

# Fundamentals of Artificial Intelligence

## Exercise 11: Making Complex Decisions *(over time)*

Jonathan Külz

Technical University of Munich

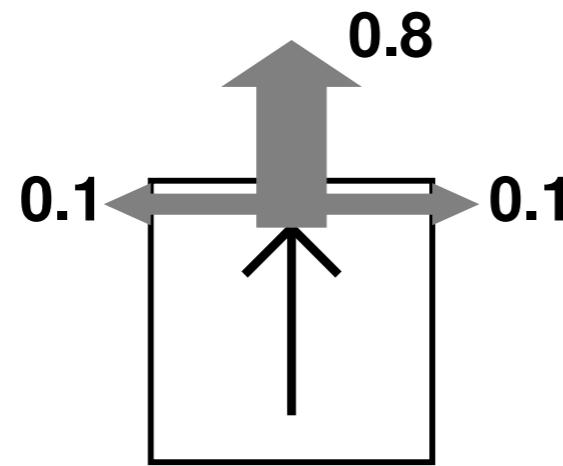
February 02nd, 2024

## Summary - Rational Decisions Over Time

$$P(s^t | a, s)$$

- Sequential decision problems in uncertain discrete environments can be modeled as **Markov decision processes (MDPs)**
- The utility of a state sequence is the sum of all the rewards over the sequence, possibly discounted over time.
- The optimal solution of an MDP is a **policy** that associates a decision with every state that the agent might reach. A solution can be obtained by **value iteration**.  
 $a$        $s \in S$
- **Policy iteration** usually converges faster, since a policy might already be optimal without knowing the exact utilities of each state.

## Problem 11.1: Roomba Problem

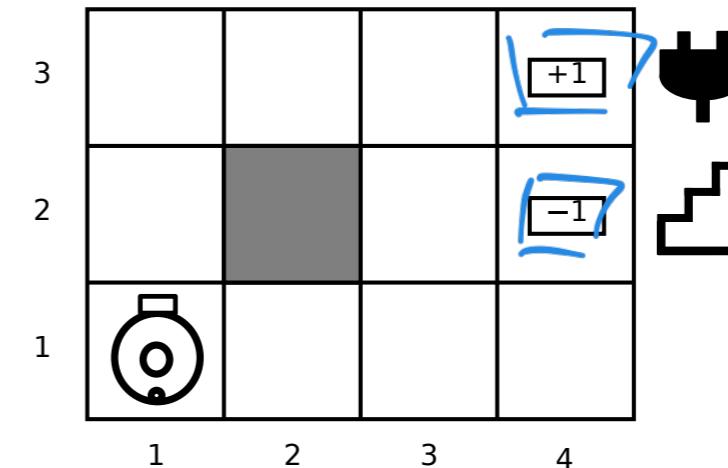


$$\gamma = 1$$

- States  $s \in S$ , actions  $a \in A = \{Up, Down, Left, Right\}$ .
- Model**  $P(s'|s, a)$  = probability that  $a$  in  $s$  leads to  $s'$ .

- Reward function** (with terminal states  $\overline{S_T} = \{s_{charge}, s_{stairs}\}$ )

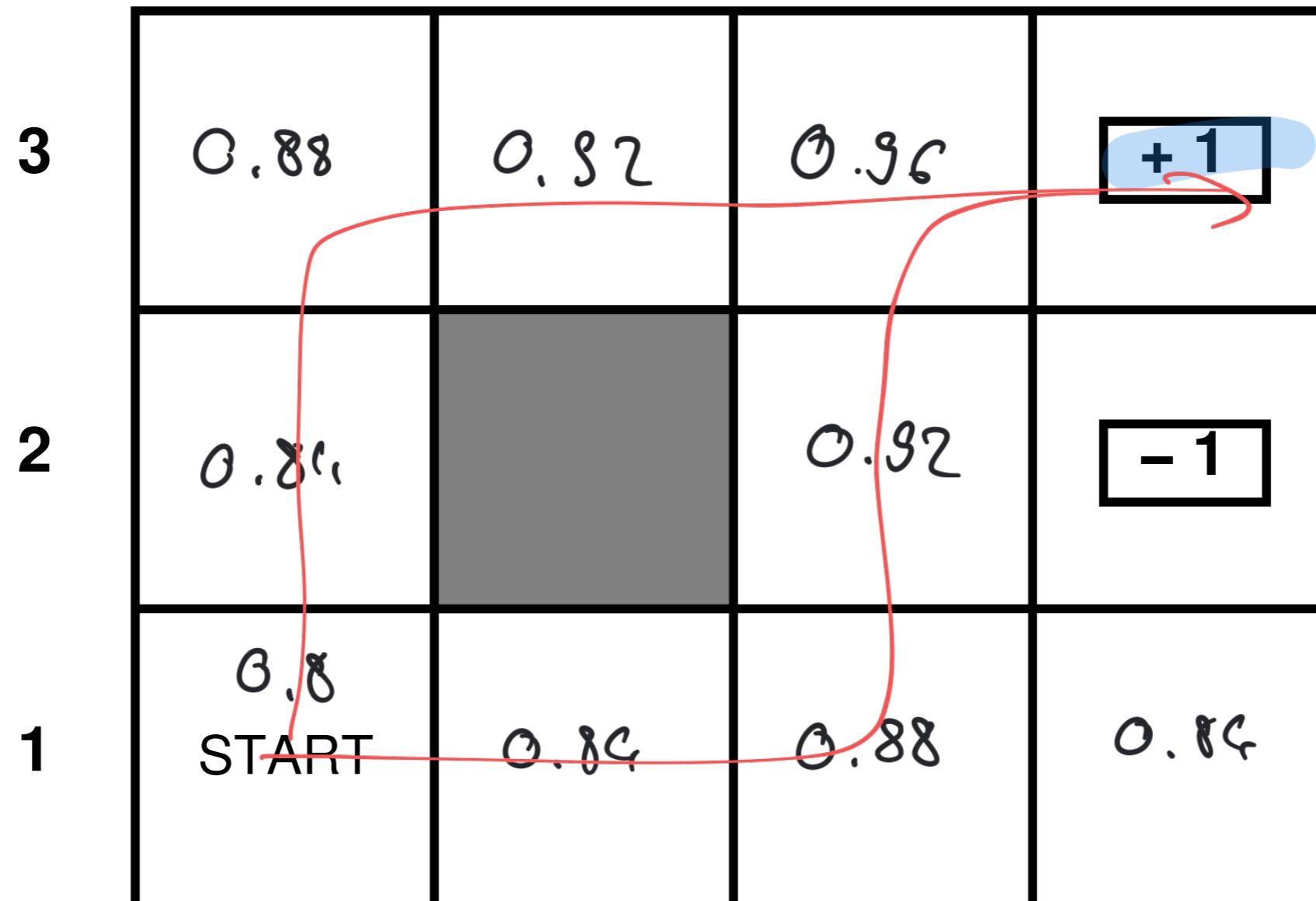
$$\overline{R(s, a, s')} = R(s) = \begin{cases} 1 & \text{if } s = \overline{s_{charge}} \\ -1 & \text{if } s = \overline{s_{stairs}} \\ -0.04 & \forall s \notin \overline{S_T} \end{cases}$$



## Problem 11.1: Roomba Problem $U(s) = R(s) + \gamma \cdot \max_a \sum P(s'|s, a) U(s')$

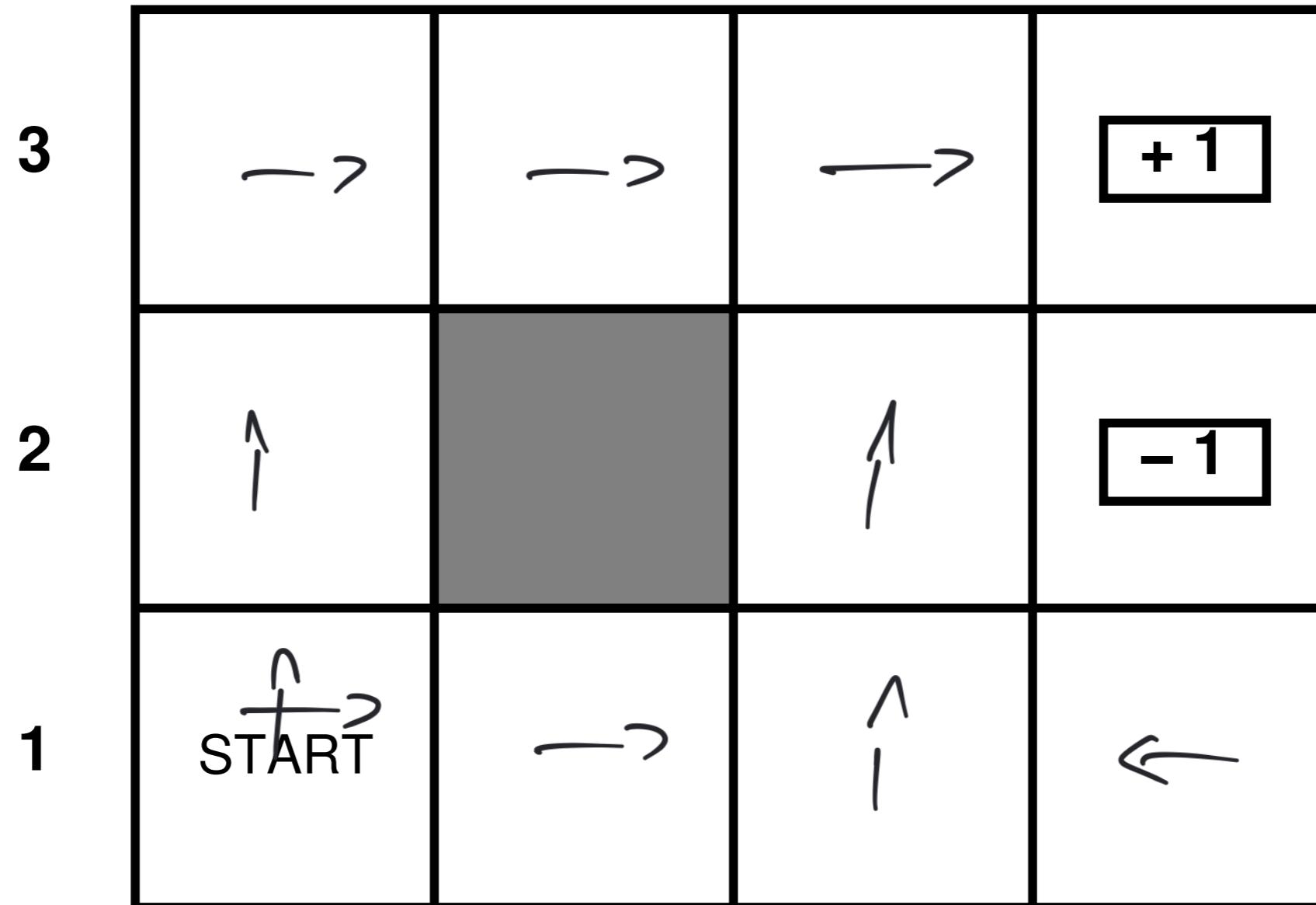
**Problem 11.1.1** Assuming the transition probability as **deterministic** and the discount factor as 1. Find the **value** of all states.

$$R(s) = \begin{cases} -0.09 \\ -1 \\ 1 \end{cases}$$



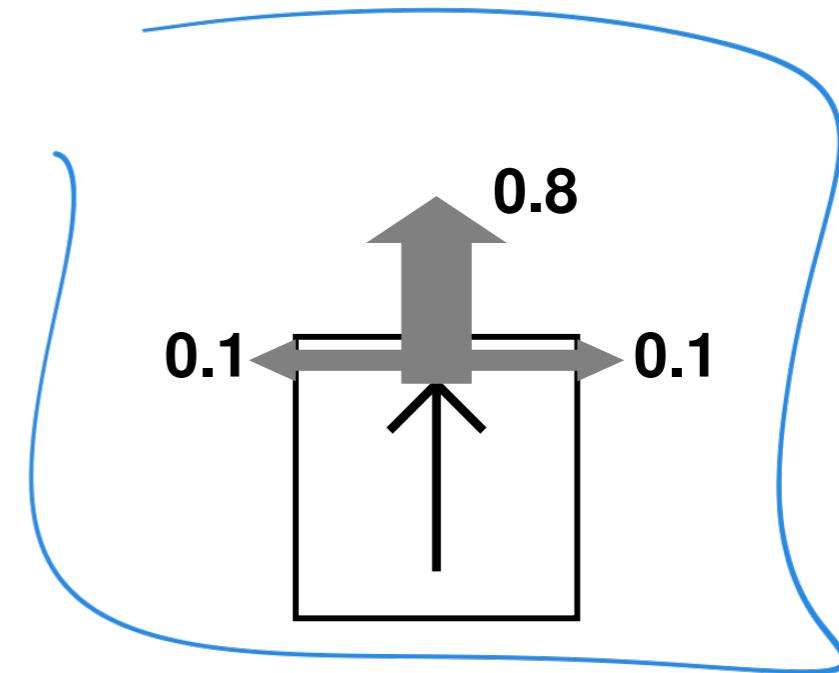
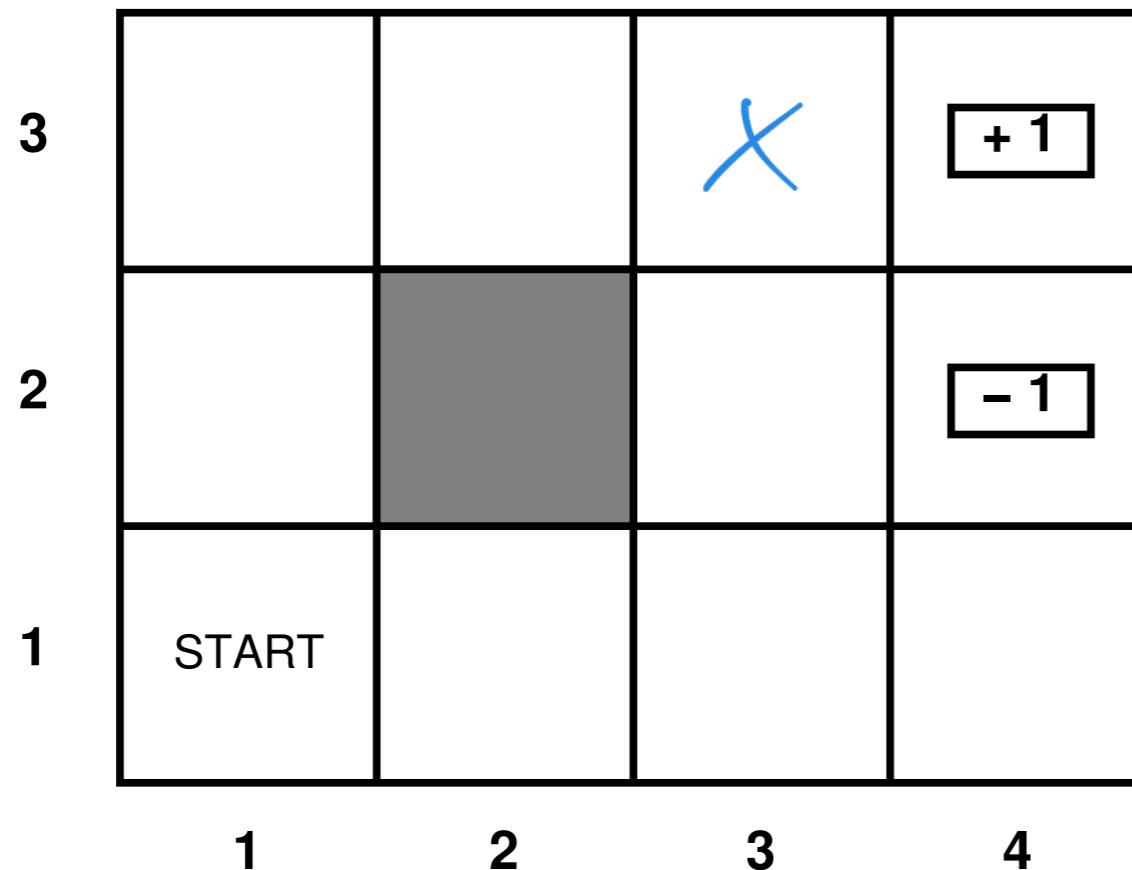
## Problem 11.1: Roomba Problem

**Problem 11.1.2** Show the corresponding **policy**.



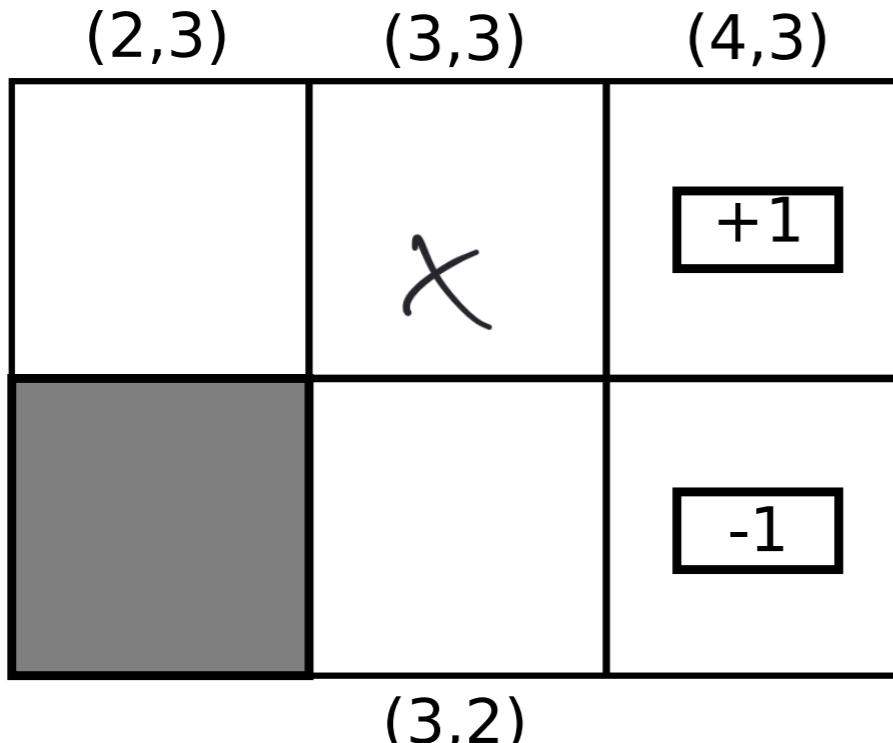
## Problem 11.1: Roomba Problem

**Problem 11.1.3** Assume that the transition probability is stochastic. Calculate the value of  $U(3, 3)$  using the **value iteration** algorithm for 2 iterations. Assume that all initial utilities are zero and  $U^1(1, 3) = -0.04$ ,  $U^1(2, 3) = -0.04$  and  $U^1(3, 2) = -0.04$ .



## Problem 11.1: Roomba Problem

**Problem 11.1.3** Assume that the transition probability is stochastic. Calculate the value of  $U(3, 3)$  using the **value iteration** algorithm for 2 iterations. Assume that all initial utilities are zero and  $U^1(1, 3) = -0.04$ ,  $U^1(2, 3) = -0.04$  and  $U^1(3, 2) = -0.04$ .



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

Bellman equation (if reward depends on state only)

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

$$V_{\pi^*}(s) = \max_a \mathbb{E} \left[ \sum_t \gamma^t R_t | a, s_t \right]$$

## Problem 11.1: Roomba Problem $a \in A(s) ; e$

**Problem 11.1.3** Compute  $U(3, 3)$

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

**Iteration 1**  $U^1(3, 3) = R(3, 3) + \gamma \max [$

$$\begin{aligned} & P((3, 3)|(3, 3), r) \cdot U^0(3, 3) + P((4, 3)|(3, 3), r) \cdot U^0(4, 3) + P((3, 2)|(3, 3), r) \cdot U^0(3, 2), \\ & P((3, 3)|(3, 3), l) \cdot U^0(3, 3) + P((2, 3)|(3, 3), l) \cdot U^0(2, 3) + P((3, 2)|(3, 3), l) \cdot U^0(3, 2), \\ & P((2, 3)|(3, 3), u) \cdot U^0(2, 3) + P((3, 3)|(3, 3), u) \cdot U^0(3, 3) + P((4, 3)|(3, 3), u) \cdot U^0(4, 3), \\ & P((2, 3)|(3, 3), d) \cdot U^0(2, 3) + P((3, 2)|(3, 3), d) \cdot U^0(3, 2) + P((4, 3)|(3, 3), d) \cdot U^0(4, 3) \end{aligned}$$

$$U^1(3, 3) = -0.04 + \max \left[ \begin{array}{l} (0.1 \cdot 0 + 0.8 \cdot 1 + 0.1 \cdot 0), \\ (0.1 \cdot 0 + 0.8 \cdot 0 + 0.1 \cdot 0), \\ (0.1 \cdot 0 + 0.8 \cdot 0 + 0.1 \cdot 1), \\ (0.1 \cdot 0 + 0.8 \cdot 0 + 0.1 \cdot 1) \end{array} \right]$$

$$U^1(3, 3) = 0.760(\text{Right})$$

## Problem 11.1: Roomba Problem

**Problem 11.1.3** Compute  $U(3, 3)$

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

**Iteration 2**

$$\begin{aligned}
 & P((3, 3)|(3, 3), r) \cdot U^1(3, 3) + P((4, 3)|(3, 3), r) \cdot U^1(4, 3) + P((3, 2)|(3, 3), r) \cdot U^1(3, 2), \quad (\text{Right}) \\
 & P((3, 3)|(3, 3), l) \cdot U^1(3, 3) + P((2, 3)|(3, 3), l) \cdot U^1(2, 3) + P((3, 2)|(3, 3), l) \cdot U^1(3, 2), \quad (\text{Left}) \\
 & P((2, 3)|(3, 3), u) \cdot U^1(2, 3) + P((3, 3)|(3, 3), u) \cdot U^1(3, 3) + P((4, 3)|(3, 3), u) \cdot U^1(4, 3), \quad (\text{Up}) \\
 & P((2, 3)|(3, 3), d) \cdot U^1(2, 3) + P((3, 2)|(3, 3), d) \cdot U^1(3, 2) + P((4, 3)|(3, 3), d) \cdot U^1(4, 3) \quad (\text{Down})
 \end{aligned}$$

## Problem 11.1: Roomba Problem

### Problem 11.1.3 Compute $U(3, 3)$

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

**Iteration 2**

$$U^2(3, 3) = -0.04 + \max \begin{aligned} & [(0.1 \cdot (0.760) + 0.8 \cdot (1) + 0.1 \cdot (-0.04)), && (\text{Right}) \\ & (0.1 \cdot 0.76 + 0.8 \cdot (-0.04) + 0.1 \cdot (-0.04)), && (\text{Left}) \\ & (0.1 \cdot (-0.04) + 0.8 \cdot (0.760) + 0.1 \cdot 1), && (\text{Up}) \\ & (0.1 \cdot (-0.04) + 0.8 \cdot (-0.04) + 0.1 \cdot 1)] && (\text{Down}) \end{aligned}$$

$$U^2(3, 3) = 0.832 \quad (\text{Right})$$

## Problem 11.1: Roomba Problem

**Problem 11.1.4** Compute the **optimal policy** of state (3, 1) after convergence. The utilities after convergence are given.

3	0.812	0.868	0.912	+ 1
2	0.762		0.660	- 1
1	0.705	0.655	0.611	0.388

## Problem 11.1: Roomba Problem

**Problem 11.1.4** Compute the **optimal policy** of state (3, 1) after convergence. The utilities after convergence are given.

	1	2	3	4
1	0.705	0.655	0.611	0.388
2	0.762		0.660	-1
3	0.812	0.868	0.912	+1

### Optimal Policy

$$\pi^*(s) = \arg \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U(s')$$

## Problem 11.1: Roomba Problem

### Problem 11.1.4

$$\pi^*(s) = \arg \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U(s')$$

$$\begin{aligned}
 &= \arg \max_a [0.1 V(3,2) + 0.8 V(4,1) + 0.1 V(3,1), \text{ } \leftarrow r \\
 &\quad 0.1 V(3,2) + 0.8 V(2,1) + 0.1 V(3,1), \text{ } \leftarrow l \\
 &\quad 0.1 V(2,1) + 0.8 V(3,2) + 0.1 V(4,1), \text{ } \leftarrow u \\
 &\quad 0.1 V(4,1) + 0.8 V(3,1) + 0.1 V(2,1)] \text{ } \leftarrow d
 \end{aligned}$$

$$\begin{aligned}
 &= \arg \max_a [0.4375, \text{ } \leftarrow r \\
 &\quad 0.6511, \text{ } \leftarrow l \\
 &\quad 0.6323, \text{ } \leftarrow u \\
 &\quad 0.5931] \text{ } \leftarrow d
 \end{aligned}$$

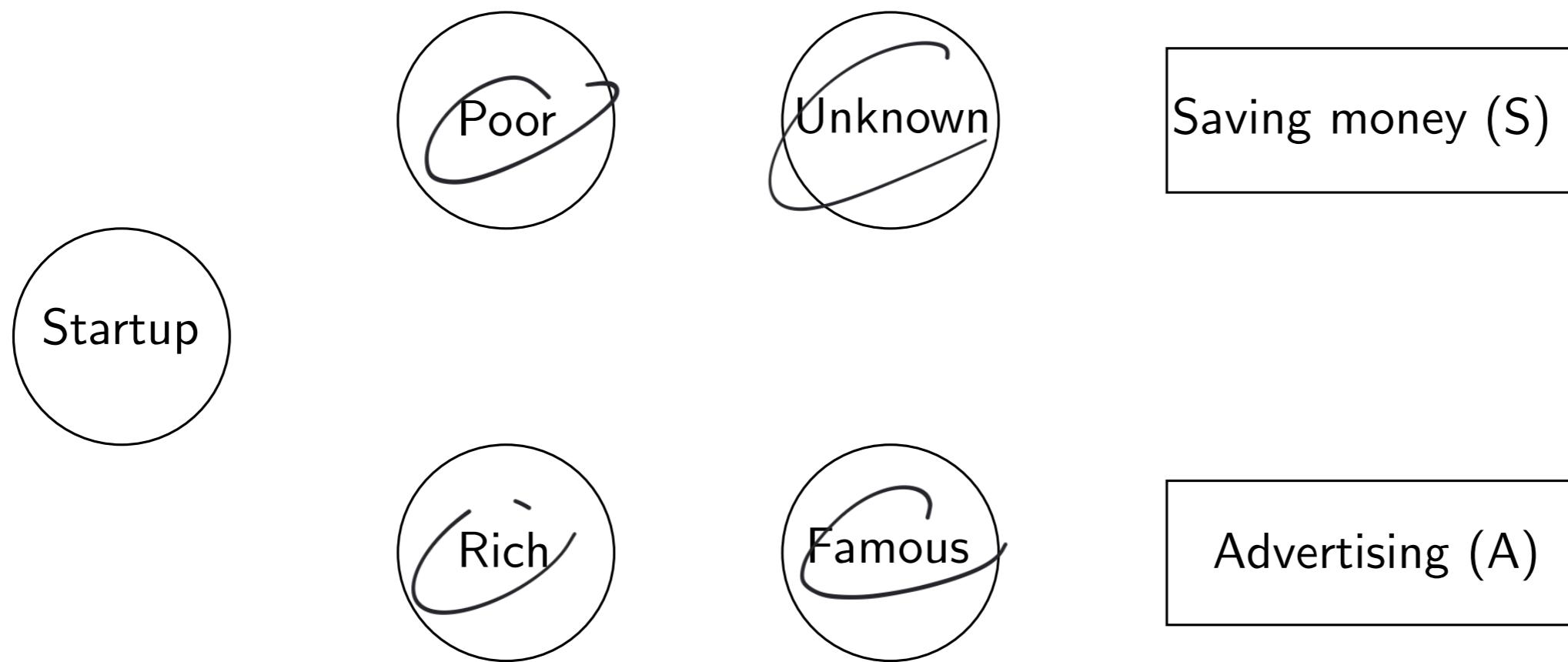
$$\pi^*(s) = l$$

0.812	0.868	0.912	+1
0.762		0.660	-1
0.705	0.655	0.611	0.388

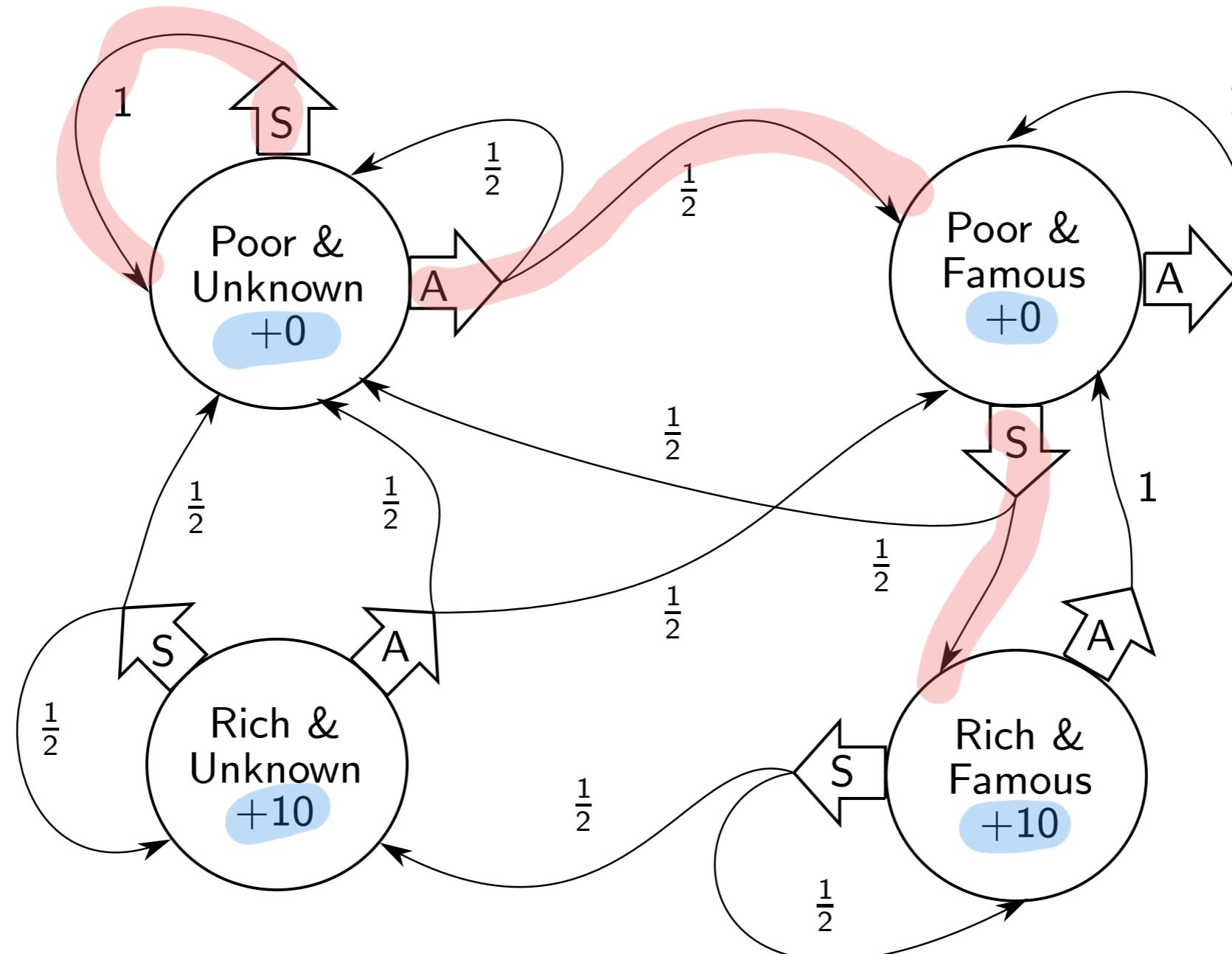
1      2      3      4

## Problem 11.2: Startup Dilemma

Assume that you run a startup company. In every decision period, you must choose between Saving money (S) or Advertising (A). If you advertise, you may become famous (f ) (50%) but because of spending money you may become poor (p). If you save money, you may become rich (r) with probability 50% but you may become also unknown (u) because you don't advertise.



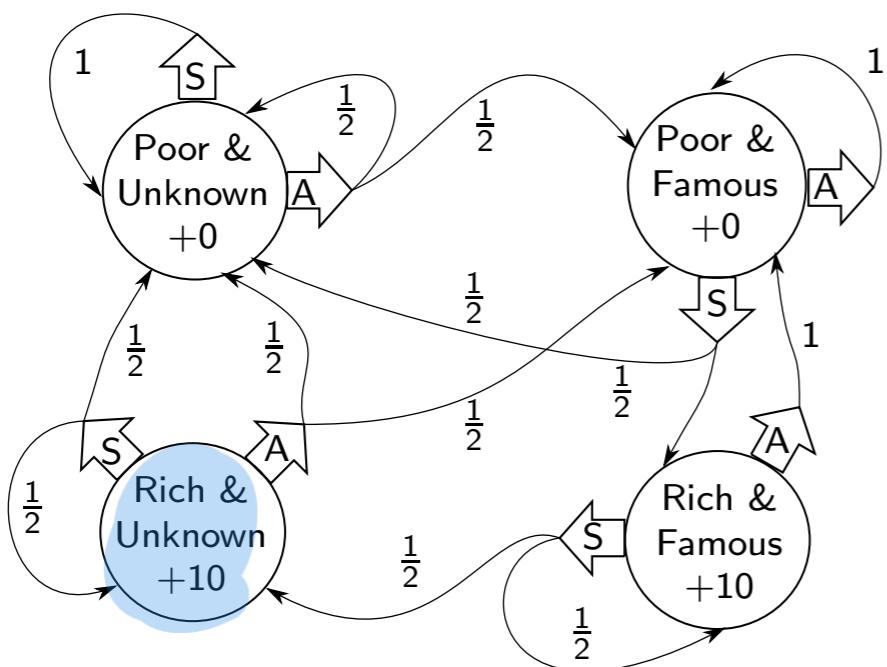
## Problem 11.2: Startup Dilemma



## Problem 11.2: Startup Dilemma

**Problem 11.2.1** Calculate the utility value for state  $U(r, u)$  for 2 iterations using value iteration.  
 Assume that the discount factor is 0.9 and that all initial states are zero. Furthermore use  
 $U^1(p, f) = 0, U^1(p, u) = 0.$

Locality



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

## Problem 11.2: Startup Dilemma

**Problem 11.2.1** Calculate the utility value for state  $U(r, u)$  for 2 iterations using value iteration. Assume that the discount factor is 0.9 and that all initial states are zero. Furthermore use  $U^1(p, f) = 0, U^1(p, u) = 0$ .

$$U^1(r, u) = 10$$

Iteration 1

$$U(s) = R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U^1(s')$$

↗ 0.9

$$U^2(r, u) = R(r, u) + \gamma \max_{(p, u)} [P((p, u)|(r, u), A) \cdot U^1(p, u) + P((p, f)|(r, u), A) \cdot U^1(p, f) + P((p, u)|(r, u), S) \cdot U^1(p, u) + P((r, u)|(r, u), S) \cdot U^1(r, u)],$$

(A)      (S)

$$= R(r, u) = 10$$

## Problem 11.2: Startup Dilemma

**Problem 11.2.1** Calculate the utility value for state  $U(r, u)$  for 2 iterations using value iteration. Assume that the discount factor is 0.9 and that all initial states are zero. Furthermore use  $U^1(p, f) = 0, U^1(p, u) = 0$ .

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

**Iteration 2:**

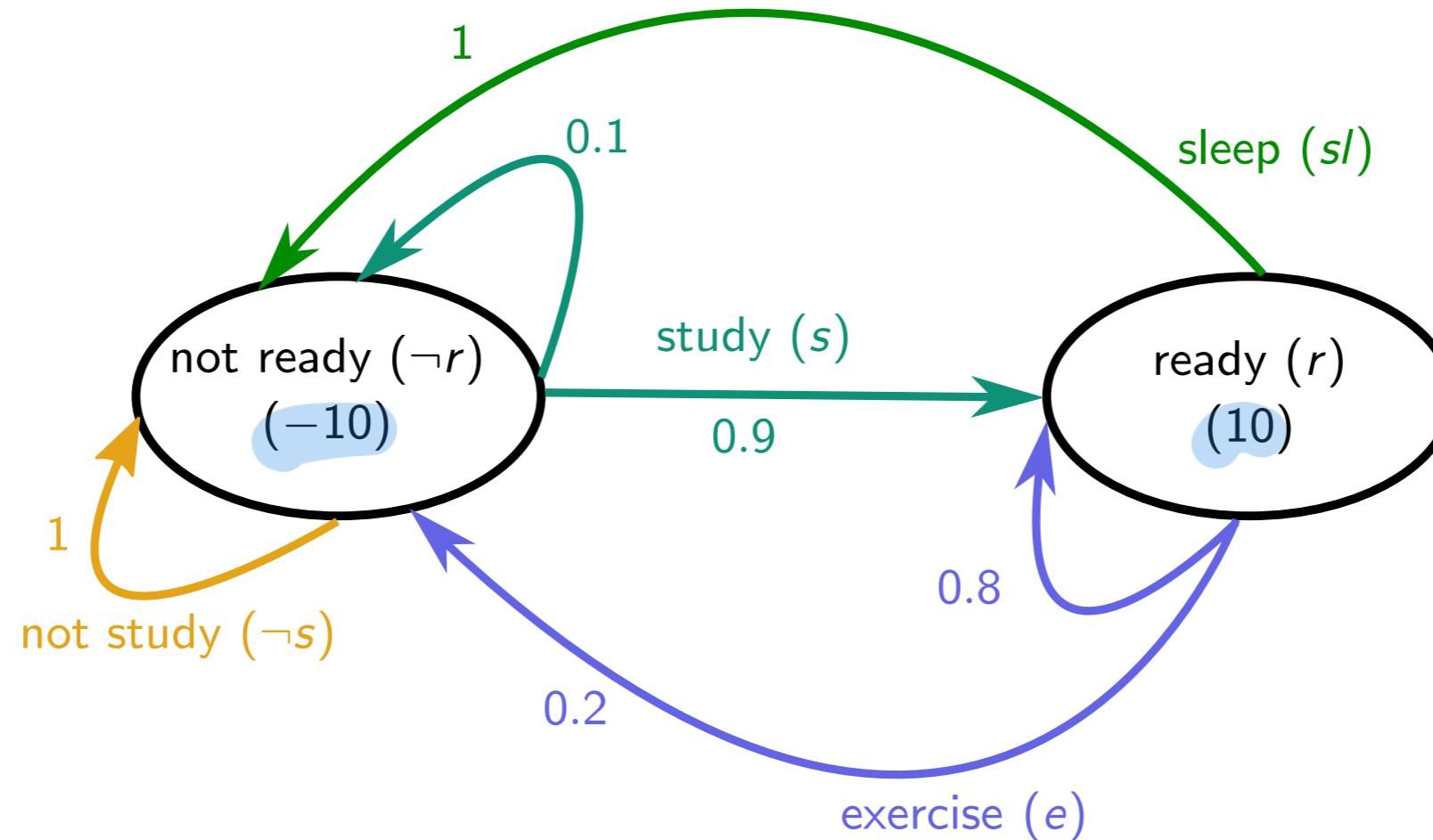
$$U^2(r, u) = R(r, u) + \gamma \max [P((p, u)|(r, u), A) \cdot U^1(p, u) \quad (A) \\ + P((p, f)|(r, u), A) \cdot U^1(p, f), \\ P((p, u)|(r, u), S) \cdot U^1(p, u) \quad (S) \\ + P((r, u)|(r, u), S) \cdot U^1(r, u)],$$

$$U^2(r, u) = 10 + 0.9 \max [0.5 \cdot 0 + 0.5 \cdot 0, \quad (A) \\ 0.5 \cdot 0 + 0.5 \cdot 10, \quad (S)]$$

$\boxed{U^2(r, u) = 14.5}$

## Problem 11.3: AI Exam

$$R(s) = \begin{cases} 10, & s = r \\ -10, & s = \neg r \end{cases}$$



## Problem 11.3: AI Exam

$\pi_1(s)$

Apply the **policy iteration** algorithm for one iteration in order to determine the policies  $\pi_1(\neg r)$  and  $\pi_1(r)$ . Assume that the discount factor is  $\gamma = 0.9$  and the initial policies are  $\pi_0(\neg r) = s$  and  $\pi_0(r) = e$ . The rewards for  $\neg r$  and  $r$  are  $-10$  and  $10$ , respectively.

## Problem 11.3: AI Exam

Apply the **policy iteration** algorithm for one iteration in order to determine the policies  $\pi_1(\neg r)$  and  $\pi_1(r)$ . Assume that the discount factor is  $\gamma = 0.9$  and the initial policies are  $\pi_0(\neg r) = s$  and  $\pi_0(r) = e$ . The rewards for  $\neg r$  and  $r$  are  $-10$  and  $10$ , respectively.

### Policy iteration

- **Policy evaluation:** Given a policy  $\pi_i$ , calculate  $U_i = U^{\pi_i}$ , the utility of each state if  $\pi_i$  were to be executed.
- **Policy improvement:** Calculate a new policy  $\pi_{i+1}$  using a one-step look-ahead based on  $U_i$  using  $\pi_{i+1}(s) = \arg \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U_i(s')$ .

## Problem 11.3: AI Exam

Apply the **policy iteration** algorithm for one iteration in order to determine the policies  $\pi_1(\neg r)$  and  $\pi_1(r)$ . Assume that the discount factor is  $\gamma = 0.9$  and the initial policies are  $\pi_0(\neg r) = s$  and  $\pi_0(r) = e$ . The rewards for  $\neg r$  and  $r$  are  $-10$  and  $10$ , respectively.

**Step 1. Policy evaluation**      
$$U_i(s) = R(s) + \gamma \sum_{s'} P(s'|s| \underbrace{\pi_i(s)}_{\pi_i}) U_i(s')$$

Compute  $U_0(r)$  and  $U_0(\neg r)$ :

$$\begin{aligned} U_0^{\pi_0}(r) &= R(r) + \gamma [P(r|r,e) U_0^{\pi_0}(r) + P(\neg r|r,e) U_0^{\pi_0}(\neg r)] \\ &= 10 + 0.9 [0.8 U_0^{\pi_0}(r) + 0.2 U_0^{\pi_0}(\neg r)] \end{aligned}$$

---


$$10 = 0.28 U_0^{\pi_0}(r) - 0.18 U_0^{\pi_0}(\neg r)$$

## Problem 11.3: AI Exam

Apply the **policy iteration** algorithm for one iteration in order to determine the policies  $\pi_1(\neg r)$  and  $\pi_1(r)$ . Assume that the discount factor is  $\gamma = 0.9$  and the initial policies are  $\pi_0(\neg r) = s$  and  $\pi_0(r) = e$ . The rewards for  $\neg r$  and  $r$  are  $-10$  and  $10$ , respectively.

**Step 1. Policy evaluation**  $U_i(s) = R(s) + \gamma \sum_{s'} P(s'|s, \pi_i(s)) U_i(s')$

$$\text{Compute } U_0(r) \text{ and } U_0(\neg r): U_0^{\pi^*}(\neg r) = 2(\neg r) + 0.9 [P(r | \neg r, s) U_0^{\pi^*}(r) + P(\neg r | \neg r, s) U_0^{\pi^*}(\neg r)]$$

$$= -10 + 0.9 [0.9 U_0^{\pi^*}(r) + 0.1 U_0^{\pi^*}(\neg r)]$$

$$-10 = 0.91 U_0^{\pi^*}(r) - 0.09 U_0^{\pi^*}(\neg r)$$

## Problem 11.3: AI Exam

Apply the **policy iteration** algorithm for one iteration in order to determine the policies  $\pi_1(\neg r)$  and  $\pi_1(r)$ . Assume that the discount factor is  $\gamma = 0.9$  and the initial policies are  $\pi_0(\neg r) = s$  and  $\pi_0(r) = e$ . The rewards for  $\neg r$  and  $r$  are  $-10$  and  $10$ , respectively.

**Step 1. Policy evaluation** 
$$U_i(s) = R(s) + \gamma \sum_{s'} P(s'|s, \pi_i(s)) U_i(s')$$

Summarize the linear equations:

$$\left| \begin{array}{l} 0.91 \cdot U_0(\neg r) - 0.81 \cdot U_0(r) = -10 \\ (-0.18) \cdot U_0(\neg r) + 0.28 \cdot U_0(r) = 10 \end{array} \right.$$

Solution:

$$\begin{aligned} U_0^{\pi}(r) &= 66.7, \\ U_0^{\pi}(\neg r) &= 48.4. \end{aligned}$$

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

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

$$= \max_{a \in A} \left[ \sum_{s'} P(s' | s, a) \underbrace{R(s)}_{\text{---}} + \sum_{s'} \gamma P(s' | s, a) V(s') \right]$$

$$= \left[ \max_{a \in A} \left( \sum_{s'} P(s' | s, a) \right) \overbrace{R(s)}^{\text{---}} + \dots \right]$$

$$= 1 \circ R(s) + \max_{a \in A} \sum_{s'} \gamma P(s' | a, s) V(s')$$

## Problem 11.3: AI Exam

Apply the **policy iteration** algorithm for one iteration in order to determine the policies  $\pi_1(\neg r)$  and  $\pi_1(r)$ . Assume that the discount factor is  $\gamma = 0.9$  and the initial policies are  $\pi_0(\neg r) = s$  and  $\pi_0(r) = e$ . The rewards for  $\neg r$  and  $r$  are  $-10$  and  $10$ , respectively.

**Step 2. Policy improvement**  $\pi_{i+1}(s) = \arg \max_{\substack{s' \\ a \in A(s)}} P(s'|s, a) U_i(s')$

$$\text{Compute } \pi_1(\neg r): \arg \max_a \left[ P(r|\neg r, s) V_0^{\pi}(r) + P(\neg r|\neg r, s) V_0^{\pi}(\neg r); P(\neg r|\neg r, \neg s) V_0^{\pi}(\neg r) \right]$$

$$= \arg \max_a [0.9 \cdot 66.7 + 0.1 \cdot 48.4; 1 \cdot 48.4]$$

$$= [64.87; 48.4] = \begin{cases} s \\ r \end{cases}$$

## Problem 11.3: AI Exam

Apply the **policy iteration** algorithm for one iteration in order to determine the policies  $\pi_1(\neg r)$  and  $\pi_1(r)$ . Assume that the discount factor is  $\gamma = 0.9$  and the initial policies are  $\pi_0(\neg r) = s$  and  $\pi_0(r) = e$ . The rewards for  $\neg r$  and  $r$  are  $-10$  and  $10$ , respectively.

**Step 2. Policy improvement**  $\pi_{i+1}(s) = \arg \max_{a \in A(s)} \sum_{s'} P(s'|s, a) U_i(s')$

$$\text{Compute } \pi_1(r): \pi_1(r) = \arg \max_{a \in \{e, \text{sl}\}} \left[ P(\neg r | r, \text{sl}) U_0(\neg r) + P(r | r, e) U_0(r) + P(\neg r | r, e) U_0(\neg r) \right]$$

$$= \arg \max_a [1 \cdot 48.4; 0.8 \cdot 66.7 + 0.2 \cdot 48.4]$$

$$= \arg \max_a [48.4; 63.04] = e$$