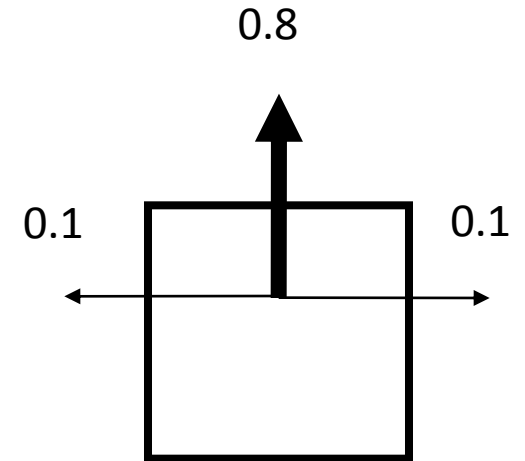


# Policy Search

<b>3</b>	-0.04	-0.04	-0.04	<b>+ 1</b>
<b>2</b>	-0.04		-0.04	<b>- 1</b>
<b>1</b>	START	-0.04	-0.04	-0.04
	<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>

3	-0.04	-0.04	-0.04	<b>+ 1</b>
2	-0.04		-0.04	<b>- 1</b>
1	START	-0.04	-0.04	-0.04
	1	2	3	4



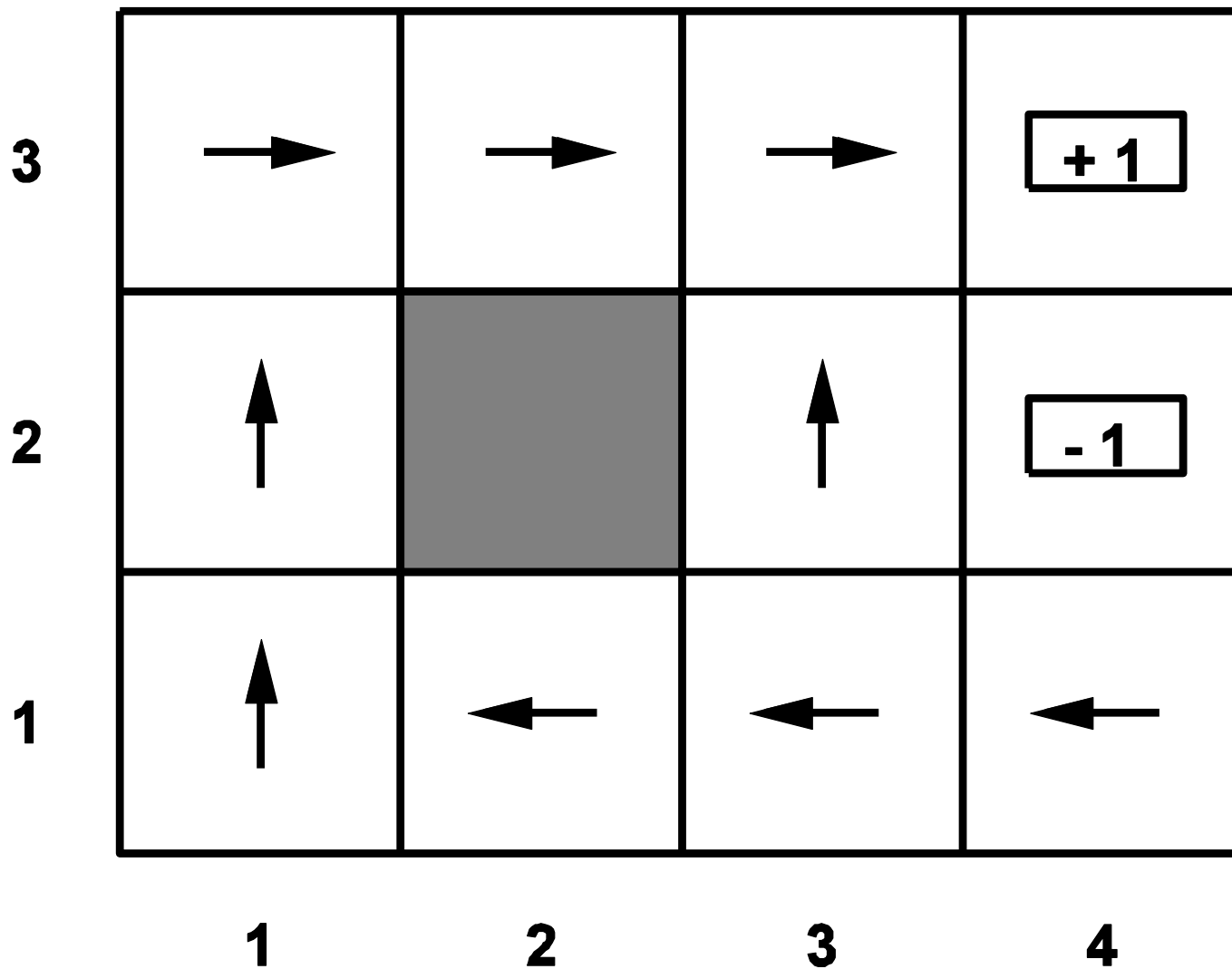
- move into desired direction with prob 80%
- move 90 degrees to left with prob 10%
- move 90 degrees to right with prob 10%

**What is the probability that [up, up, right, right, right] ends in (4,3)?**

**A=1    B=0.32768    C=0.32776    D=0.5**

We still assume the problem is fully observable – the agent knows where it is and what the rewards are...

# A Policy



# Other representations?

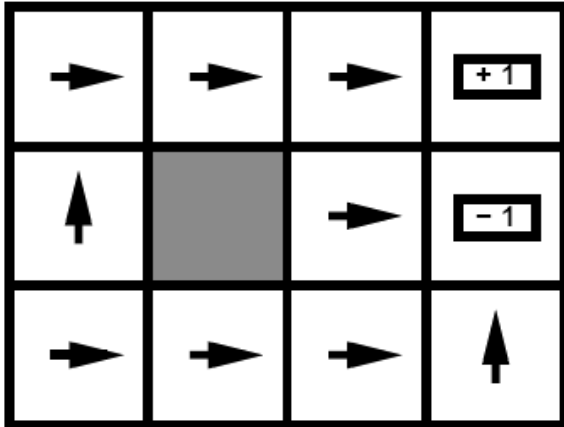
- How can a policy be represented?
  - An arrow for every grid cell
  - What about a continuous world?
    - $\Theta(x,y)$
    - NN
    - C++ function

$R(s)$

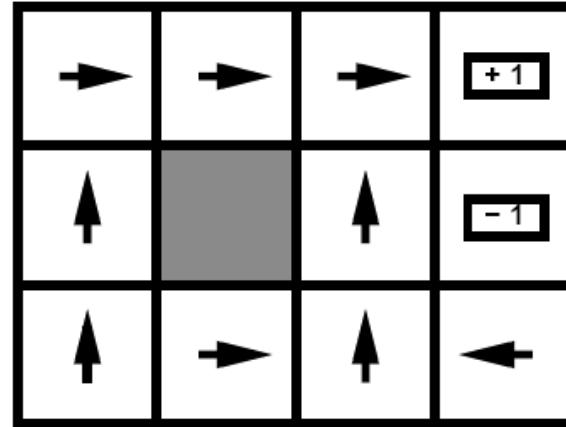
<b>3</b>	-0.04	-0.04	-0.04	<b>+ 1</b>
<b>2</b>	-0.04		-0.04	<b>- 1</b>
<b>1</b>	START	-0.04	-0.04	-0.04
	<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>

- $R(s)$  is the short term reward in state  $s$
- $U^\pi(s)$  is the long term reward when following policy  $\pi$  from state  $s$

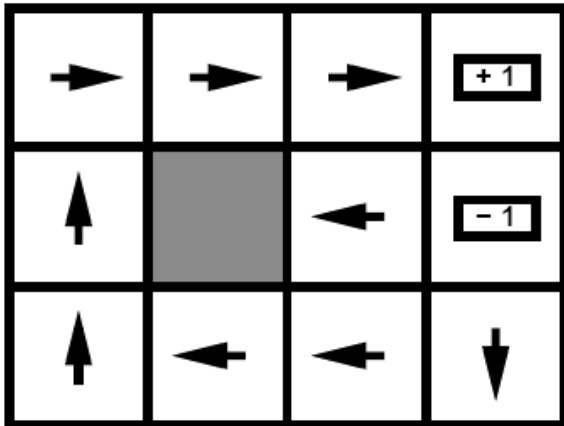
# Optimal Policies for Other Rewards



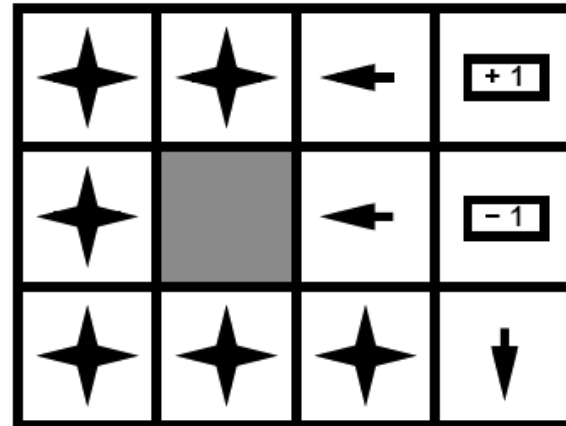
$$R(s) < -1.6284$$



$$-0.4278 < R(s) < -0.0850$$



$$-0.0221 < R(s) < 0$$



$$R(s) > 0$$

# Markovian Model

- Model:
  - Initial state:  $S_0$
  - Transition function:  $T(s,a,s')$ 
    - $T(s,a,s')$  is the probability of moving from state  $s$  to  $s'$  when executing action  $a$ .
  - Reward function:  $R(s)$ 
    - Real valued reward that the agent receives for entering state  $s$ .
- Assumptions
  - Markov property:  $T(s,a,s')$  and  $R(s)$  only depend on current state  $s$ , but not on any states visited earlier.
  - Extension: Function  $R$  may be non-deterministic as well



# Markov Decision Process

- Representation of Environment:
  - finite set of states  $S$
  - set of actions  $A$  for each state  $s$  in  $S$
- Process
  - At each discrete time step, the agent
    - observes state  $s_t$  in  $S$  and then
    - chooses action  $a_t$  in  $A$ .
  - After that, the environment
    - gives agent an immediate reward  $r_t$
    - changes state to  $s_{t+1}$  (can be probabilistic)

# Utilities

- Rating a state sequence  $[s_0, s_1, s_2, \dots]$ 
  - $U([s_0, \dots, s_N]) = \sum_i R(s_i)$
  - Additive rewards:
    - $U_h([s_0, s_1, s_2, \dots]) = R(s_0) + R(s_1) + R(s_2) + \dots$
  - We want preferences to be **stationary**
  - If  $[s_0, s_1, s_2, \dots]$  better than  $[s_0, s'_1, s'_2, \dots]$  implies  $[s_1, s_2, \dots]$  better than  $[s'_1, s'_2, \dots]$
- Reward vs. Utility

# Discounted Utility

- Problem:
  - What happens to utility value when
    - either the state space has no terminal states
    - or the policy never directs the agent to a terminal state
  - Utility becomes infinite
- Solution
  - Use discounted utility
    - closer rewards count more than awards far in the future
  - finite utility even for infinite state sequences

$$U([s_0, \dots, s_N]) = \sum_i \gamma^i R(s_i) \leq \sum_i \gamma^i R_{\max} = R_{\max} / (1 - \gamma)$$

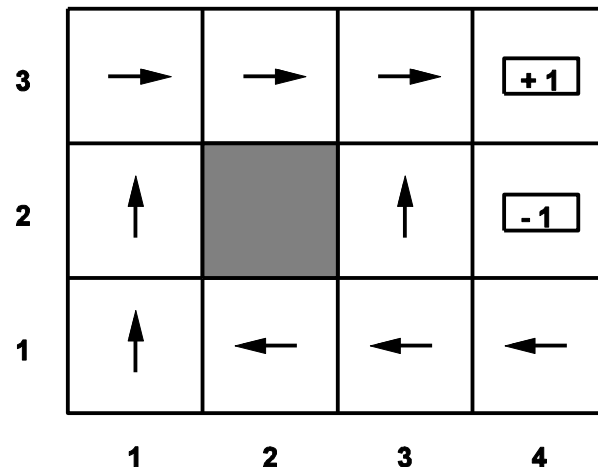
# Policy

- Definition:

- A policy  $\pi$  describes which action an agent should select in each state
- $a = \pi(s)$

- Utility of a policy

- For now additive utility
- Let  $P([s_0, \dots, s_N] \mid \pi, s_0)$  be the probability of state sequence  $[s_0, \dots, s_N]$  when following policy  $\pi$  from state  $s_0$
- Expected utility:  $U^\pi(s) = \sum U([s_0, \dots, s_N]) P([s_0, \dots, s_N] \mid \pi, s_0)$   
→ measure of quality of policy  $\pi$
- Optimal policy  $\pi^*$ : Policy with maximal  $U^\pi(s)$  in each state  $s$



# Utility $\Leftrightarrow$ Policy

- Equivalence:
  - If we know the utility  $U(s)$  of each state, we can derive the corresponding policy:

$$\pi^*(s) = \operatorname{argmax}_a \sum_{s'} T(s, a, s') U(s')$$

- If we know the policy  $\pi$ , we can compute the corresponding utility of each state:

PolicyEvaluation algorithm

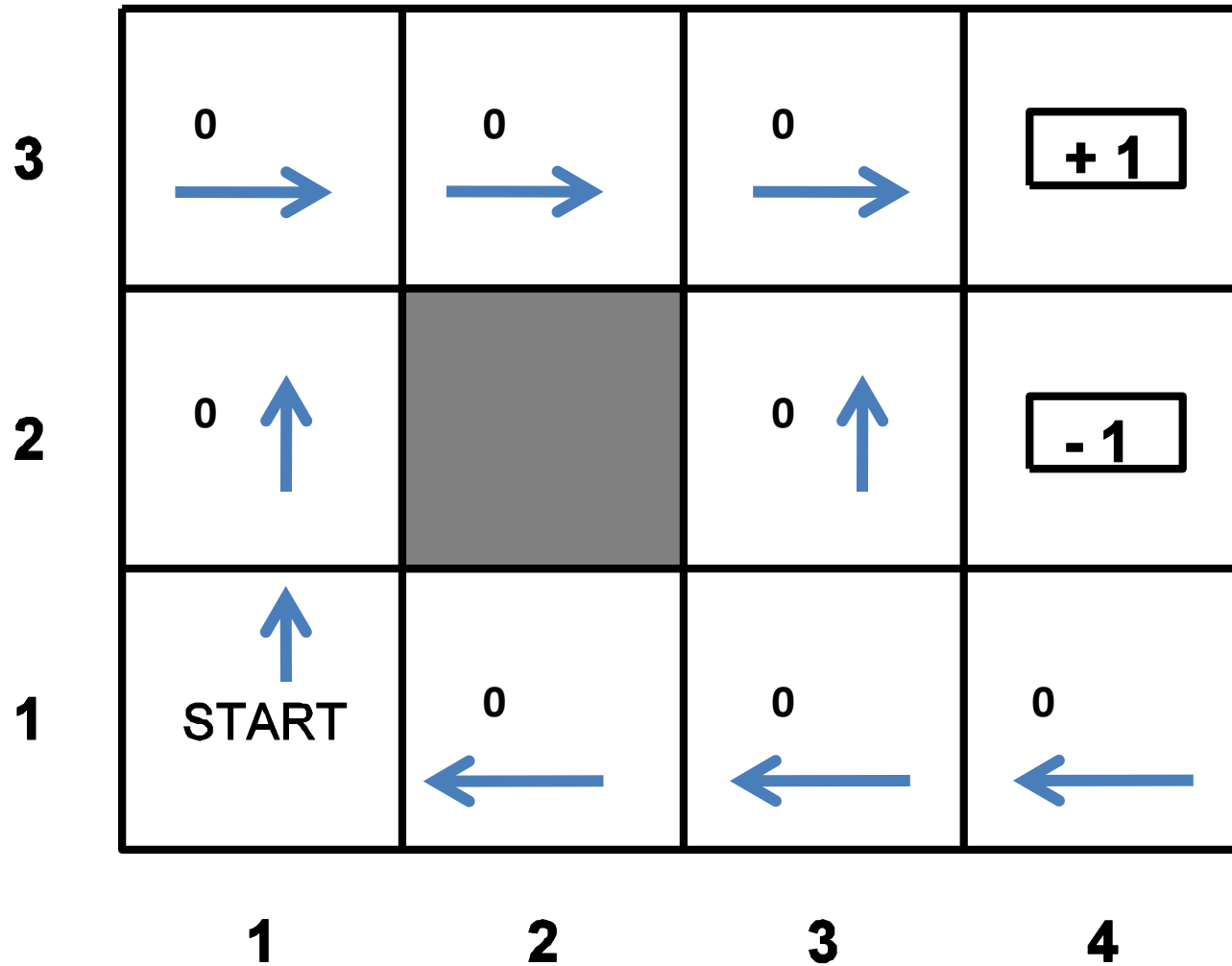
Bellman Equation:

$$U(s) = R(s) + \gamma \max_a \sum_{s'} T(s, a, s') U(s')$$

Optimal Utility  $\Leftrightarrow$  Optimal Policy

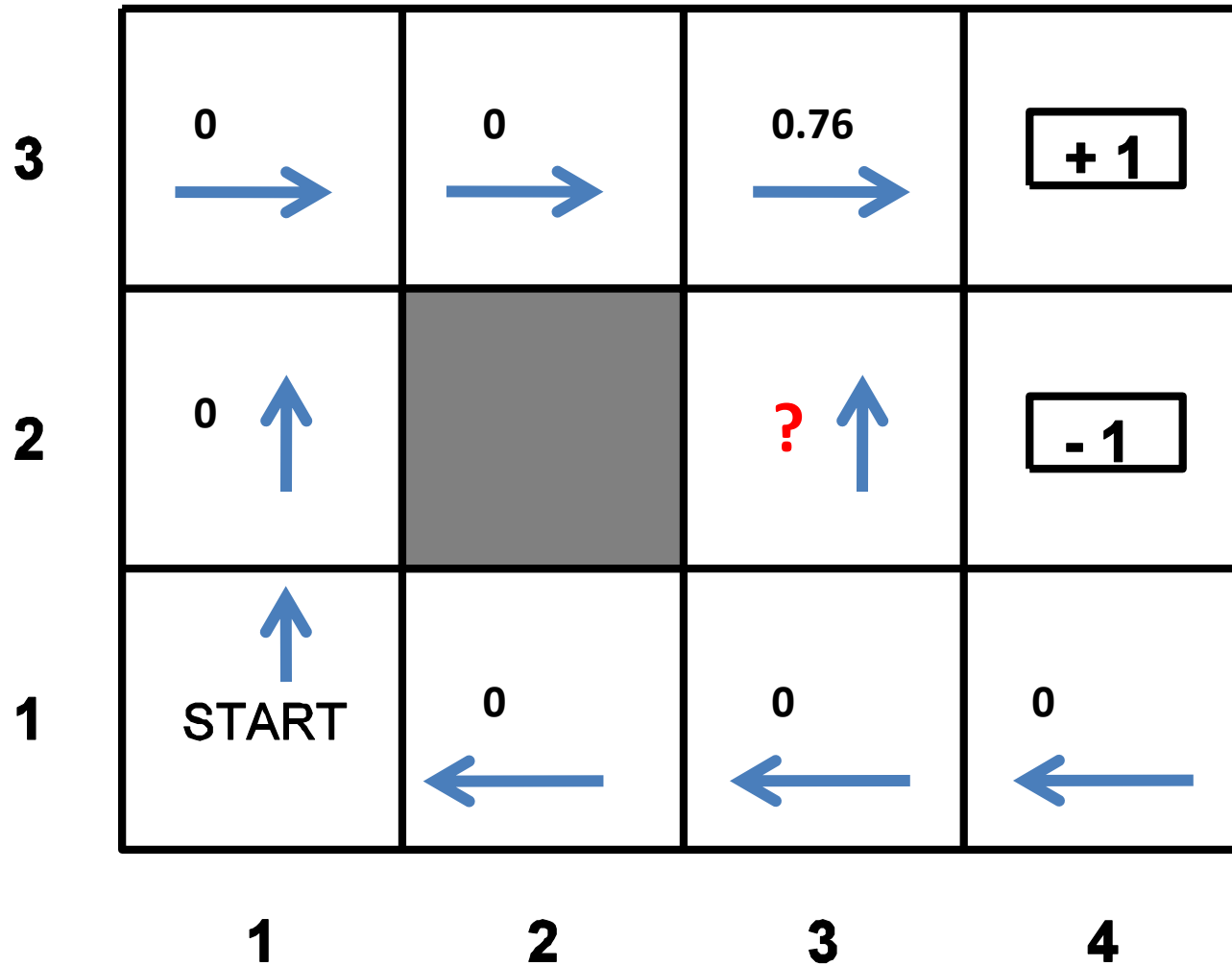
# How to Compute the Utility for a given Policy?

- Definition: Utility of path Likelihood of path
  - $U^\pi(s) = \sum [ \boxed{\sum_i \gamma^i R(s_i)} ] P( \boxed{[s_0, s_1, \dots]} \mid \pi, s_0=s)$
- Recursive computation:
  - $U^\pi(s) = R(s) + \gamma \sum_{s'} T(s, \pi(s), s') U^\pi(s')$
  - What is  $P([s_0, s_1, \dots] \mid \pi)$ ?
    - $P([s_0, s_1, \dots] \mid \pi) = T(s_0, \pi(s_0), s_1) * T(s_1, \pi(s_1), s_2) * \dots$



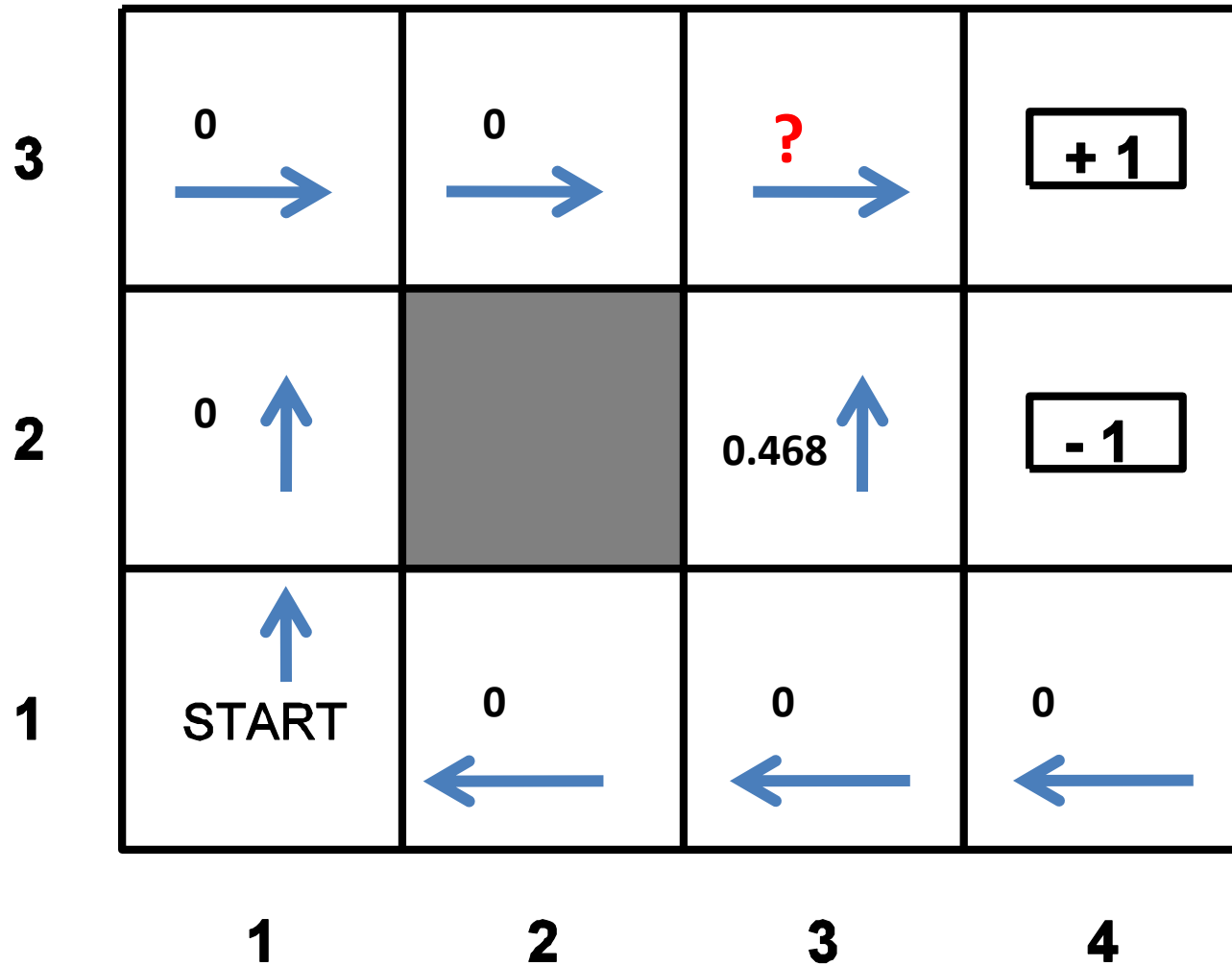
•  $U^\pi((3,3)) \leftarrow -0.04 + 0.8 * 1 + 0.1 * 0 + 0.1 * 0 = 0.76$

$R(s) = -0.04$  everywhere except goals



•  $U^\pi((3,2)) = ?$     A= -0.04    B=0.312    C=0.428    D=0.468

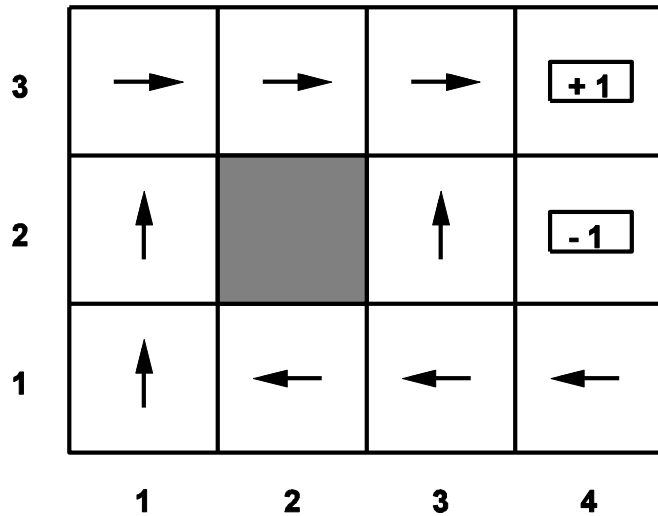




•  $U^\pi((3,3)) = ?$     A= -0.04    B=0.760    C=0.883    D=0.885

# Convergence

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



3	0.812	0.868	0.918	+1
2	0.762		0.660	-1
1	0.705	0.655	0.611	0.388
	1	2	3	4

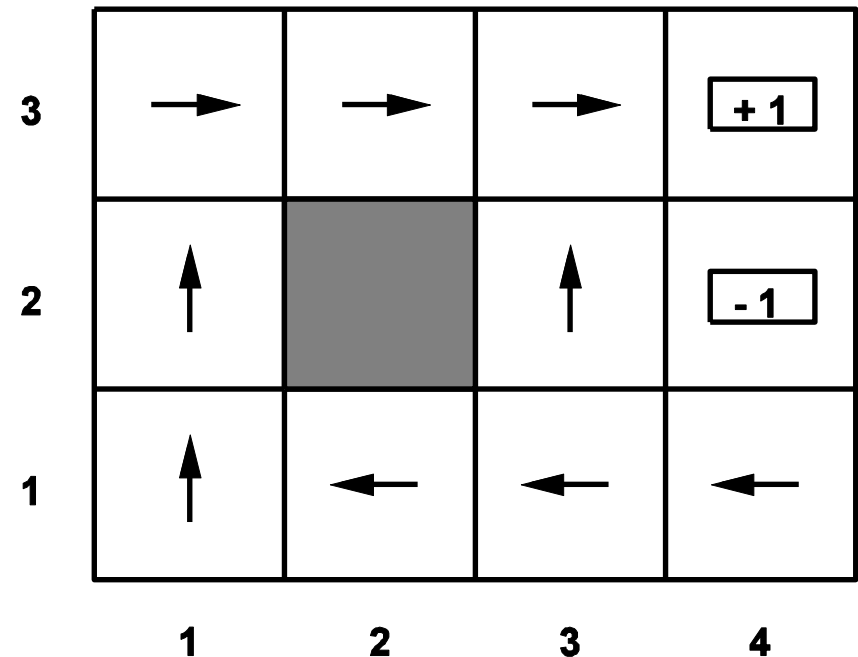
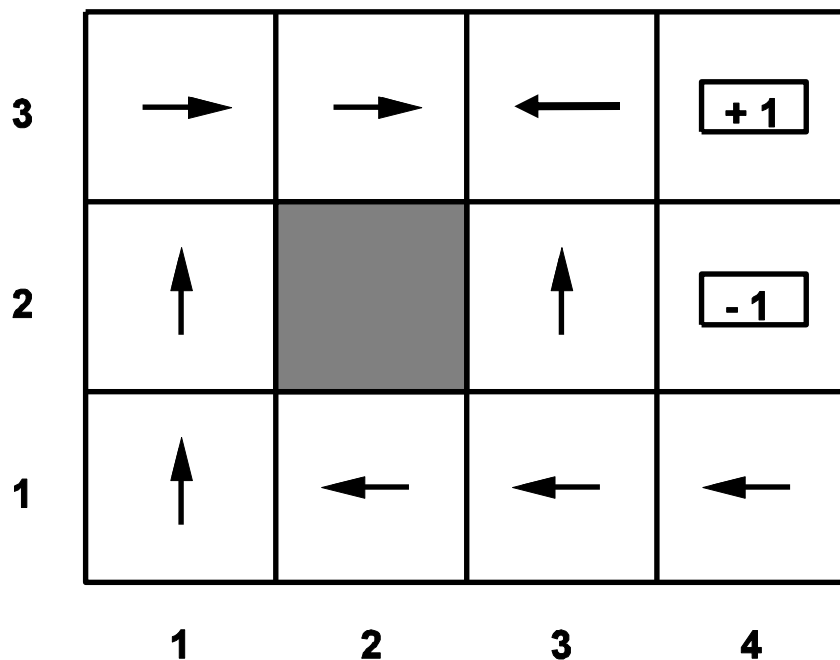
Here:  $\gamma=1.0$ ,  $R(s)=-0.04$

# Bellman Update

- Goal: Solve set of  $n=|S|$  equations (one for each state)
$$U^\pi(s_0) = R(s_0) + \gamma \sum_{s'} T(s_0, \pi(s), s') U^\pi(s')$$
$$\dots$$
$$U^\pi(s_n) = R(s_n) + \gamma \sum_{s'} T(s_n, \pi(s), s') U^\pi(s')$$
- Is this a set of linear equations? Why or why not?
- Algorithm [Policy Evaluation] (fix  $\pi$ ):
  - $i=0$ ;  $U^\pi_0(s)=0$  for all  $s$
  - repeat
    - $i = i + 1$
    - for each state  $s$  in  $S$  do
      - $U^\pi_i(s) = R(s) + \gamma \sum_{s'} T(s, \pi(s), s') U^\pi_{i-1}(s')$
    - endfor
  - until difference between  $U^\pi_i$  and  $U^\pi_{i-1}$  small enough
  - return  $U^\pi_i$

# How to Find the Optimal Policy $\pi^*$ ?

- Is policy  $\pi$  optimal? How can we tell?
  - If  $\pi$  is not optimal, then there exists some state where
$$\pi(s) \neq \operatorname{argmax}_a \sum_{s'} T(s, a, s') U^\pi(s')$$
  - How to find the optimal policy  $\pi^*$ ?



# How to Find the Optimal Policy $\pi^*$ ?

- Algorithm [Policy Iteration]:
  - repeat
    - $U^\pi = \text{PolicyEvaluation}(\pi, S, T, R)$
    - for each state  $s$  in  $S$  do
      - If [  $\max_a \sum_{s'} T(s, a, s') U^\pi(s') > \sum_{s'} T(s, \pi(s), s') U^\pi(s')$  ] then
        - »  $\pi(s) = \text{argmax}_a \sum_{s'} T(s, a, s') U^\pi(s')$
    - **endfor**
  - until  $\pi$  does not change any more
  - return  $\pi$

# Value Iteration

- Optimal Utility  $\Leftrightarrow$  Optimal Policy
- Find optimal utility directly

1	-2	5
7	?	6
3	6	5



**A**



**B**



**C**



**D**

**What is the best utility possible for the center?**

# Value Iteration

- Optimal Utility  $\Leftrightarrow$  Optimal Policy
- Find optimal utility directly

1	-2	5
7	↑	6
3	6	5

$$-0.04 + 0.8 \cdot -2 + 0.1 \cdot 7 + 0.1 \cdot 6 = -0.34$$

1	-2	5
7	←	6
3	6	5

$$-0.04 + 0.8 \cdot 7 + 0.1 \cdot -2 + 0.1 \cdot 6 = 5.96$$

1	-2	5
7	↓	6
3	6	5

$$-0.04 + 0.8 \cdot 6 + 0.1 \cdot 7 + 0.1 \cdot 6 = 6.06$$

1	-2	5
7	→	6
3	6	5

$$-0.04 + 0.8 \cdot 6 + 0.1 \cdot 6 + 0.1 \cdot -2 = 5.16$$

# Value Iteration Algorithm

- Algorithm [Value Iteration]:
  - $i=0$ ;  $U_0(s)=0$  for all  $s$
  - repeat
    - $i = i + 1$
    - for each state  $s$  in  $S$  do
      - $U_i(s) = R(s) + \gamma \max_a \sum_{s'} T(s, a, s') U_{i-1}(s')$
    - endfor
  - until difference between  $U_i$  and  $U_{i-1}$  small enough
  - return  $U_i$

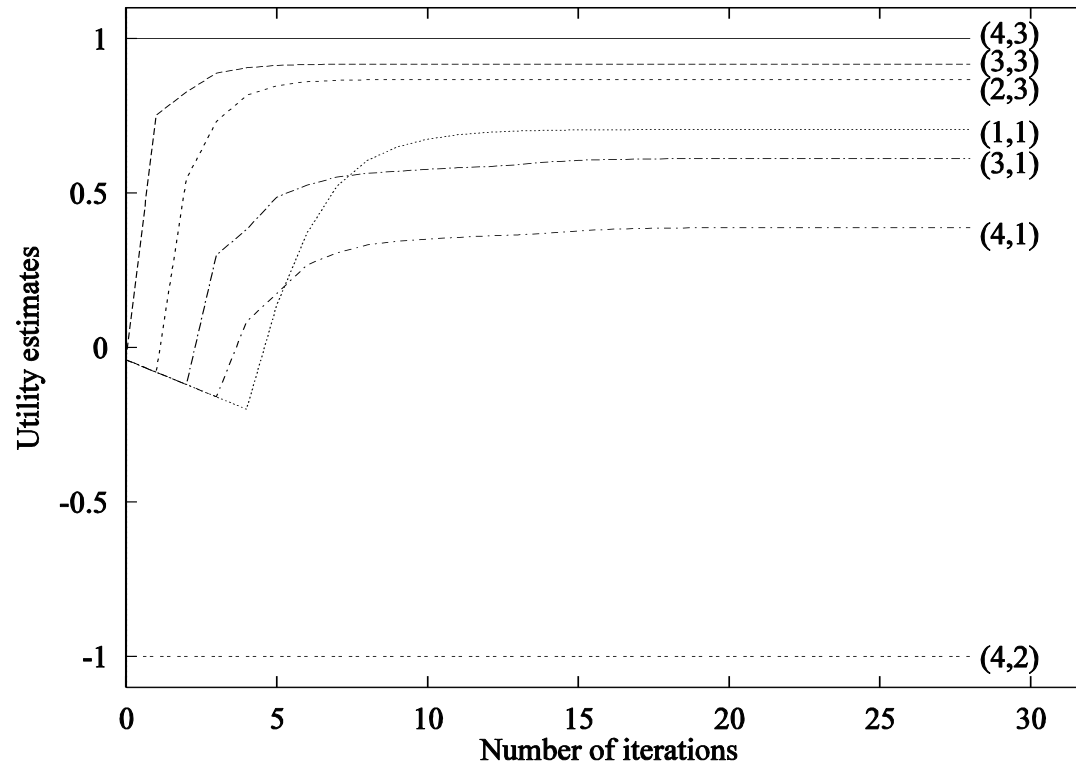
→ derive optimal policy via  $\pi^*(s) = \operatorname{argmax}_a \sum_{s'} T(s, a, s') U(s')$



# Two approaches to finding $\pi^*$

- Policy Iteration
  - Local search: Start with random policy
  - Evaluate a candidate policy using bellman update
  - $U^\pi(s_0) = R(s_0) + \gamma \sum_{s'} T(s_0, \pi(s), s') U^\pi(s')$
  - Linear set of equations
  - Confirm policy is optimal or make changes
- Value Iteration
  - find optimal utilities by iteratively solving
    - $U_i(s) = R(s) + \gamma \max_a \sum_{s'} T(s, a, s') U_{i-1}(s')$
  - Derive optimal policy to follow optimal utilities

# Convergence of Value Iteration



- Value iteration is guaranteed to converge to optimal  $U$  for  $0 \leq \gamma < 1$
- Faster convergence for smaller  $\gamma$
- Guaranteed to reach equilibrium (see textbook)