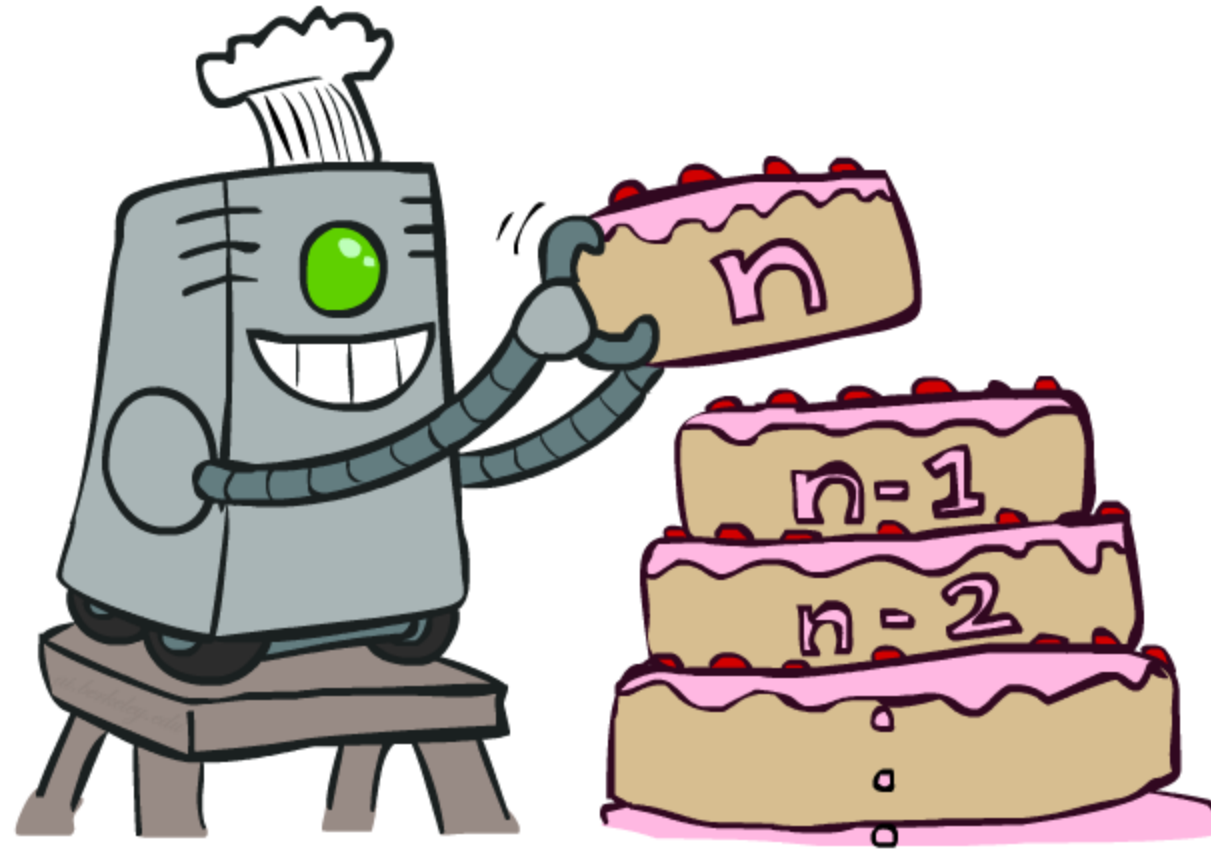
Limited Lookahead



$\gamma^2$



# Value Iteration



# Demo Value Iteration



# Value Iteration

Start with  $V_0(s) = 0$ : no time steps left means an expected reward sum of zero

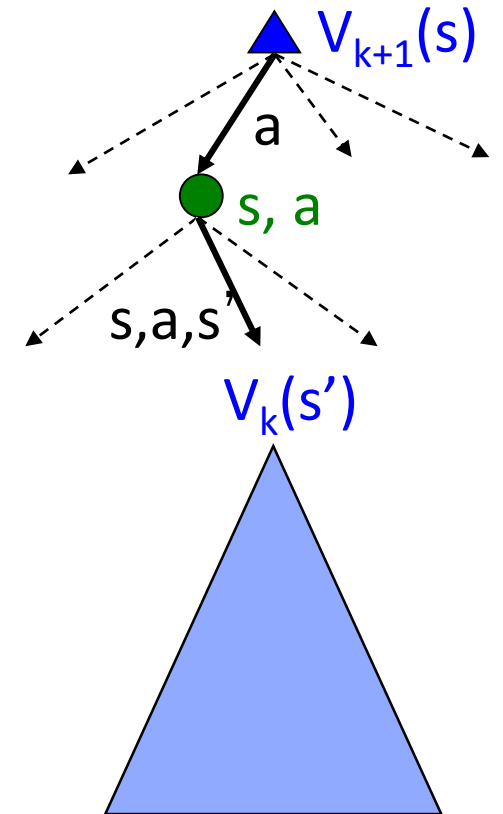
Given vector of  $V_k(s)$  values, apply **Bellman update** once (do one ply of expectimax with  $R$  and  $\gamma$  from **each** state):

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

Repeat until convergence

**Will this process converge?**

**Yes!**



# Piazza Poll 2

What is the complexity of each iteration in Value Iteration?

S -- set of states; A -- set of actions

I:  $O(|S||A|)$

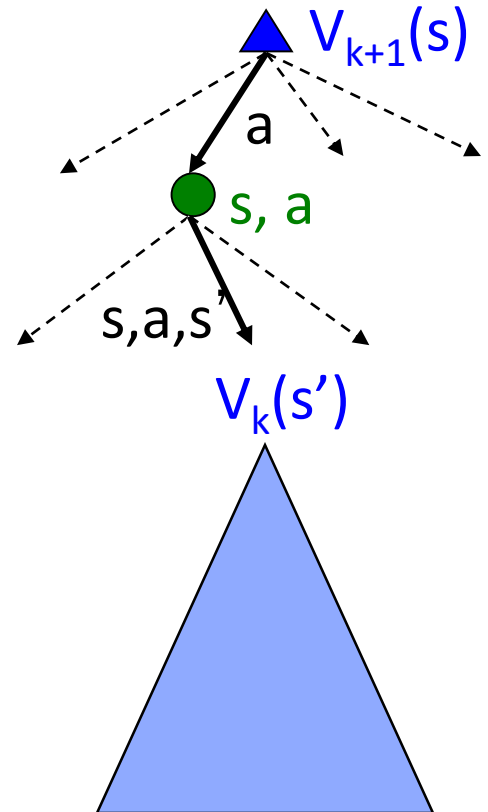
II:  $O(|S|^2|A|)$

III:  $O(|S||A|^2)$

IV:  $O(|S|^2|A|^2)$

V:  $O(|S|^2)$

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



# Piazza Poll 2

What is the complexity of each iteration in Value Iteration?

S -- set of states; A -- set of actions

I:  $O(|S||A|)$

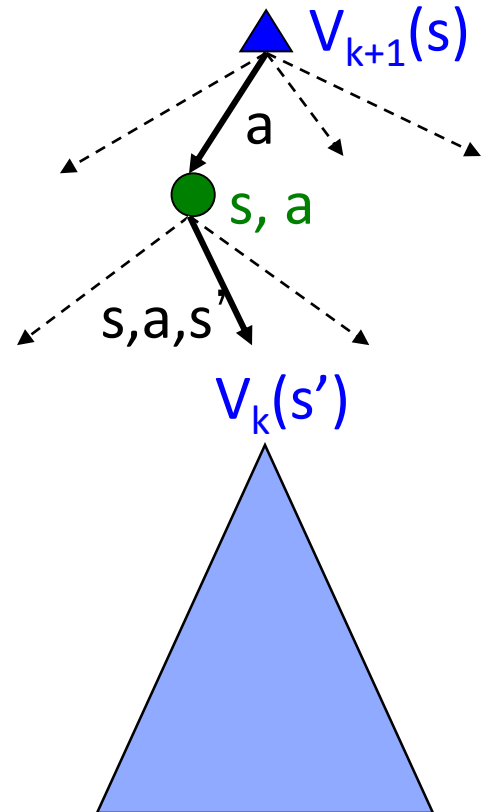
II:  $O(|S|^2|A|)$

III:  $O(|S||A|^2)$

IV:  $O(|S|^2|A|^2)$

V:  $O(|S|^2)$

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



# Value Iteration

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

function **VALUE-ITERATION**(MDP=(S,A,T,R, $\gamma$ ), *threshold*) returns a state value function

for  $s$  in  $S$

$V_0(s) \leftarrow 0$

$k \leftarrow 0$

repeat

$\delta \leftarrow 0$

for  $s$  in  $S$

$V_{k+1}(s) \leftarrow -\infty$

for  $a$  in  $A$

$v \leftarrow 0$

for  $s'$  in  $S$

$v \leftarrow v + T(s, a, s')(R(s, a, s') + \gamma V_k(s'))$

$V_{k+1}(s) \leftarrow \max\{V_{k+1}(s), v\}$

$\delta \leftarrow \max\{\delta, |V_{k+1}(s) - V_k(s)|\}$

$k \leftarrow k + 1$

until  $\delta < \textit{threshold}$

return  $V_{k-1}$

Do we really need to store the value of  $V_k$  for each  $k$ ?

Does  $V_{k+1}(s) \geq V_k(s)$  always hold?



# Value Iteration

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

function **VALUE-ITERATION**(MDP=(S,A,T,R, $\gamma$ ), *threshold*) returns a state value function

for  $s$  in  $S$

$V_0(s) \leftarrow 0$

$k \leftarrow 0$

repeat

$\delta \leftarrow 0$

for  $s$  in  $S$

$V_{k+1}(s) \leftarrow -\infty$

for  $a$  in  $A$

$v \leftarrow 0$

for  $s'$  in  $S$

$v \leftarrow v + T(s, a, s')(R(s, a, s') + \gamma V_k(s'))$

$V_{k+1}(s) \leftarrow \max\{V_{k+1}(s), v\}$

$\delta \leftarrow \max\{\delta, |V_{k+1}(s) - V_k(s)|\}$

$k \leftarrow k + 1$

until  $\delta < \textit{threshold}$

return  $V_{k-1}$

Do we really need to store the value of  $V_k$  for each  $k$ ?

No. Use  $V = V_{last}$  and  $V' = V_{current}$

Does  $V_{k+1}(s) \geq V_k(s)$  always hold?

No. If  $T(s, a, s') = 1$  and  $R(s, a, s') < 0$ , then  $V_1(s) = R(s, a, s') < 0$

# Bellman Equation vs Value Iteration vs Bellman Update

Bellman equations **characterize** the optimal values:

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

Value iteration **computes** them by applying Bellman update repeatedly

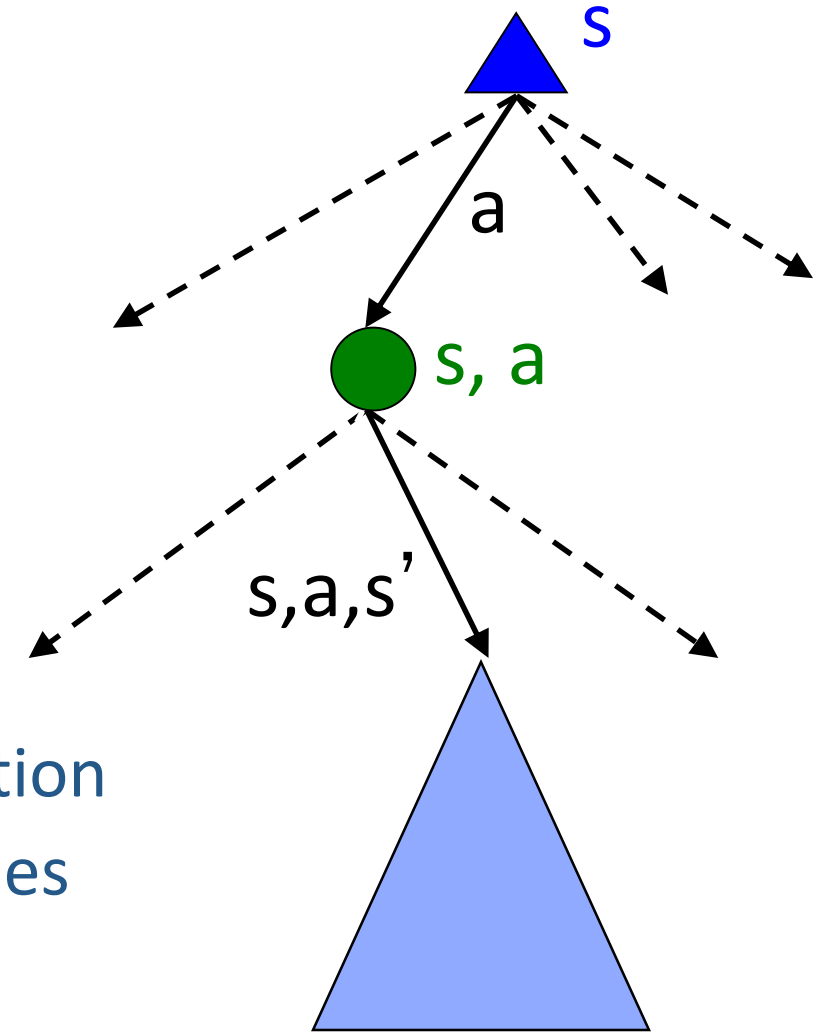
$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k(s')]$$

Value iteration is a method for solving Bellman Equation

$V_k$  vectors are also interpretable as time-limited values

Value iteration finds the fixed point of the function

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



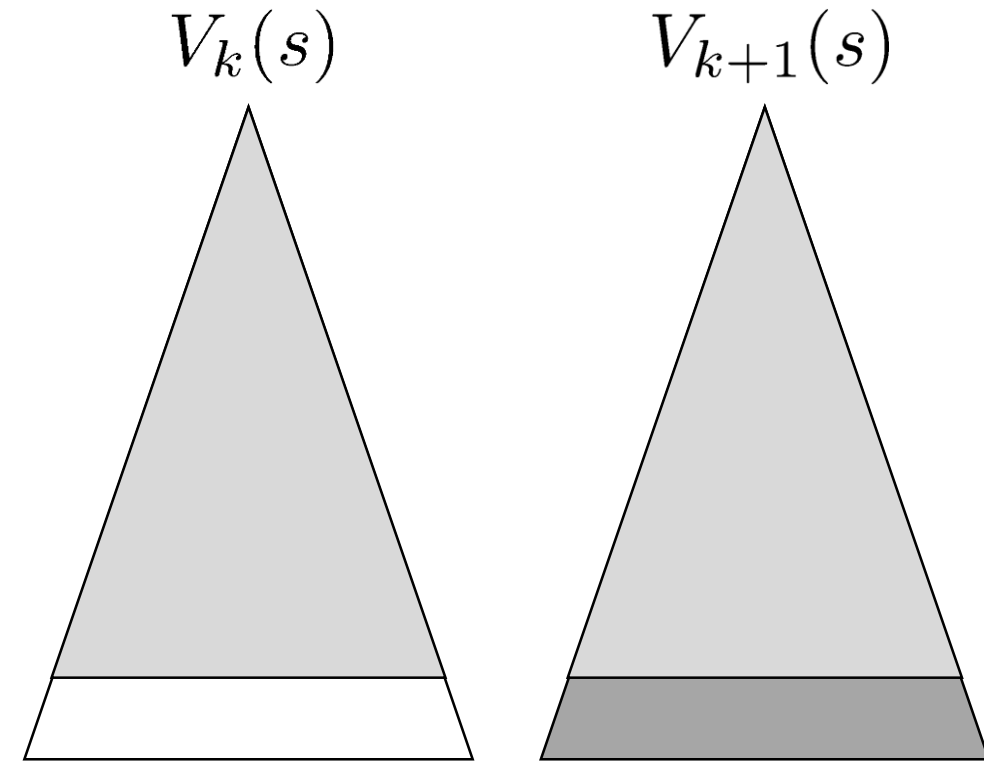
# Value Iteration Convergence

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

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

Case 2: If  $\gamma < 1$  and  $|R(s, a, s')| \leq R_{max} < \infty$

- Intuition: For any state  $V_k$  and  $V_{k+1}$  can be viewed as depth  $k+1$  expectimax results (with  $R$  and  $\gamma$ ) in nearly identical search trees
- The difference is that on the bottom layer,  $V_{k+1}$  has actual rewards while  $V_k$  has zeros
- $|R(s, a, s')| \leq R_{max}$
- $|V_1(s) - V_0(s)| = |V_1(s) - 0| \leq R_{max}$
- $|V_{k+1}(s) - V_k(s)| \leq \gamma^k R_{max}$
- So  $|V_{k+1}(s) - V_k(s)| \rightarrow 0$  as  $k \rightarrow \infty$



If we initialized  $V_0(s)$  differently, what would happen?

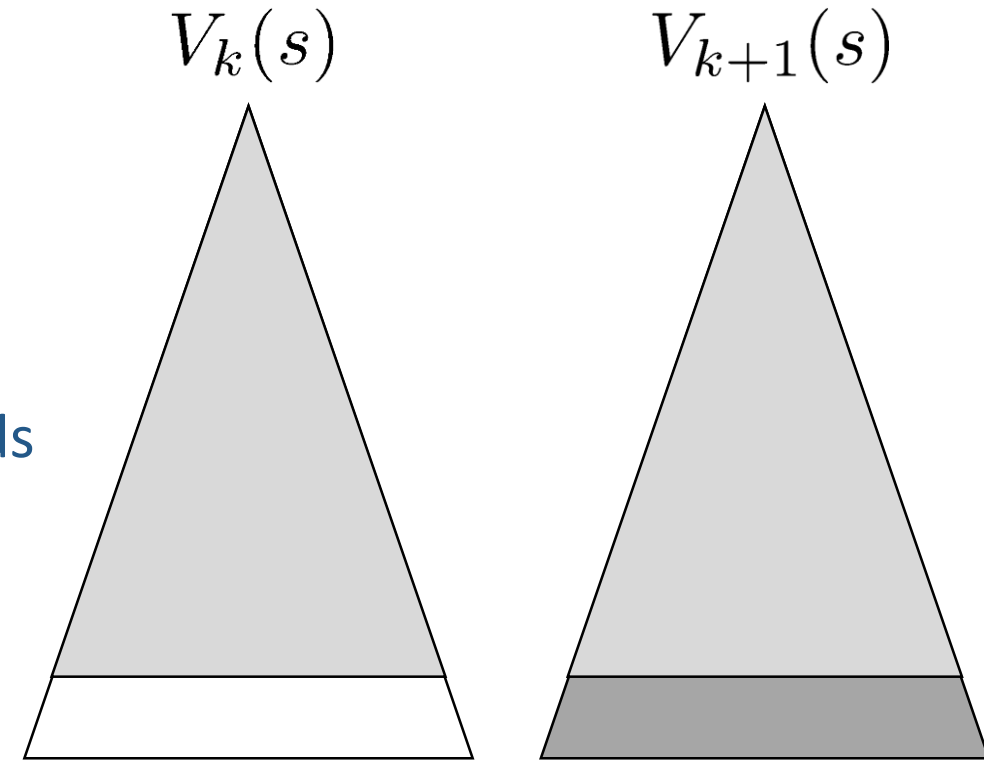
# Value Iteration Convergence

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

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

Case 2: If  $\gamma < 1$  and  $|R(s, a, s')| \leq R_{max} < \infty$

- Intuition: For any state  $V_k$  and  $V_{k+1}$  can be viewed as depth  $k+1$  expectimax results (with  $R$  and  $\gamma$ ) in nearly identical search trees
- The difference is that on the bottom layer,  $V_{k+1}$  has actual rewards while  $V_k$  has zeros
- Each value at last layer of  $V_{k+1}$  tree is at most  $R_{max}$  in magnitude
- But everything is discounted by  $\gamma^k$  that far out
- So  $V_k$  and  $V_{k+1}$  are at most  $\gamma^k R_{max}$  different
- So as  $k$  increases, the values converge



If we initialized  $V_0(s)$  differently, what would happen?

Still converge to  $V^*(s)$  as long as  $|V_0(s)| < +\infty$ , but may be slower

# Other ways to solve Bellman Equation?

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

Treat  $V^*(s)$  as variables

Solve Bellman Equation through Linear Programming

# Other ways to solve Bellman Equation?

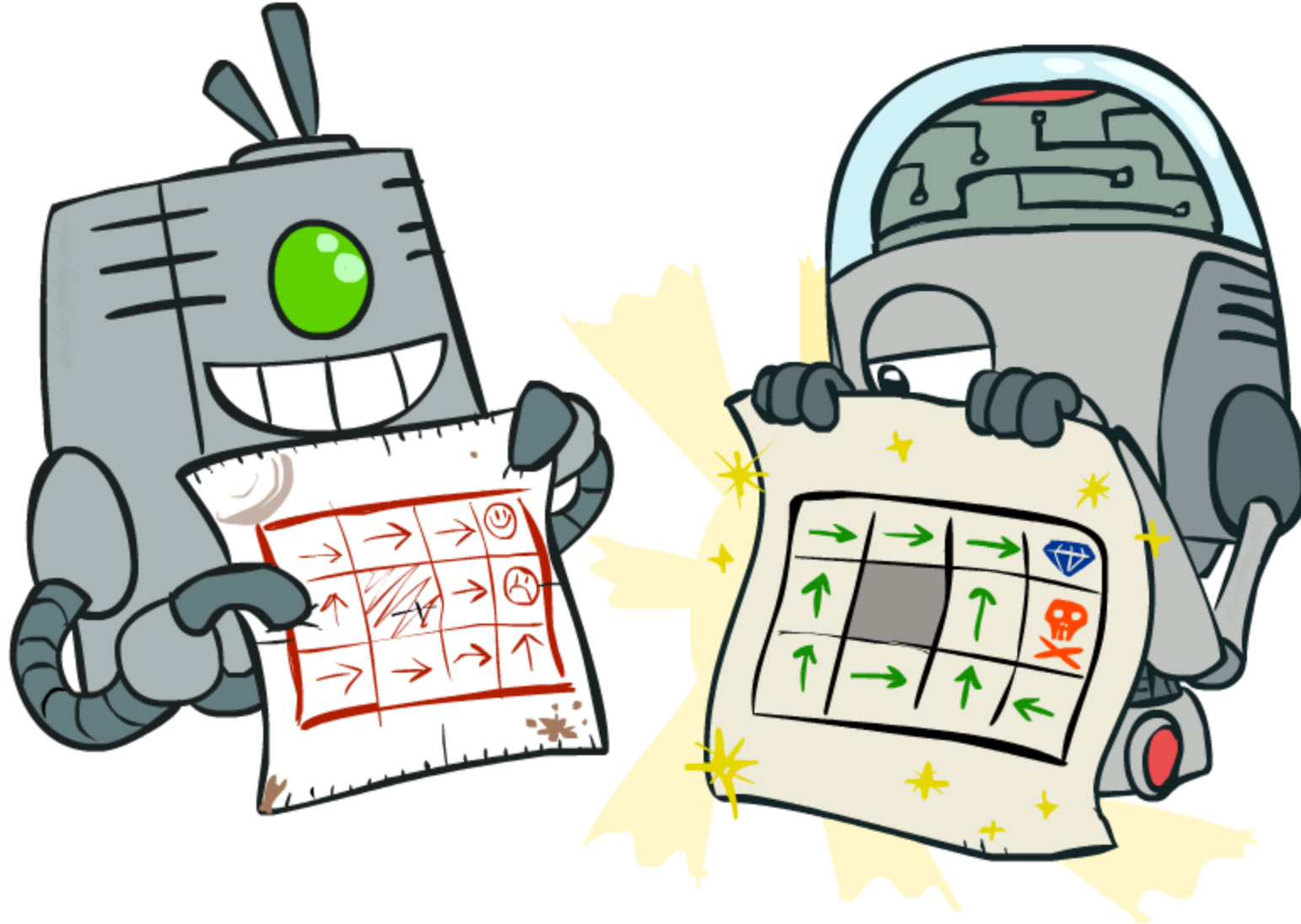
$$V^*(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$$

Treat  $V^*(s)$  as variables

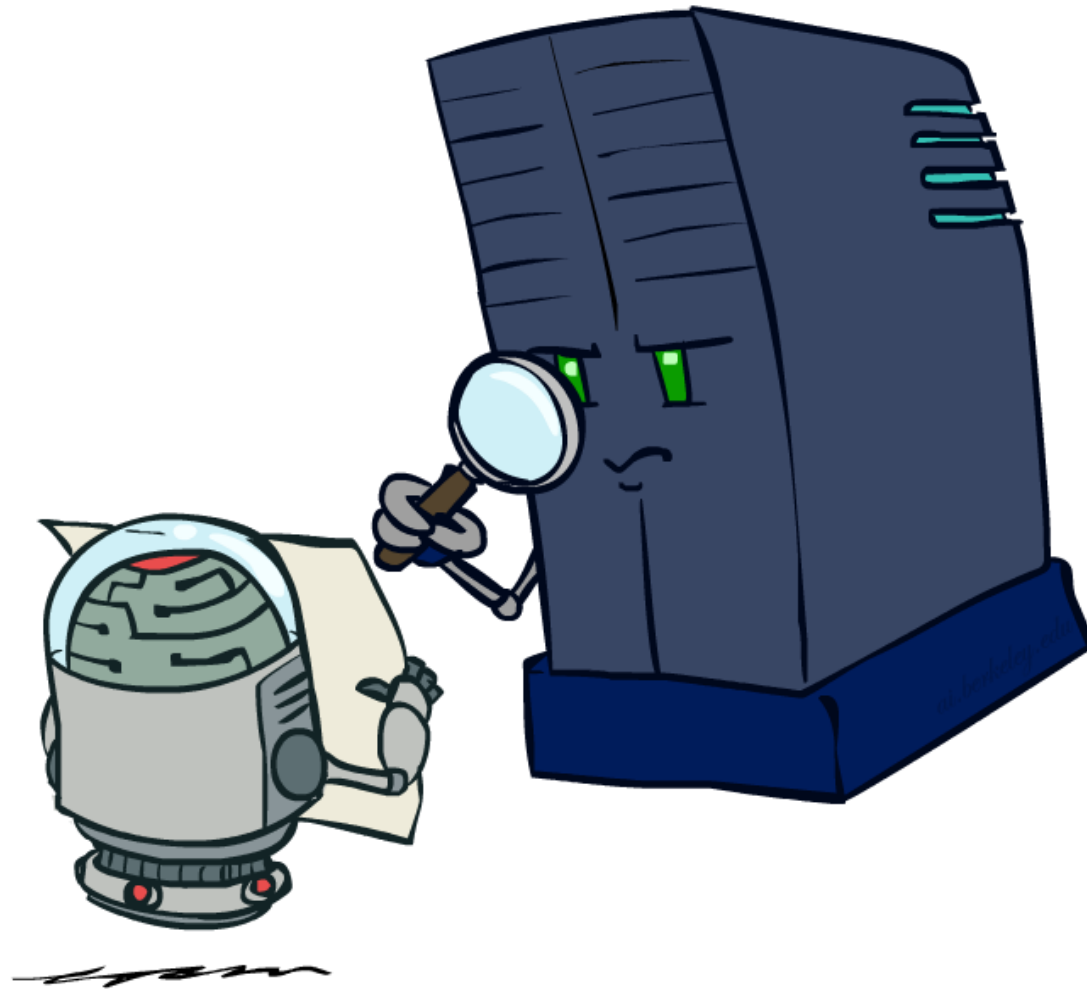
Solve Bellman Equation through Linear Programming

$$\begin{aligned} & \min_{V^*} \sum_s V^*(s) \\ & \text{s.t. } V^*(s) \geq \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')], \forall s, a \end{aligned}$$

# Policy Iteration for Solving MDPs



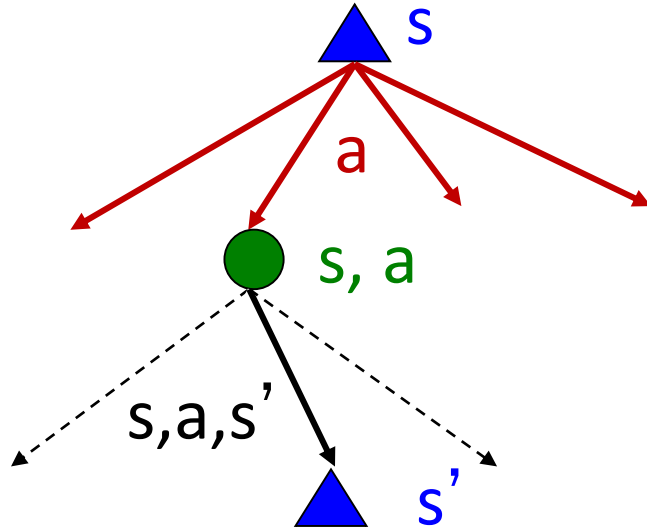
# Policy Evaluation



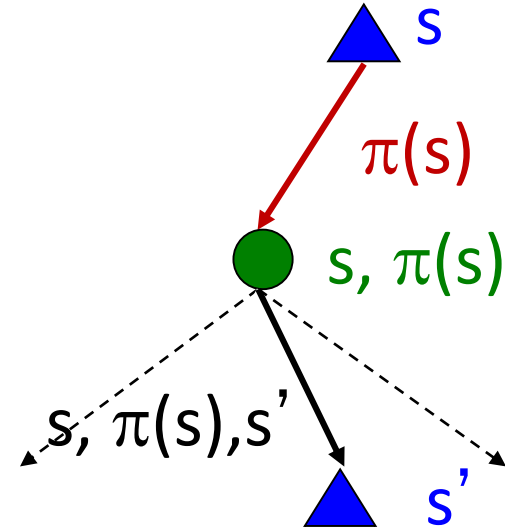


# Fixed Policies

Do the optimal action



Do what  $\pi$  says to do



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

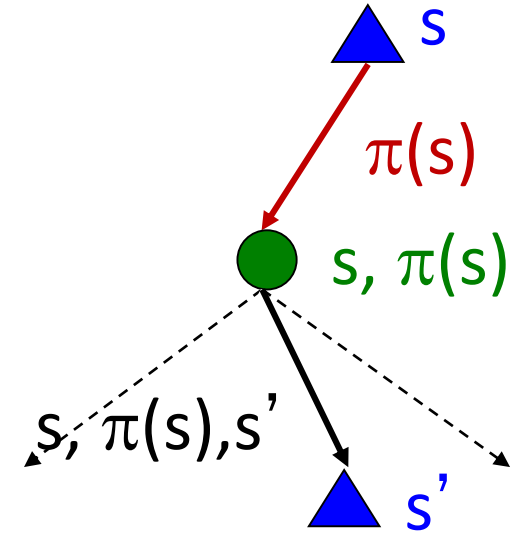
- ... though the tree's value would depend on which policy we fixed

# Utilities for a Fixed Policy

Another 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 (one-step look-ahead / Bellman equation):

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

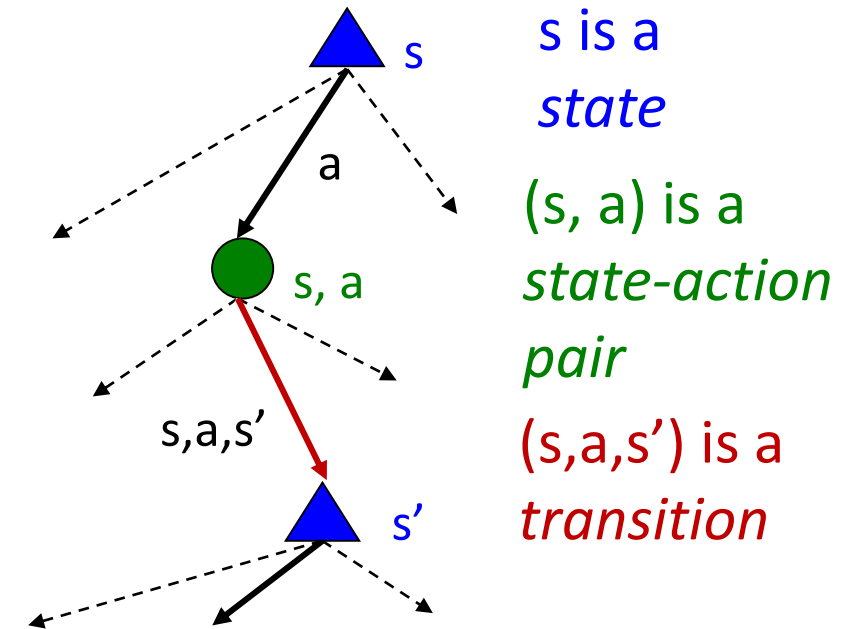
# Compare

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

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

# MDP Quantities

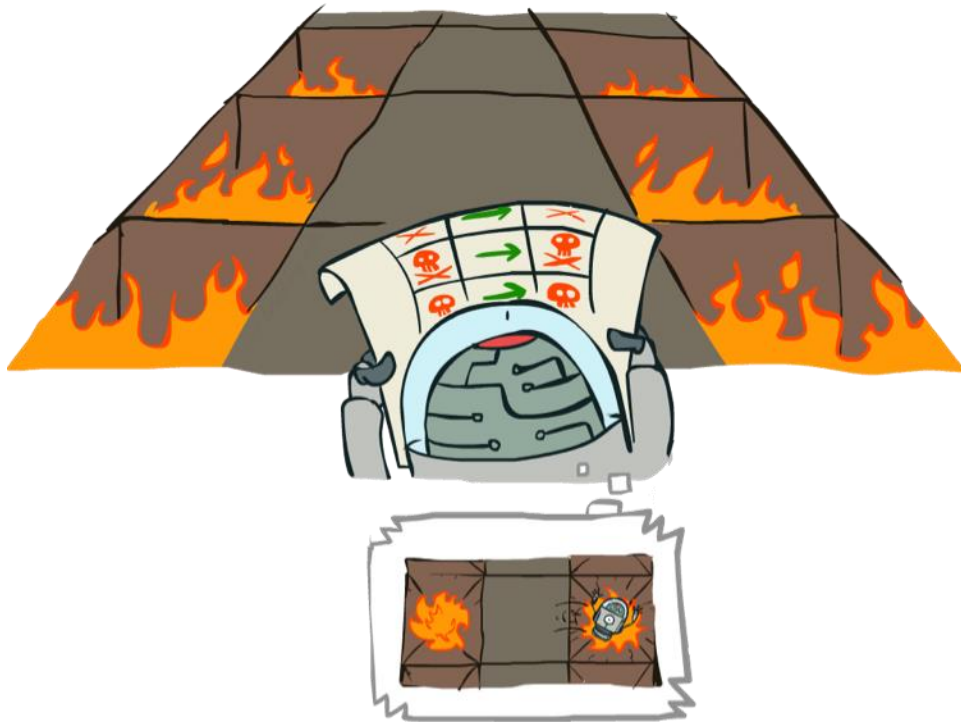
- A policy  $\pi$ : map of states to actions
- The optimal policy  $\pi^*$ :  $\pi^*(s)$  = optimal action from state  $s$
- Value function of a policy  $V^\pi(s)$ : expected utility starting in  $s$  and acting according to  $\pi$
- Optimal value function  $V^*$ :  $V^*(s) = V^{\pi^*}(s)$
- Q function of a policy  $Q^\pi(s)$ : expected utility starting out having taken action  $a$  from state  $s$  and (thereafter) acting according to  $\pi$
- Optimal Q function  $Q^*$  :  $Q^*(s,a) = Q^{\pi^*}(s,a)$



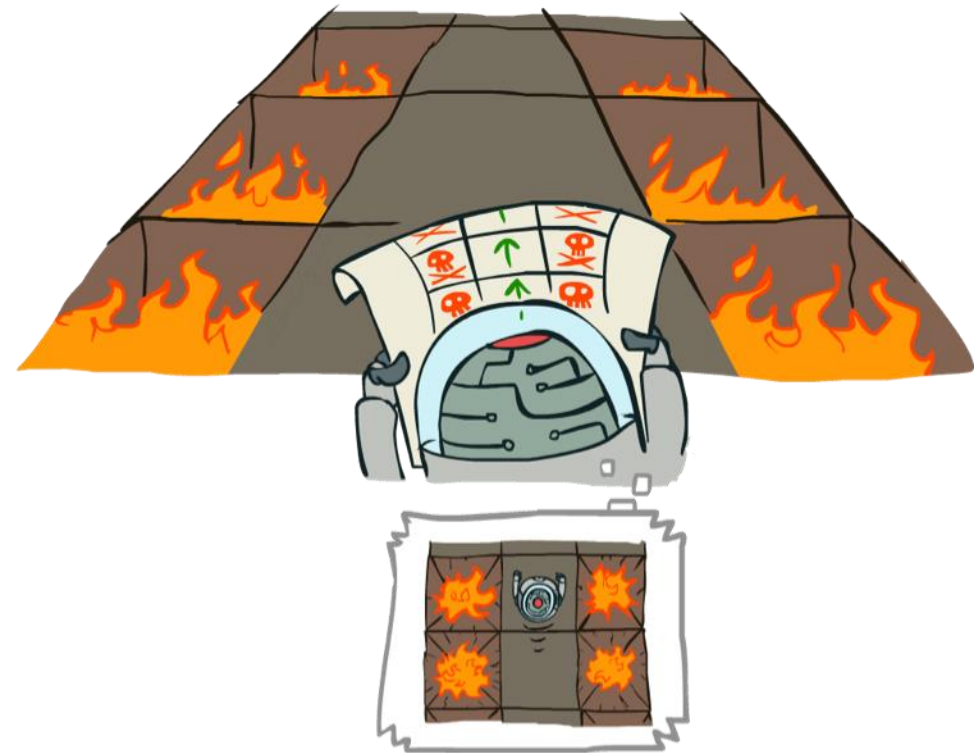
Solve MDP: Find  $\pi^*$ ,  $V^*$  and/or  $Q^*$

# Example: Policy Evaluation

Always Go Right



Always Go Forward



# Example: Policy Evaluation

Always Go Right



Always Go Forward



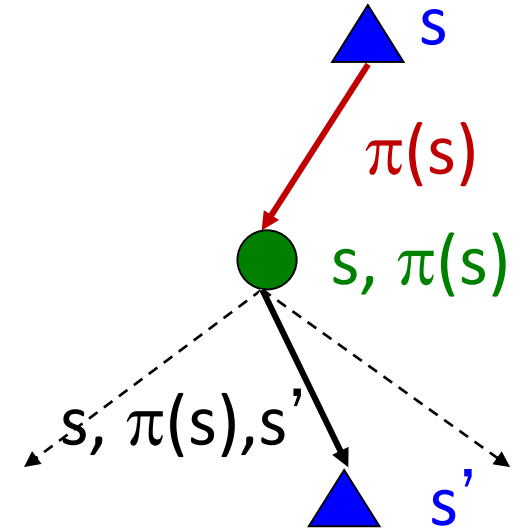
# Policy Evaluation

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

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

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

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



# Piazza Poll 3

What is the complexity of each iteration in Policy Evaluation?

S -- set of states; A -- set of actions

I:  $O(|S||A|)$

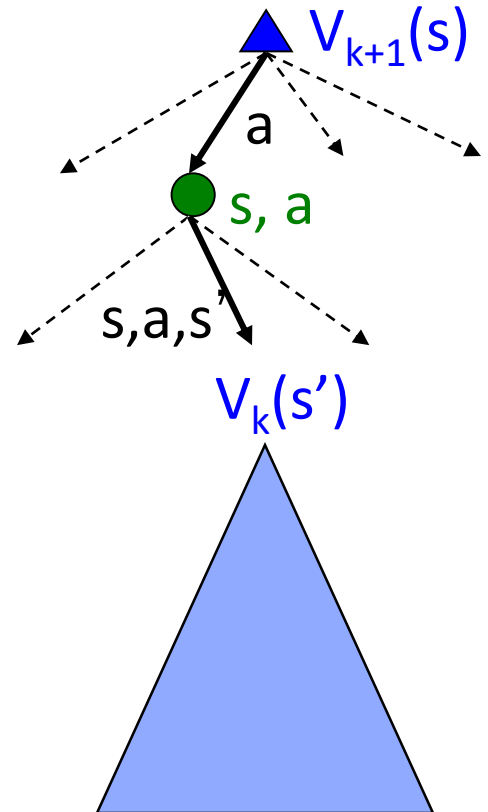
II:  $O(|S|^2|A|)$

III:  $O(|S||A|^2)$

IV:  $O(|S|^2|A|^2)$

V:  $O(|S|^2)$

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





# Piazza Poll 3

What is the complexity of each iteration in Policy Evaluation?

S -- set of states; A -- set of actions

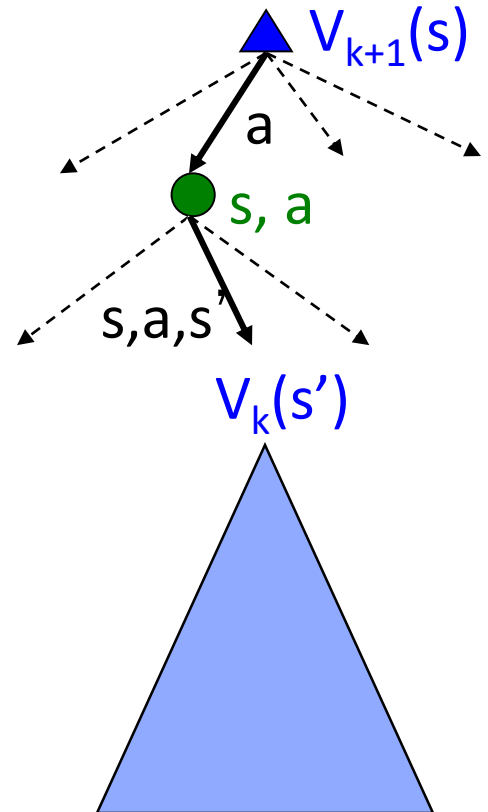
I:  $O(|S||A|)$

II:  $O(|S|^2|A|)$

III:  $O(|S||A|^2)$

IV:  $O(|S|^2|A|^2)$

V:  $O(|S|^2)$



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

# Policy Evaluation

Idea 2: Bellman Equation w.r.t. a given policy  $\pi$  defines a linear system

- Solve with your favorite linear system solver

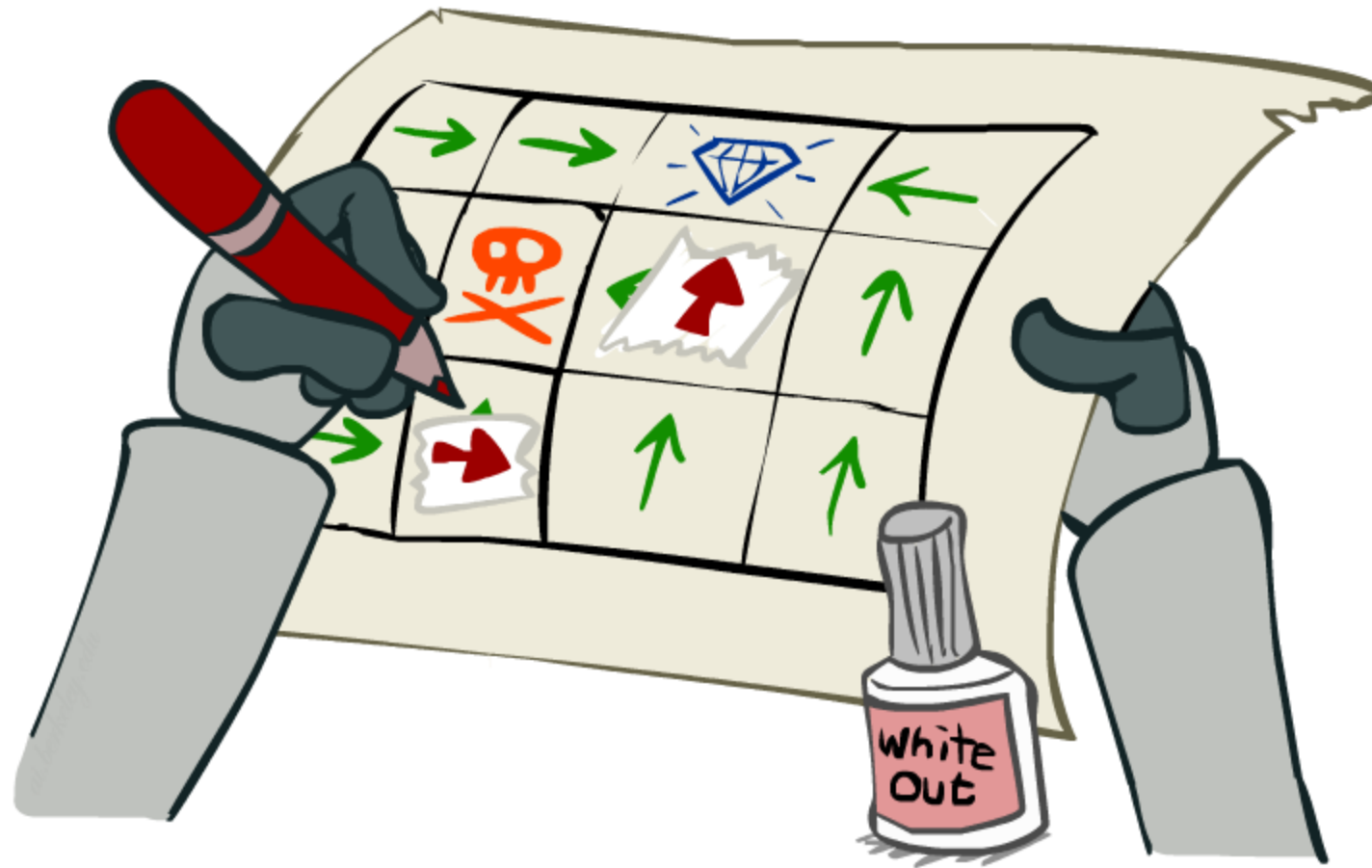
$$V^\pi(s) = \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V^\pi(s')]$$

Treat  $V^\pi(s)$  as variables

How many variables?

How many constraints?

# Policy Iteration



# Problems with Value Iteration

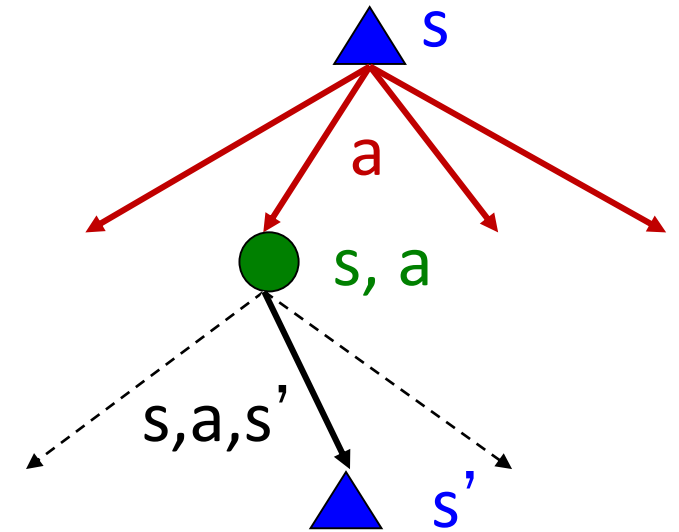
Value iteration repeats the Bellman updates:

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

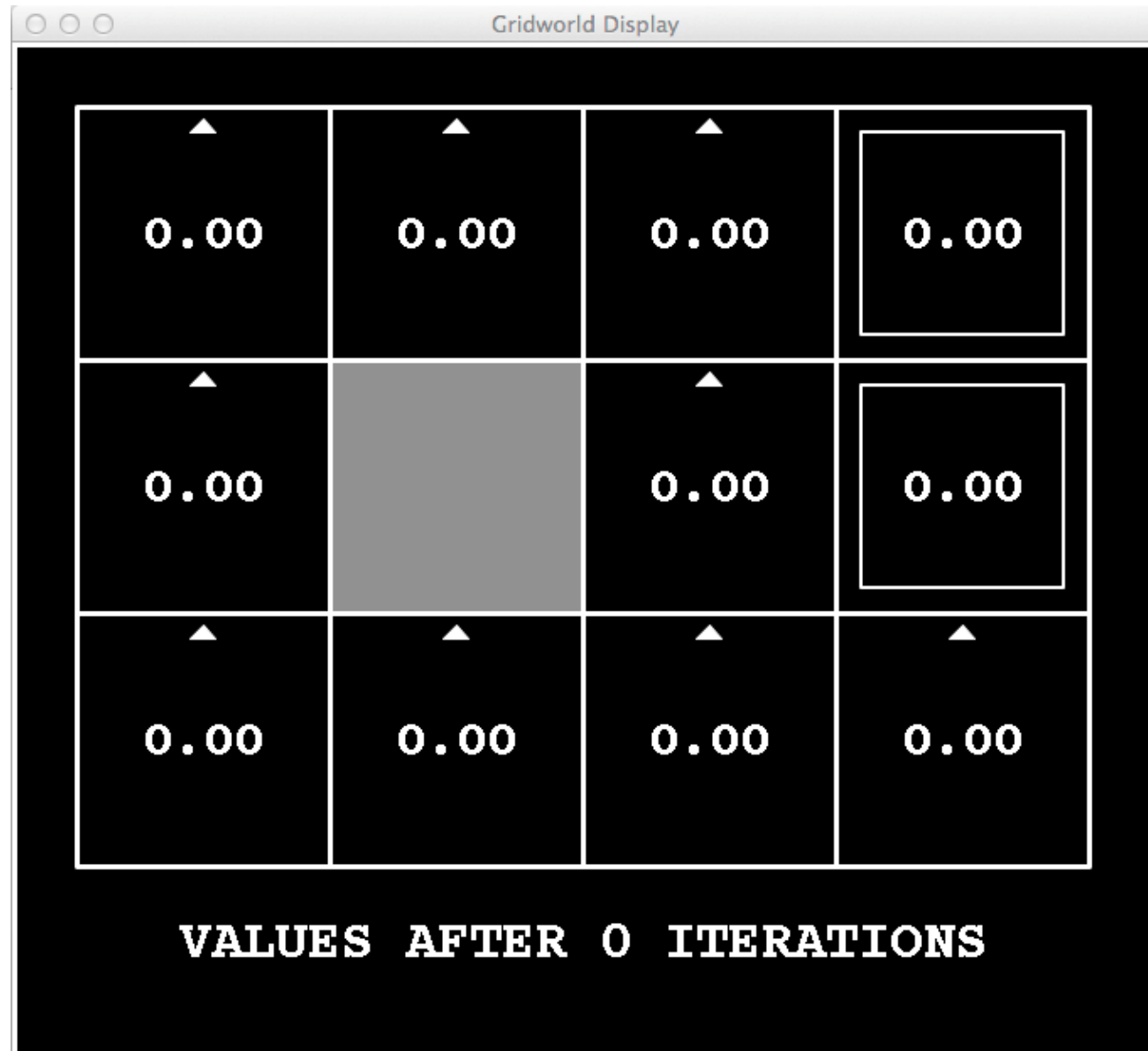
Problem 1: It's slow –  $O(|S|^2|A|)$  per iteration

Problem 2: The “max” at each state rarely changes

Problem 3: The policy often converges long before the values

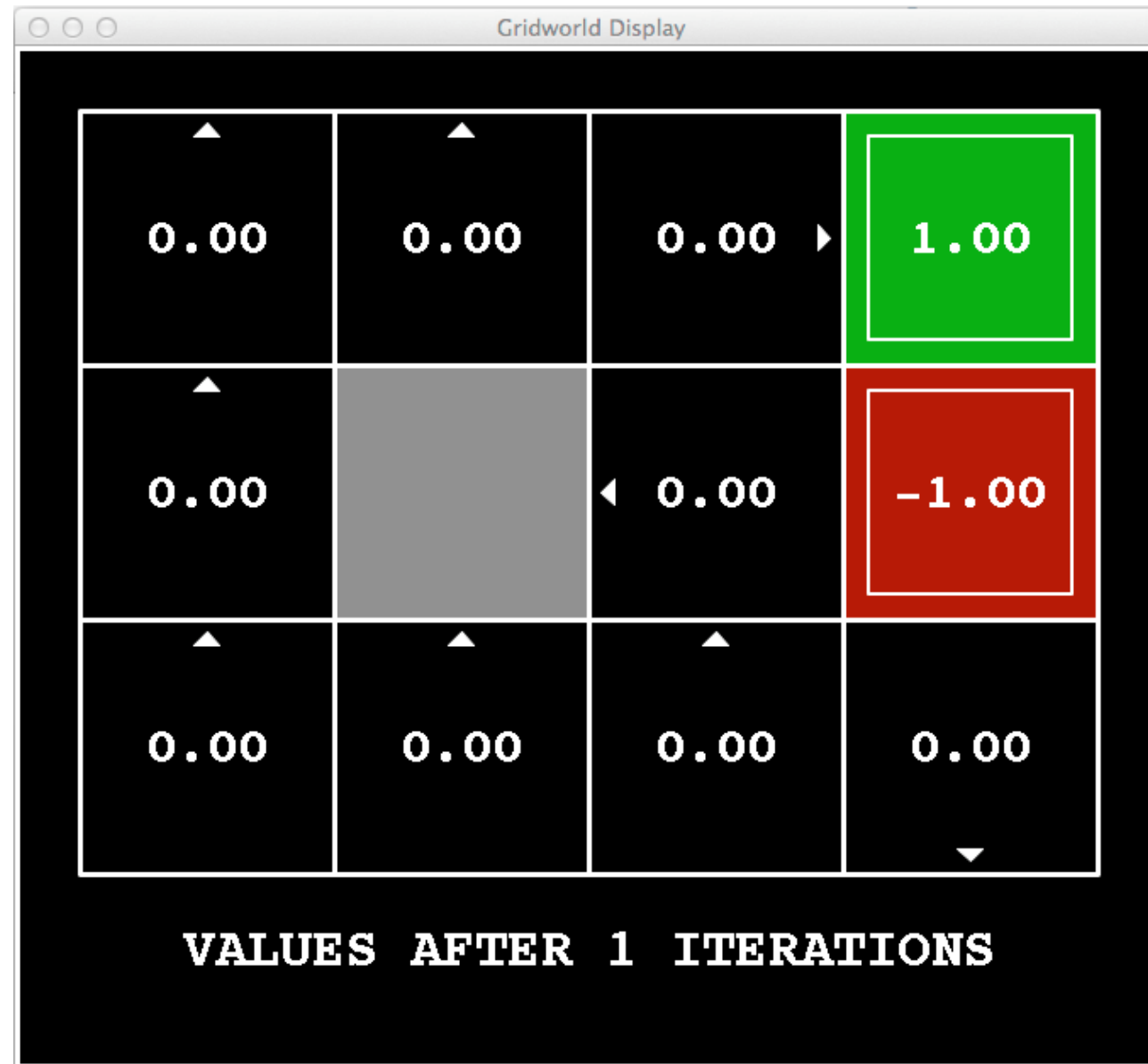


$k=0$



Noise = 0.2  
Discount = 0.9  
Living reward = 0

$k=1$



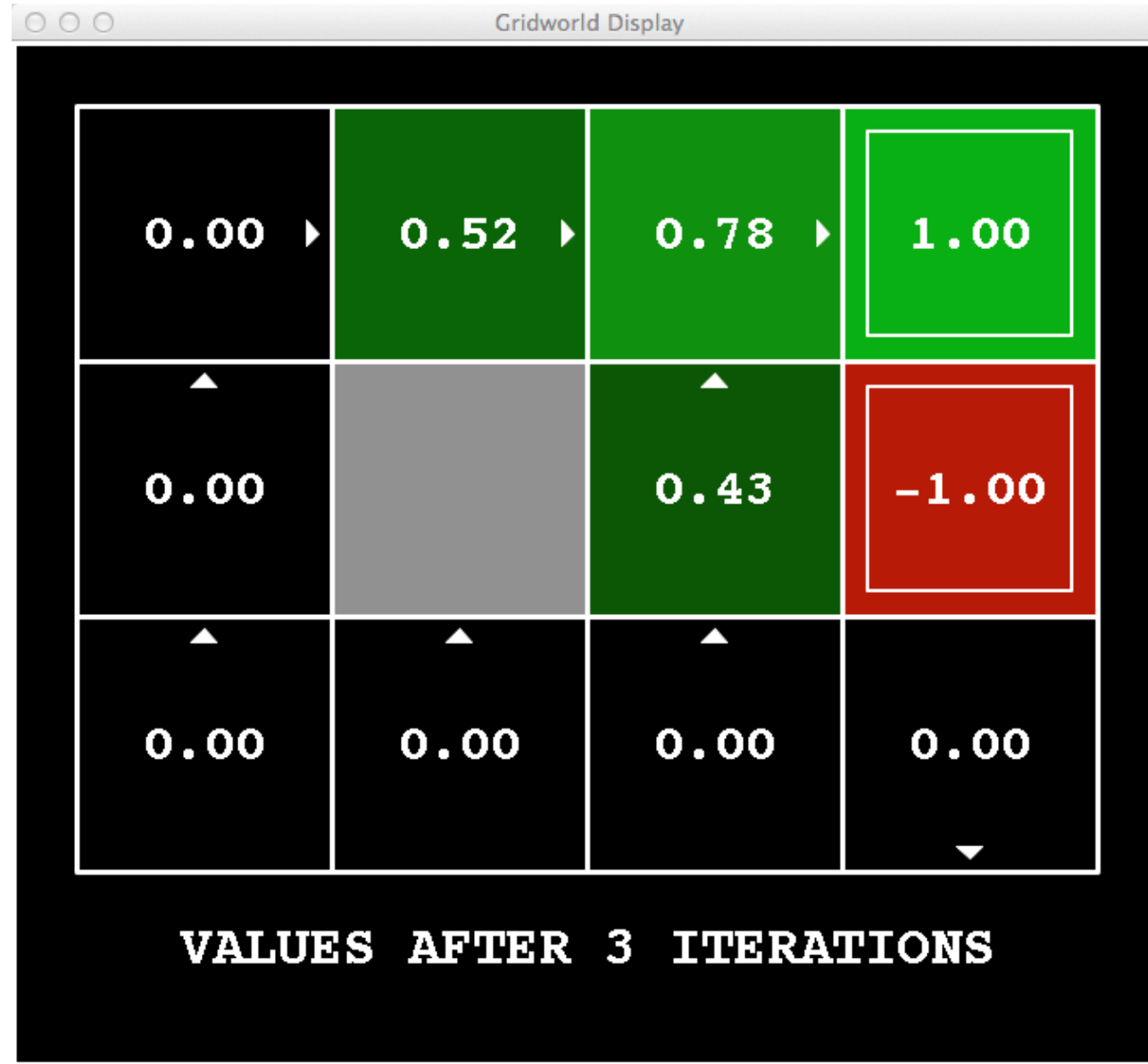
Noise = 0.2  
Discount = 0.9  
Living reward = 0

$k=2$



Noise = 0.2  
Discount = 0.9  
Living reward = 0

$k=3$



Noise = 0.2  
Discount = 0.9  
Living reward = 0



$k=4$



Noise = 0.2  
Discount = 0.9  
Living reward = 0

k=5



Noise = 0.2  
Discount = 0.9  
Living reward = 0

k=6



Noise = 0.2  
Discount = 0.9  
Living reward = 0

$k=7$



Noise = 0.2  
Discount = 0.9  
Living reward = 0

$k=8$



Noise = 0.2  
Discount = 0.9  
Living reward = 0

k=9



Noise = 0.2  
Discount = 0.9  
Living reward = 0

k=10



Noise = 0.2  
Discount = 0.9  
Living reward = 0

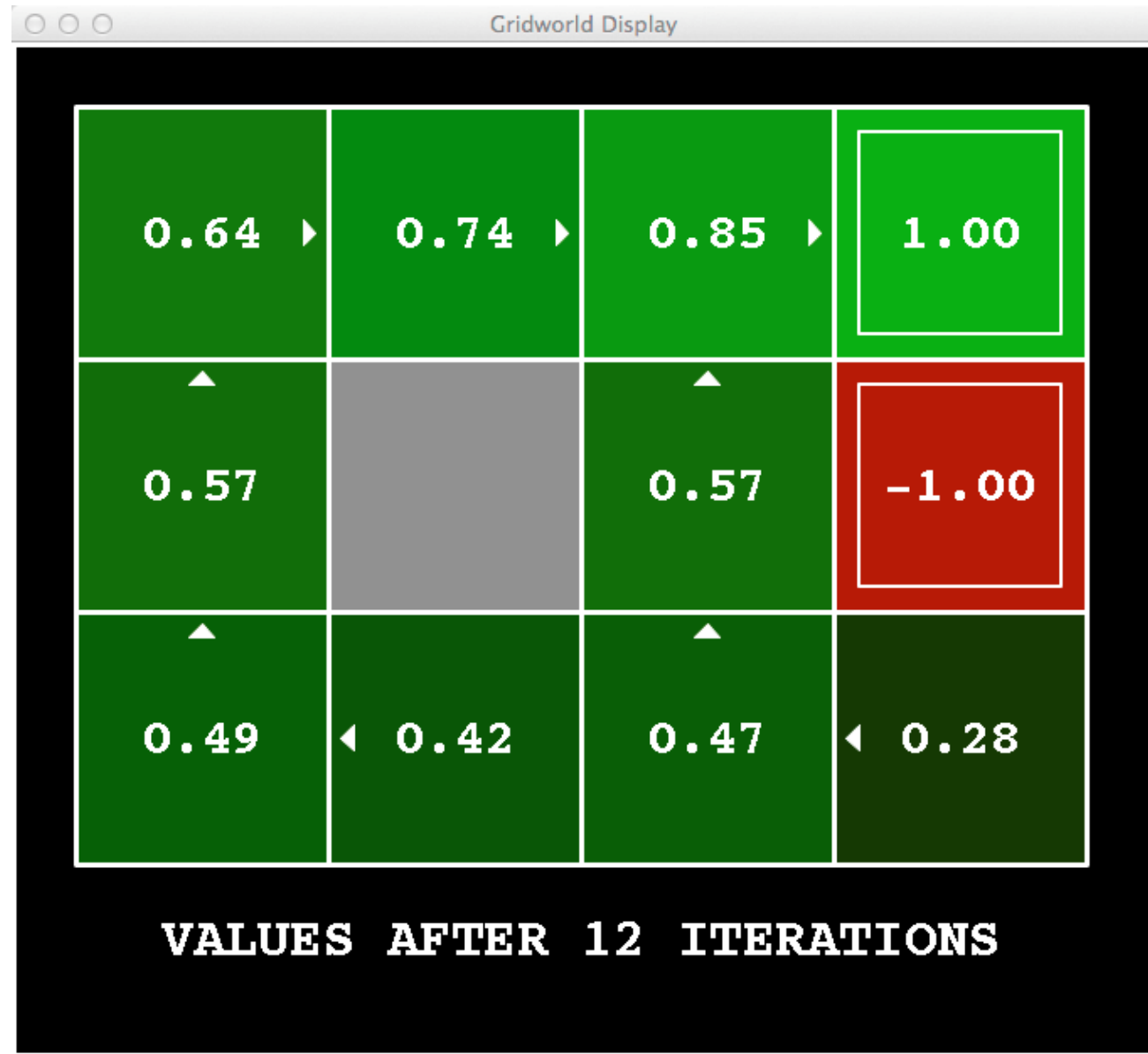
k=11



Noise = 0.2  
Discount = 0.9  
Living reward = 0



k=12



Noise = 0.2  
Discount = 0.9  
Living reward = 0

k=100



Noise = 0.2  
Discount = 0.9  
Living reward = 0

# Policy Iteration

Alternative approach for optimal values:

- **Step 1: Policy evaluation**: calculate utilities for some fixed policy (may not be optimal!) until convergence
- **Step 2: Policy improvement**: update policy using one-step look-ahead with resulting converged (but not optimal!) utilities as future values
- **Repeat** steps until policy converges

This is **policy iteration**

- It's still optimal!
- Can converge (much) faster under some conditions

# Policy Iteration

Policy Evaluation: For fixed current policy  $\pi$ , find values w.r.t. the policy

- Iterate until values converge:

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

Policy Improvement: For fixed values, get a better policy with one-step look-ahead:

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

Similar to how you derive optimal policy  $\pi^*$  given optimal value  $V^*$

# Piazza Poll 4

True/False:  $V^{\pi_{i+1}}(s) \geq V^{\pi_i}(s), \forall s$

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

# Piazza Poll 4

True/False:  $V^{\pi_{i+1}}(s) \geq V^{\pi_i}(s), \forall s$

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

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

If I take first step according to  $\pi_{i+1}$  and then follow  $\pi_i$ , we get an expected utility of

$$\tilde{V}_1(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^{\pi_i}(s')]$$

Which is  $\geq V^{\pi_i}(s)$

What if I take two steps according to  $\pi_{i+1}$ ?

# 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

In policy iteration:

- We do several passes that update utilities with fixed policy (each pass is fast because we consider only one action, not all of them)
- After the policy is evaluated, a new policy is chosen (slow like a value iteration pass)
- The new policy will be better (or we're done)

(Both are **dynamic programs** for solving MDPs)

# Summary: MDP Algorithms

So you want to....

- Turn **values** into a **policy**: use one-step lookahead
- Compute optimal **values**: use **value iteration** or **policy iteration**
- Compute **values** for a particular **policy**: use **policy evaluation**

These all look the same!

- They basically are – they are all variations of Bellman updates
- They all use one-step lookahead expectimax fragments
- They differ only in whether we plug in a fixed policy or max over actions



# MDP Notation

Standard expectimax:  $V(s) = \max_a \sum_{s'} P(s'|s, a) V(s')$

Bellman equations:  $V(s) = \max_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V(s')]$

Value iteration:  $V_{k+1}(s) = \max_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V_k(s')], \quad \forall s$

Q-iteration:  $Q_{k+1}(s, a) = \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma \max_{a'} Q_k(s', a')], \quad \forall s, a$

Policy extraction:  $\pi_V(s) = \operatorname{argmax}_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V(s')], \quad \forall s$

Policy evaluation:  $V_{k+1}^\pi(s) = \sum_{s'} P(s'|s, \pi(s)) [R(s, \pi(s), s') + \gamma V_k^\pi(s')], \quad \forall s$

Policy improvement:  $\pi_{new}(s) = \operatorname{argmax}_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V^{\pi_{old}}(s')], \quad \forall s$

# MDP Notation

Standard expectimax:  $V(s) = \max_a \sum_{s'} P(s'|s, a) V(s')$

Bellman equations:  $V(s) = \max_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V(s')]$

Value iteration:  $V_{k+1}(s) = \max_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V_k(s')], \quad \forall s$

Q-iteration:  $Q_{k+1}(s, a) = \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma \max_{a'} Q_k(s', a')], \quad \forall s, a$

Policy extraction:  $\pi_V(s) = \operatorname{argmax}_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V(s')], \quad \forall s$

Policy evaluation:  $V_{k+1}^\pi(s) = \sum_{s'} P(s'|s, \pi(s)) [R(s, \pi(s), s') + \gamma V_k^\pi(s')], \quad \forall s$

Policy improvement:  $\pi_{new}(s) = \operatorname{argmax}_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V^{\pi_{old}}(s')], \quad \forall s$

# MDP Notation

Standard expectimax:  $V(s) = \max_a \sum_{s'} P(s'|s, a) V(s')$

Bellman equations:  $V(s) = \max_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V(s')]$

Value iteration:  $V_{k+1}(s) = \max_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V_k(s')], \quad \forall s$

Q-iteration:  $Q_{k+1}(s, a) = \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma \max_{a'} Q_k(s', a')], \quad \forall s, a$

Policy extraction:  $\pi_V(s) = \operatorname{argmax}_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V(s')], \quad \forall s$

Policy evaluation:  $V_{k+1}^\pi(s) = \sum_{s'} P(s'|s, \pi(s)) [R(s, \pi(s), s') + \gamma V_k^\pi(s')], \quad \forall s$

Policy improvement:  $\pi_{new}(s) = \operatorname{argmax}_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V^{\pi_{old}}(s')], \quad \forall s$

# MDP Notation

Standard expectimax:  $V(s) = \max_a \sum_{s'} P(s'|s, a) V(s')$

Bellman equations:  $V(s) = \max_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V(s')]$

Value iteration:  $V_{k+1}(s) = \max_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V_k(s')], \quad \forall s$

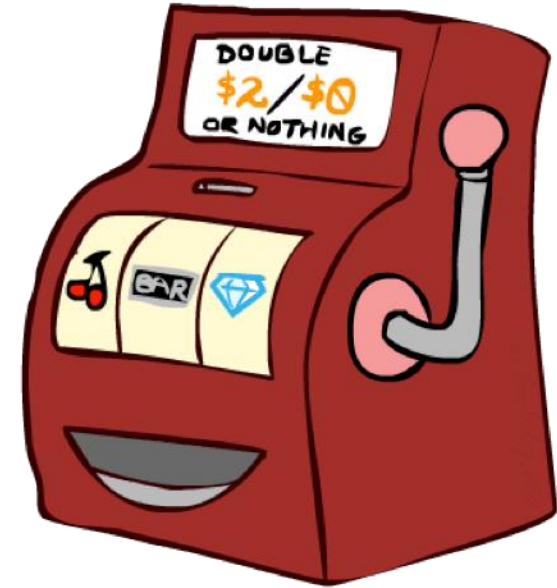
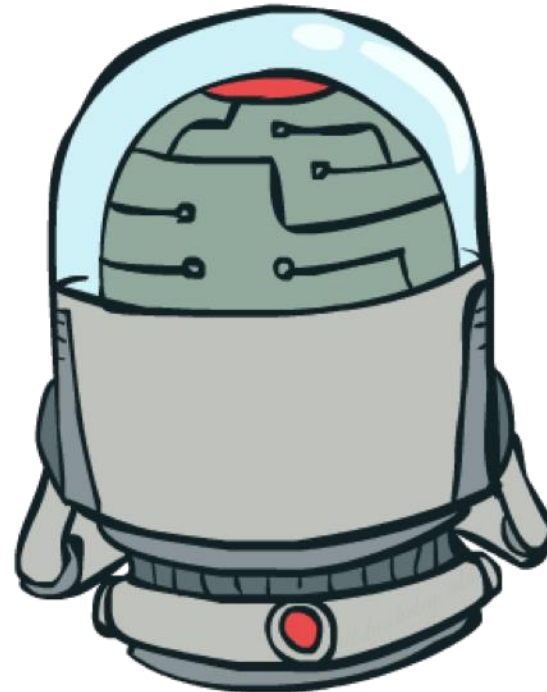
Q-iteration:  $Q_{k+1}(s, a) = \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma \max_{a'} Q_k(s', a')], \quad \forall s, a$

Policy extraction:  $\pi_V(s) = \arg\max_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V(s')], \quad \forall s$

Policy evaluation:  $V_{k+1}^\pi(s) = \sum_{s'} P(s'|s, \pi(s)) [R(s, \pi(s), s') + \gamma V_k^\pi(s')], \quad \forall s$

Policy improvement:  $\pi_{new}(s) = \arg\max_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V^{\pi_{old}}(s')], \quad \forall s$

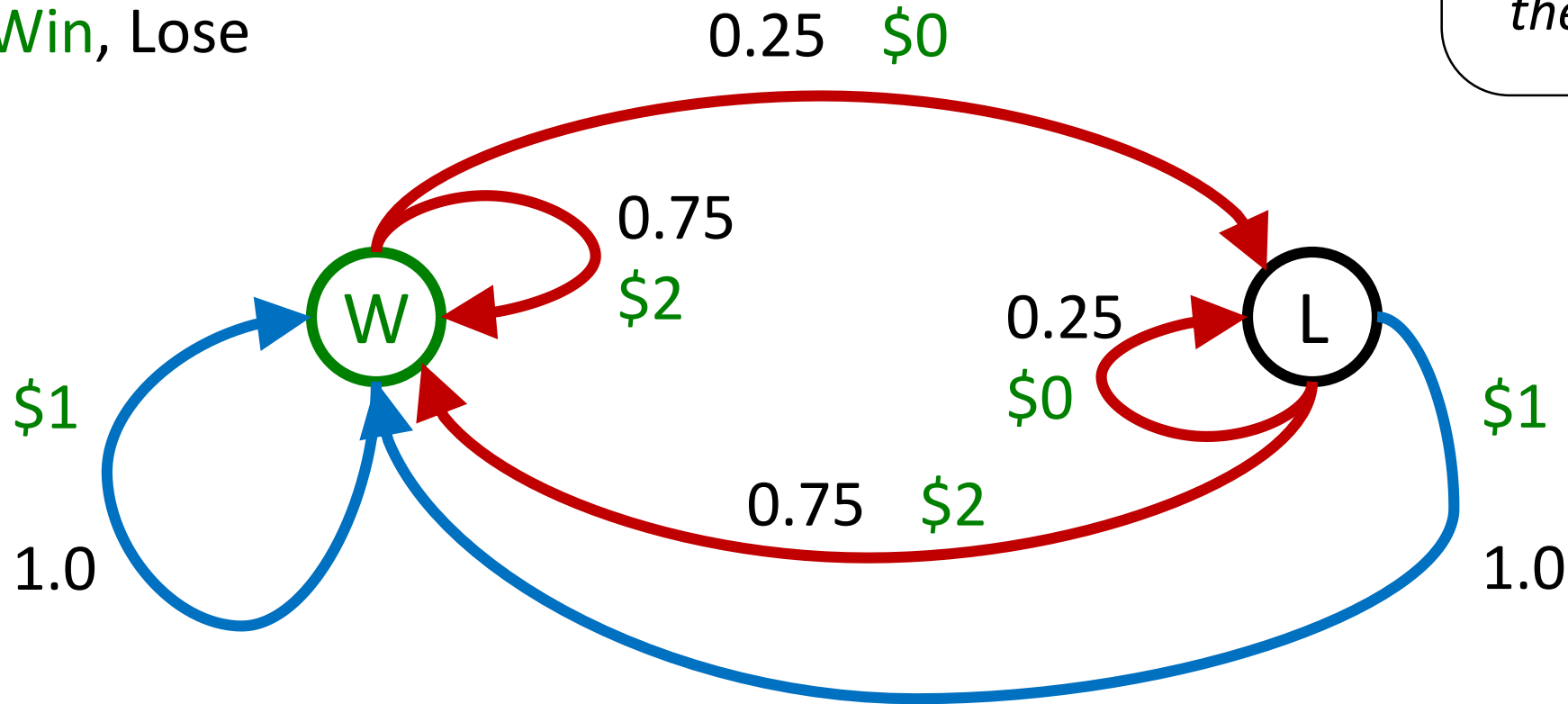
# Double Bandits



# Double-Bandit MDP

Actions: *Blue*, *Red*

States: *Win*, Lose



*No discount*  
*100 time steps*  
*Both states have the same value*

Actually a simple MDP where the current state does not impact transition or reward:

$$P(s'|s, a) = P(s'|a) \text{ and } R(s, a, s') = R(a, s')$$

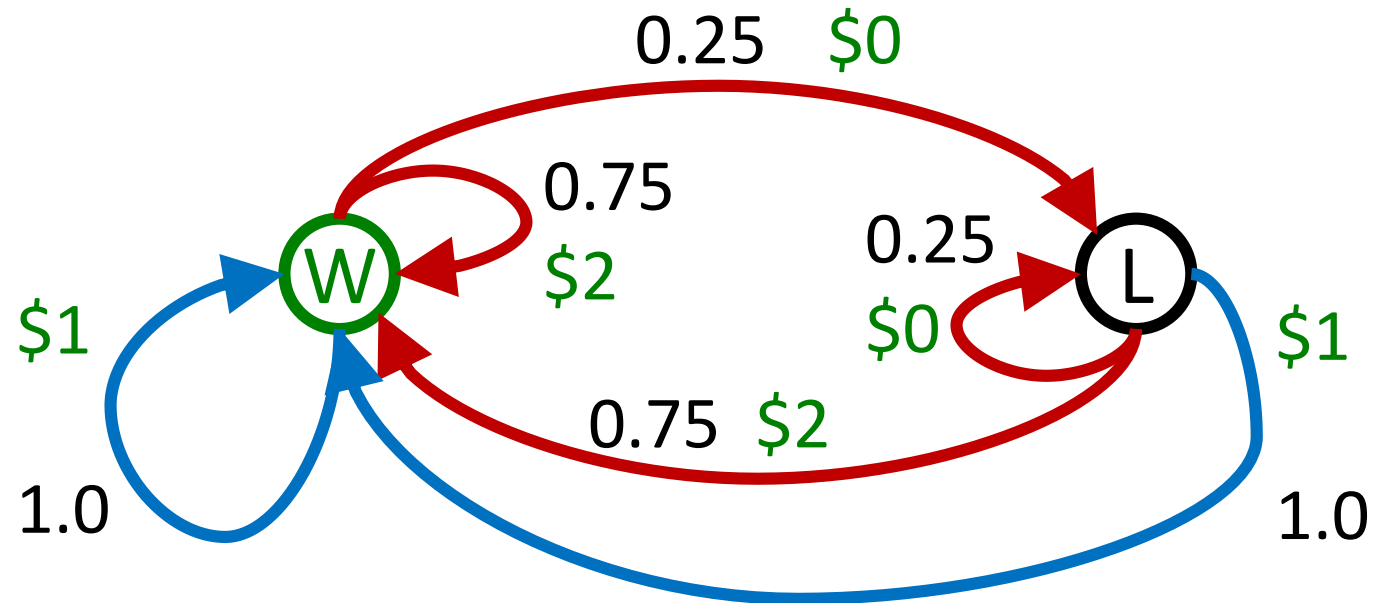
# Offline Planning

## Solving MDPs is offline planning

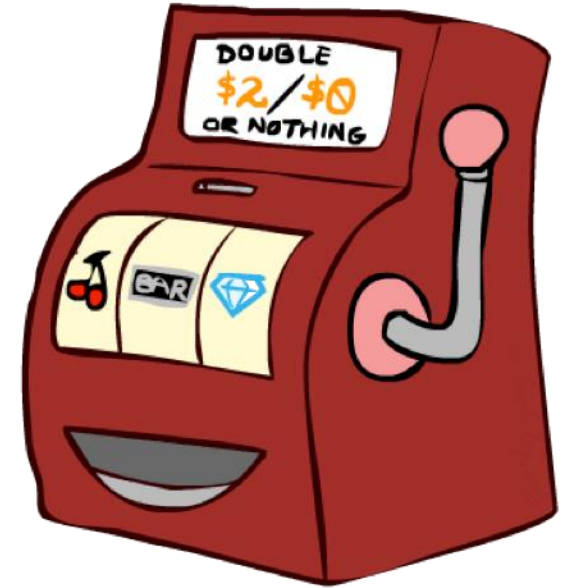
- You determine all quantities through computation
- You need to know the details of the MDP
- You do not actually play the game!

*No discount*  
*100 time steps*  
*Both states have the same value*

	Value
Play Red	150
Play Blue	100



# Let's Play!



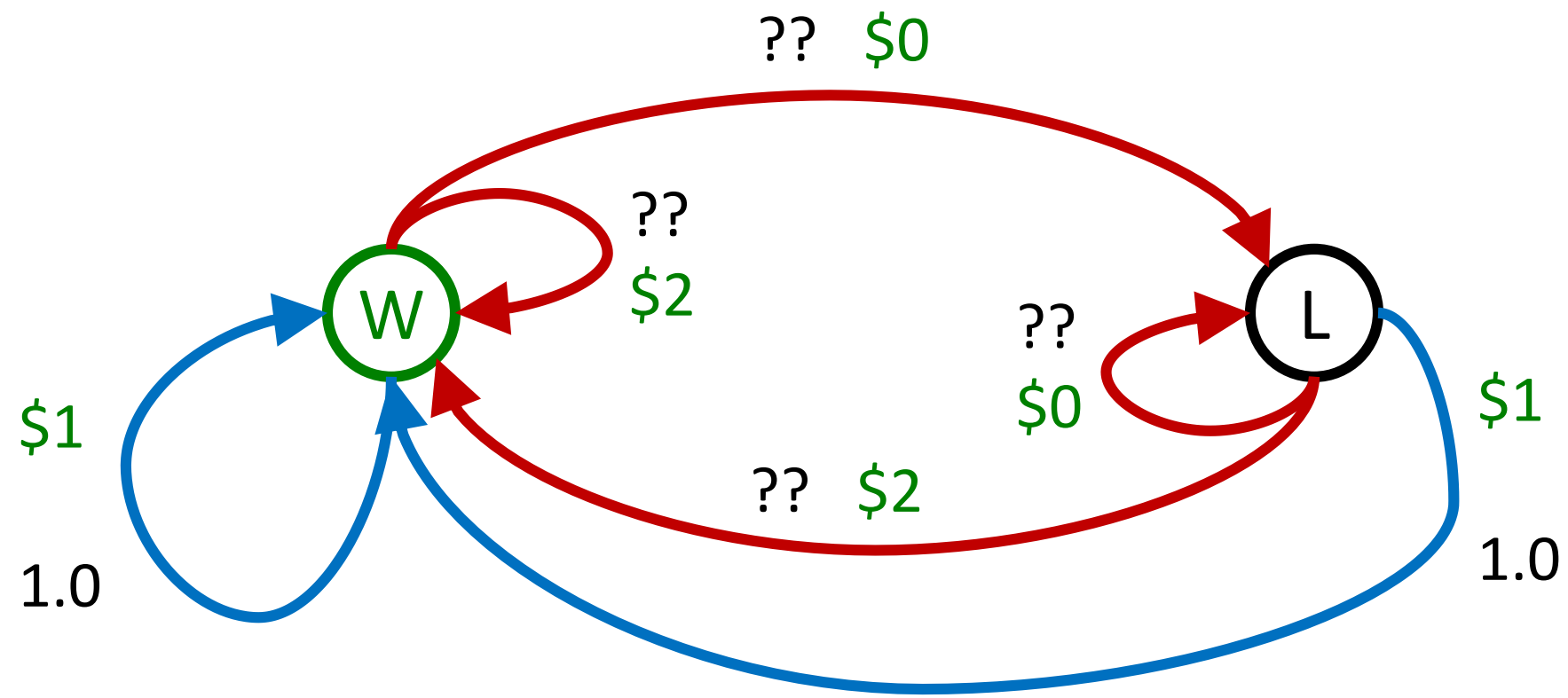
\$2 \$2 \$0 \$2 \$2

\$2 \$2 \$0 \$0 \$0

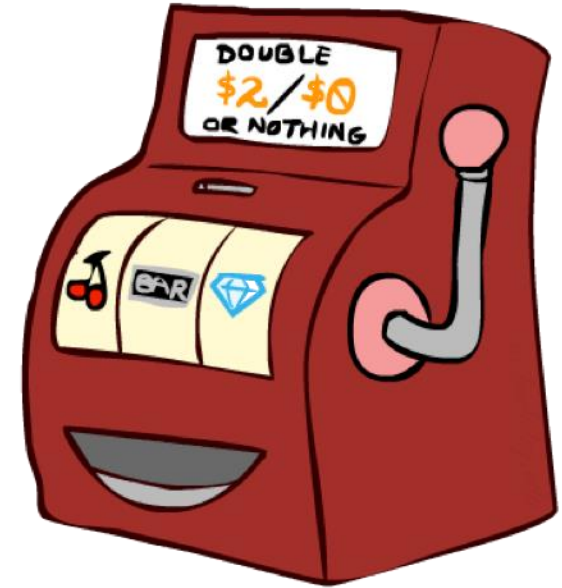


# Online Planning

Rules changed! Red's win chance is different.



# Let's Play!



\$0 \$0 \$0 \$2 \$0  
\$2 \$0 \$0 \$0 \$0

# What Just Happened?



That wasn't planning, it was learning!

- Specifically, reinforcement learning
- There was an MDP, but you couldn't solve it with just computation
- You needed to actually act to figure it out

Important ideas in reinforcement learning that came up

- **Exploration**: you have to 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**: because of chance, you have to try things repeatedly
- **Difficulty**: learning can be much harder than solving a known MDP

Next Time: Reinforcement Learning!