

CS 188: Artificial Intelligence

MDP II: Value/Policy Iteration



Instructor: Stuart Russell and Dawn Song

University of California, Berkeley

Recap: Markov Decision Process (MDP)

- What is a Markov Decision Process?



Andrey Markov
(1856-1922)

Recap: Markov Decision Process (MDP)

- What is a Markov Decision Process?
 - State transition model is markov
 - Utility function is additive discounted rewards
- An MDP is defined by:
 - A set of states $s \in S$
 - A set of actions $a \in A$
 - A transition model $T(s, a, s')$
 - Probability that a from s leads to s' , i.e., $P(s' | s, a)$
 - A reward function $R(s, a, s')$ for each transition
 - A start state
 - Possibly a terminal state (or absorbing state)
 - Utility function which is additive discounted rewards

$$U_h([s_0, a_0, s_1, a_1, s_2, \dots]) = R(s_0, a_0, s_1) + \gamma R(s_1, a_1, s_2) + \gamma^2 R(s_2, a_2, s_3) + \dots$$

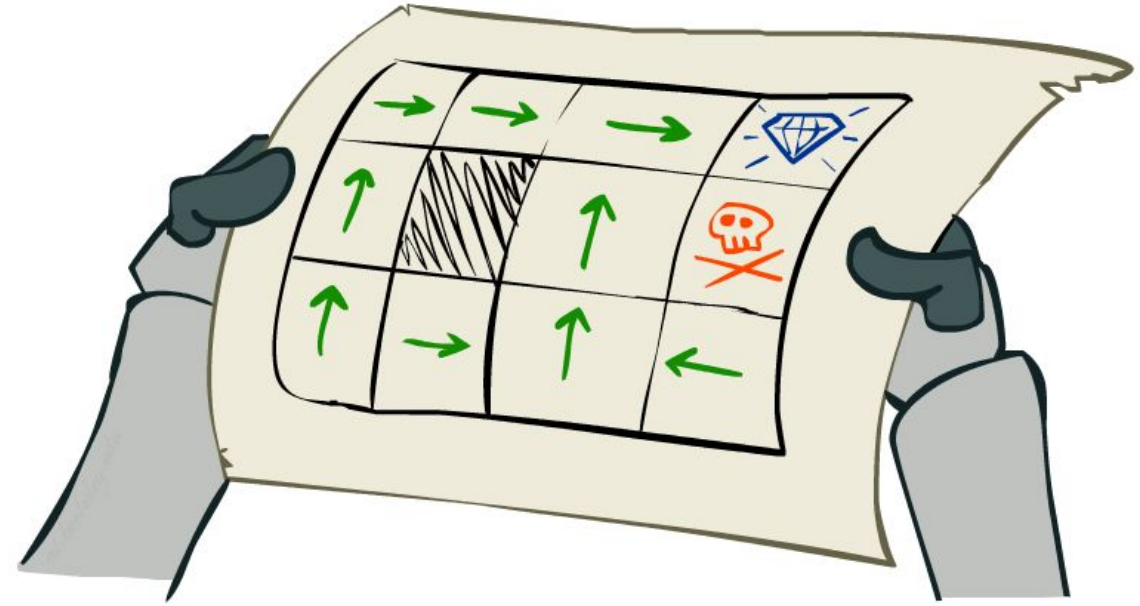
where $\gamma \in [0,1]$ is the **discount factor**



Andrey Markov
(1856-1922)

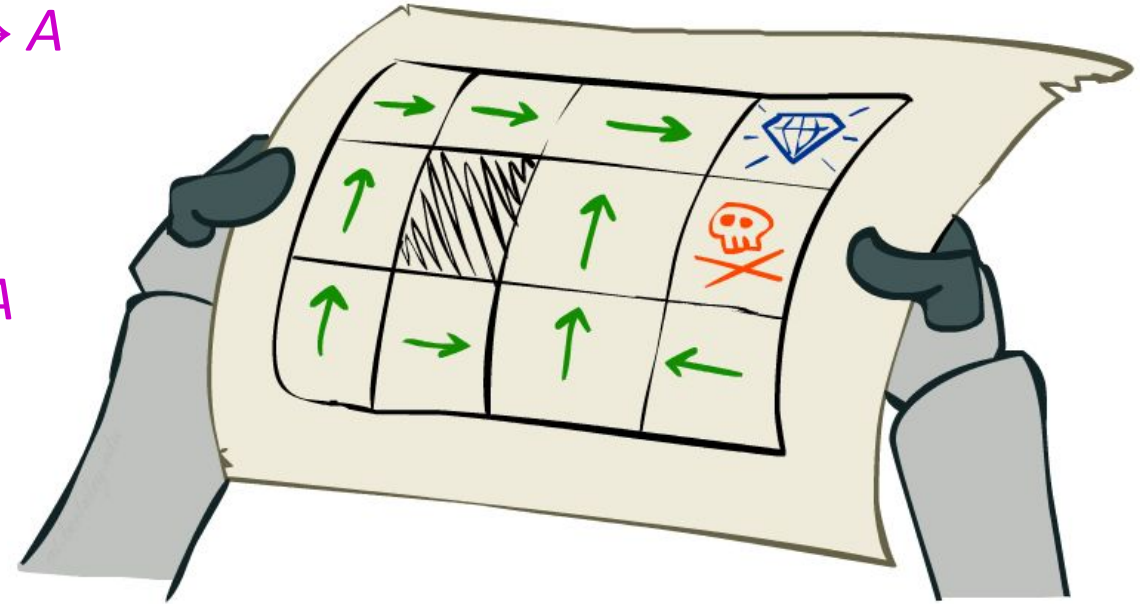
Recap: Policies

- What is a policy?



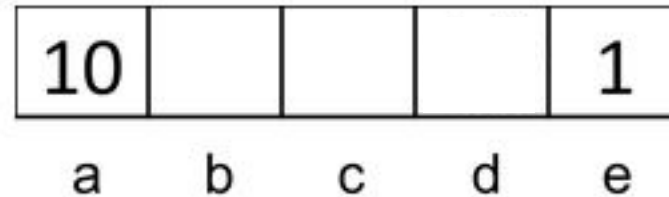
Recap: Policies

- A policy π gives an action for each state, $\pi: S \rightarrow A$
- For MDPs, we want an optimal **policy** $\pi^*: S \rightarrow A$
 - An optimal policy maximizes expected utility



Quiz: Discounting

- Given:



- Actions: East, West, and Exit (only available in exit states a, e)
- Transitions: deterministic

- Quiz 1: For $\gamma = 1$, what is the optimal policy?



- Quiz 2: For $\gamma = 0.1$, what is the optimal policy?



- Quiz 3: For which γ are West and East equally good when in state d?

The utility of a policy

- Executing a policy π from any state s_0 generates a sequence

$s_0, \pi(s_0), s_1, \pi(s_1), s_2, \dots$

- This corresponds to a sequence of rewards

$R(s_0, \pi(s_0), s_1), R(s_1, \pi(s_1), s_2), \dots$

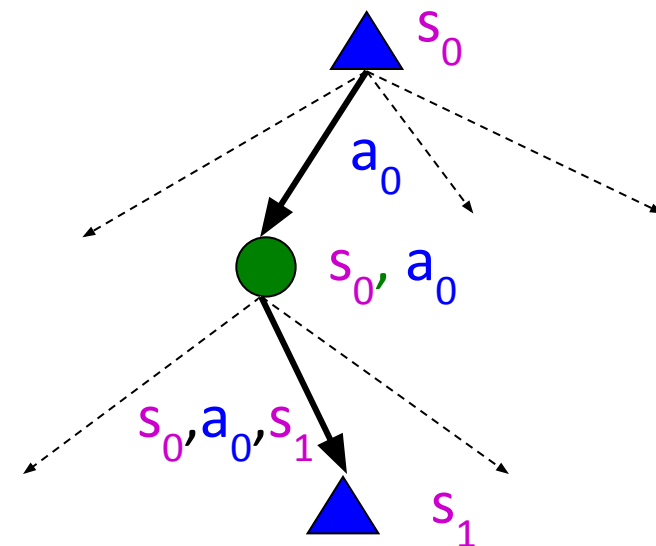
- This reward sequence happens with probability

$P(s_1 | s_0, \pi(s_0)) \times P(s_2 | s_1, \pi(s_1)) \times \dots$

- The value (expected utility) of π in s_0 is written $U^\pi(s_0)$

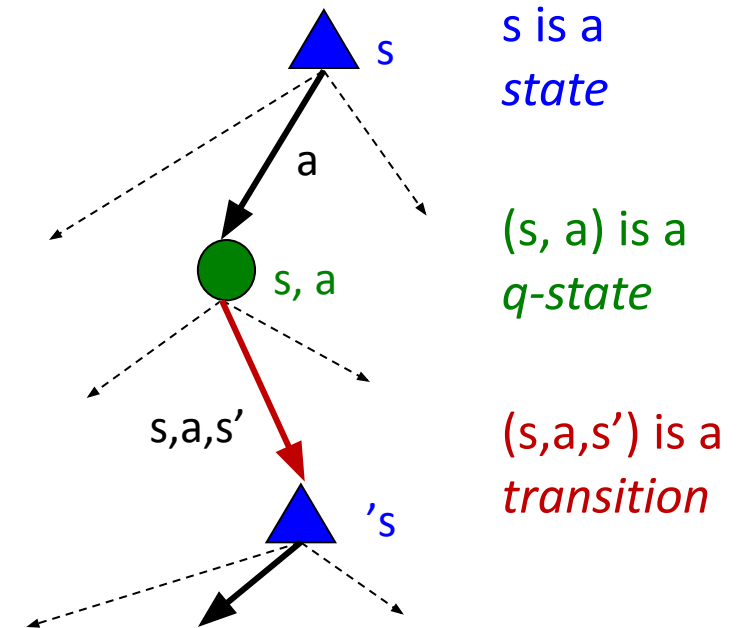
- It's the sum over all possible state sequences of
(discounted sum of rewards) \times (probability of state sequence)

$$U^\pi(s_0) = E \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(s_t), s_{t+1}) \right]$$



Optimal Quantities

- The optimal policy:
 $\pi^*(s)$ = optimal action from state s
Gives highest $U^\pi(s)$ for any π
- The value (utility) of a state s :
 $U^*(s) = U^{\pi^*}(s)$ = expected utility starting in s and acting optimally
- The value (utility) of a q-state (s,a) :
 $Q^*(s,a)$ = expected utility of taking action a in state s and (thereafter) acting optimally
 $U^*(s) \stackrel{?}{=} Q^*(s,a)$



Optimal Quantities

- The optimal policy:

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

Gives highest $U^\pi(s)$ for any π

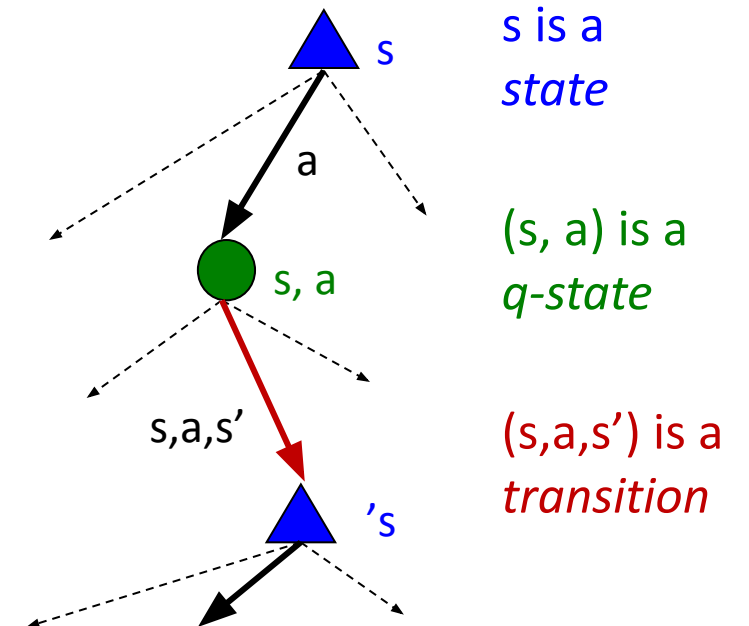
- The value (utility) of a state s :

$U^*(s) = U^{\pi^*}(s)$ = expected utility starting in s and acting optimally

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

$Q^*(s,a)$ = expected utility of taking action a in state s and (thereafter) acting optimally

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



Bellman equations (Shapley, 1953)

- The value/utility of a state is
 - The expected reward for the next transition plus the discounted value/utility of the next state, assuming the agent chooses the optimal action
- Hence we have a recursive definition of value (Bellman equation):

$$U(s) = \max_{a \in A(s)} \sum_{s'} P(s' | s, a) [R(s, a, s') + \gamma U(s')]$$

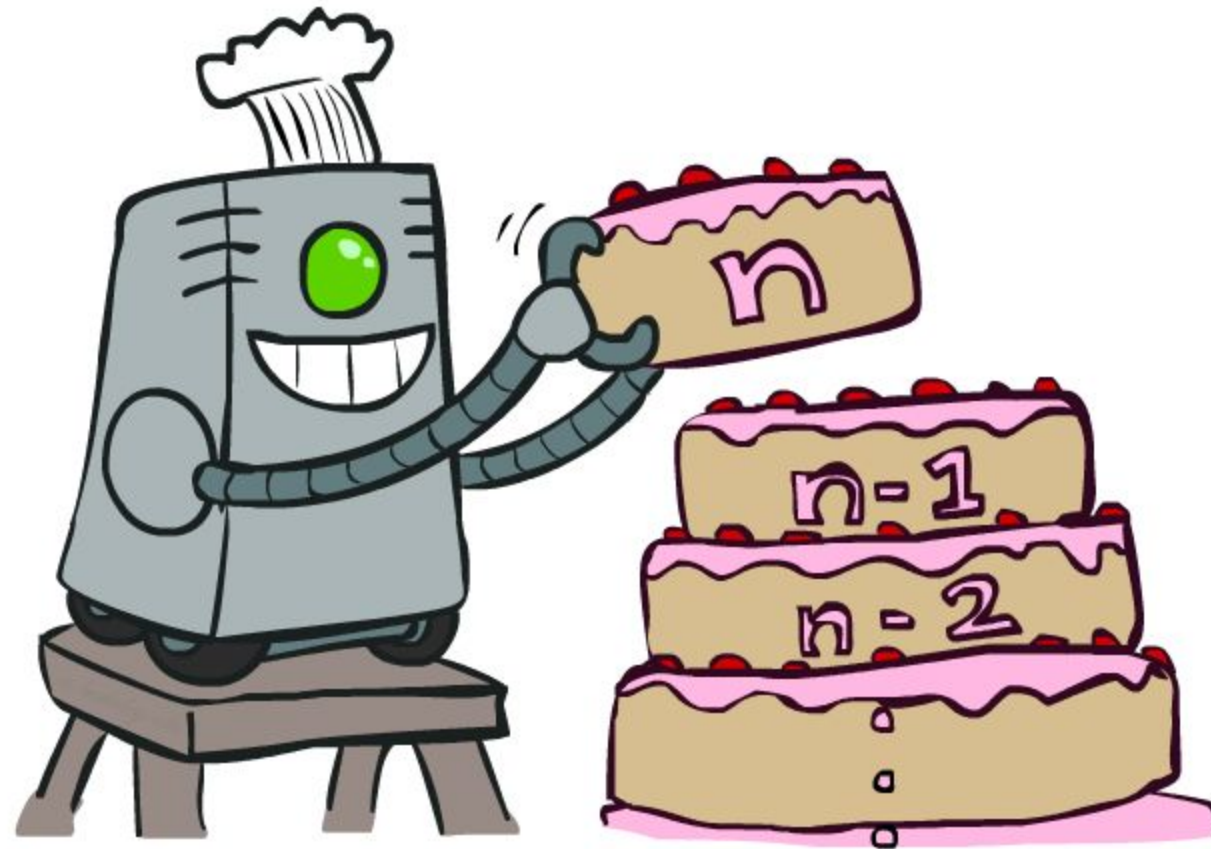
- Similarly, Bellman equation for Q-functions

$$\begin{aligned} Q(s, a) &= \sum_{s'} P(s' | s, a) [R(s, a, s') + \gamma U(s')] \\ &= \sum_{s'} P(s' | s, a) [R(s, a, s') + \gamma \max_{a'} Q(s', a')] \end{aligned}$$

Solving MDPs

- Value iteration
- Policy iteration

Value Iteration



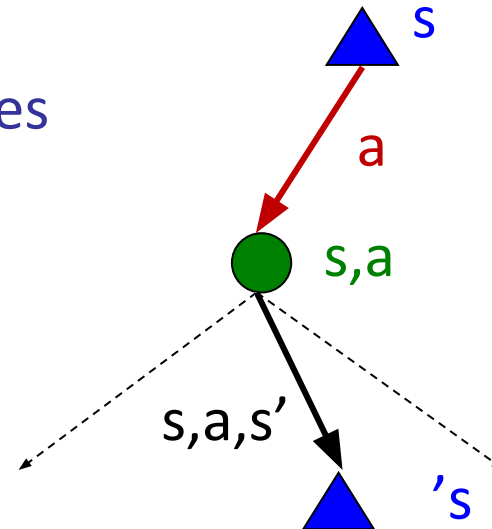
Value Iteration

- Start with (say) $U_0(s) = 0$ and some termination parameter ϵ
- Repeat until convergence (i.e., until all updates smaller than ϵ)
 - Do a **Bellman update** (essentially one ply of expectimax) from each state:

$$U_{k+1}(s) \leftarrow \max_a \sum_{s'} P(s' | a, s) [R(s, a, s') + \gamma U_k(s')]$$

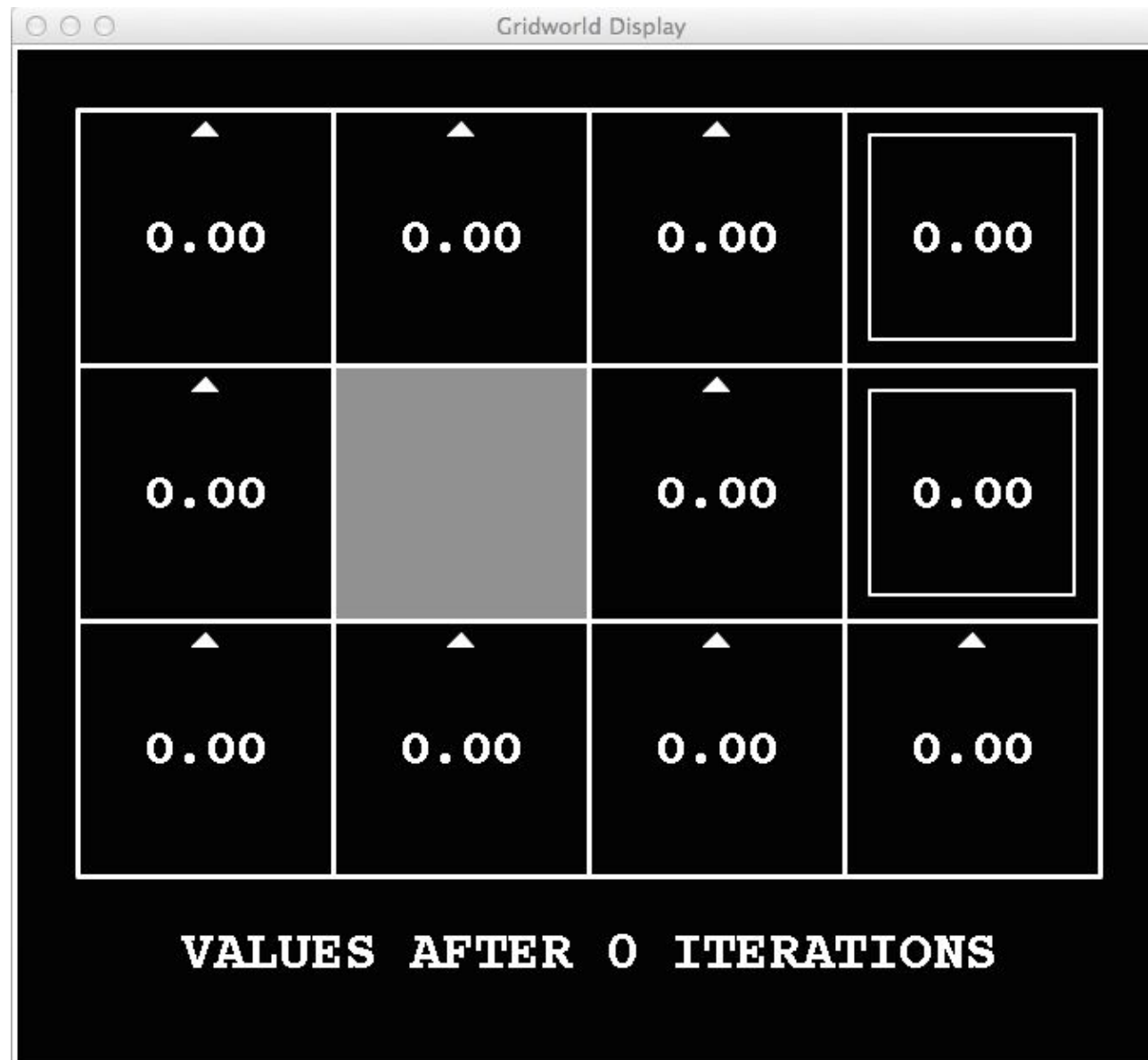
$$U \leftarrow BU$$

- Theorem: will converge to unique optimal values
- Runtime per iteration?
 - Complexity of each iteration: $O(S^2A)$



k=0

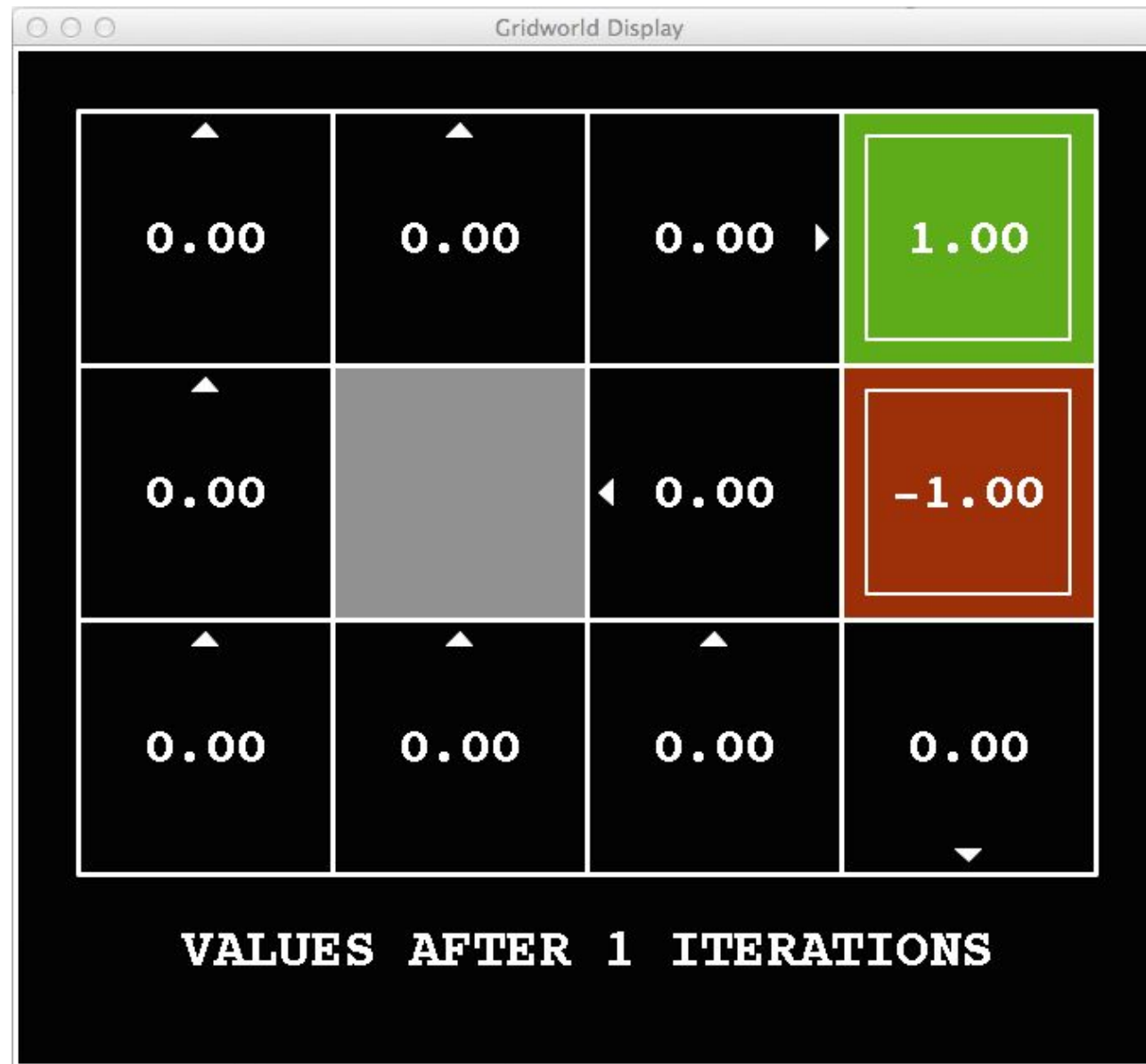
$$U_{i+1}(s) \leftarrow \max_{a \in A(s)} \sum_{s'} P(s' | s, a) [R(s, a, s') + \gamma U_i(s')]$$



Noise = 0.2
Discount = 0.9
Living reward = 0

k=1

$$U_{i+1}(s) \leftarrow \max_{a \in A(s)} \sum_{s'} P(s' | s, a) [R(s, a, s') + \gamma U_i(s')]$$



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

How do we know it will converge?*

- New concept: **contraction**

- If some operator F is a contraction by a factor, it brings any pair of objects **closer** to each other (according to some metric $d(,)$)

- For any x, y $d(Fx, Fy) \leq c d(x, y)$ where $c < 1$

E.g., $Fx = x/2$

- If F is a contraction it has a unique fixed point z (i.e., $Fz=z$)

- Proof: Suppose it had two fixed points z and z'
 - Then $d(Fz, Fz') = d(z, z')$ violating the contraction property

- Reminder: Value iteration is just $U_{k+1} \leftarrow BU_k$

- The Bellman update B is a contraction by γ

- Metric is the **max norm**: $\|V - W\| = \max_s |V(s) - W(s)|$
- Proof: follows from definition of B , i.e., Bellman equation