

PROBLEM SESSION 4: EXACT SOLUTION METHODS

October 18, 2023 4:00pm PT

Topic 1. MDP Overview

a) Markov Decision Process (MDP): defined by the tuple $(\mathcal{S}, \mathcal{A}, T, R, \gamma)$

- \mathcal{S} - State Space: the environment, the *minimum information set* required to make a decision

- \mathcal{A} - Action Space: what the agent can do

- T - Transition model: system dynamics (how the system evolves)

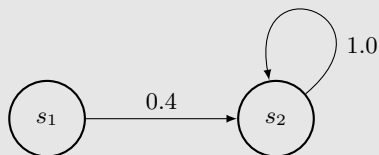


Figure 1: Simple MDP

	s_1	s_2
s_1		
s_2		

Figure 2: Transition Model

- R - Reward Function: expected reward from taking action a in state s and transitioning to state s'

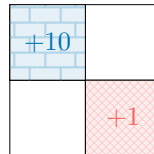
- γ - Discount factor: used to weight future rewards

b) Utility: a discounted sequence of rewards

- Utility of a sequence of states **without** discounting: why is this problematic?

$$U([s_1, s_2, \dots, s_n]) = \sum_{t=1}^n r_t$$

- Thought exercise: would an agent want to collect rewards in the blue cell (bricks) or the red cell (crosshatch) forever?



Is there a preference for $10 + 10 + 10 + \dots$ or $1 + 1 + 1 + \dots$ as $n \rightarrow \infty$ (“infinite horizon”)?

- Solution: discount with γ !

c) Policy π : a function of the state that tells you *what to do* in every state

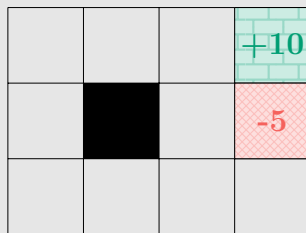


Figure 3: Optimal Path

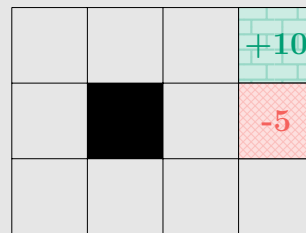
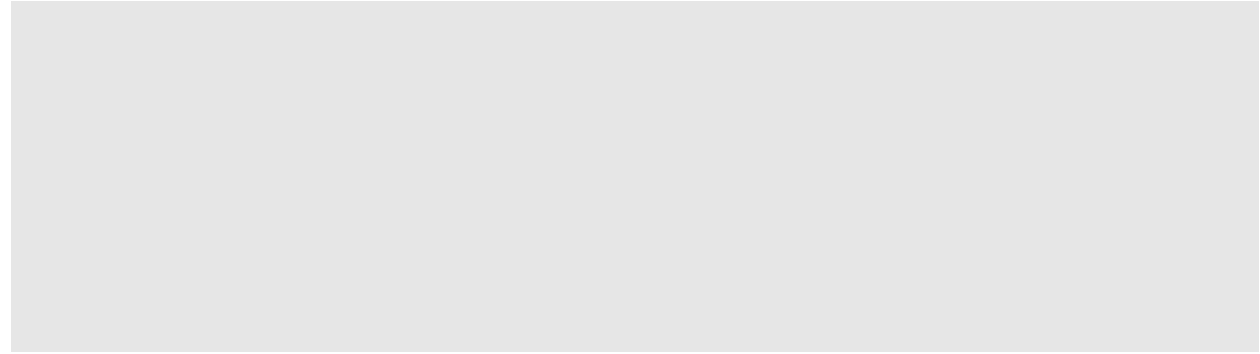


Figure 4: Optimal Policy

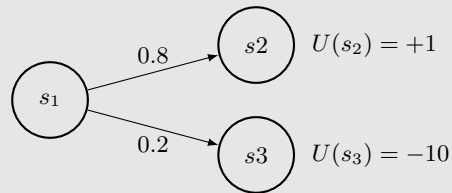
- $U^\pi(s) \rightarrow$ utility from executing policy π from state s (the *value function*)
- $\pi^*(s) = \arg \max_{\pi} U^\pi(s)$

- d) Bellman Equation: “The expected utility of a state is the reward at that state plus the discounted sum of expected future rewards.”



- e) A Note On Expectation

An expected value is a *weighted average*



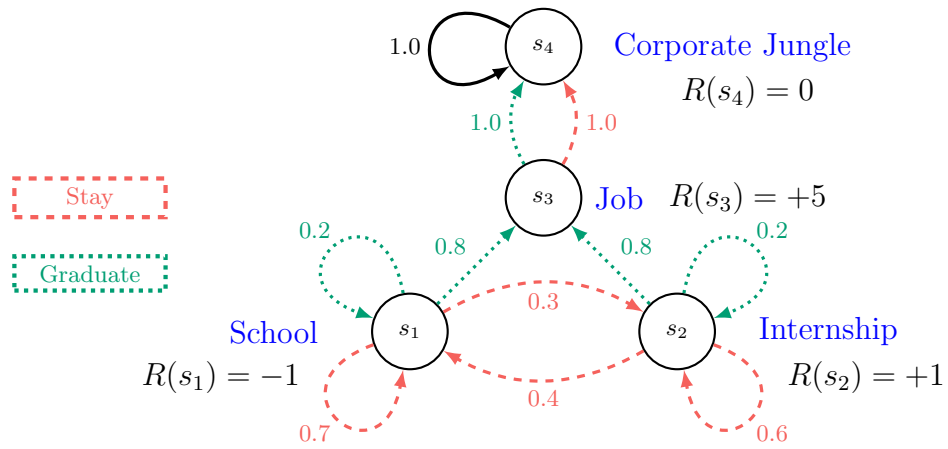
What is the expected utility when transitioning out of s_1 ?

Topic 2. Value Iteration Example

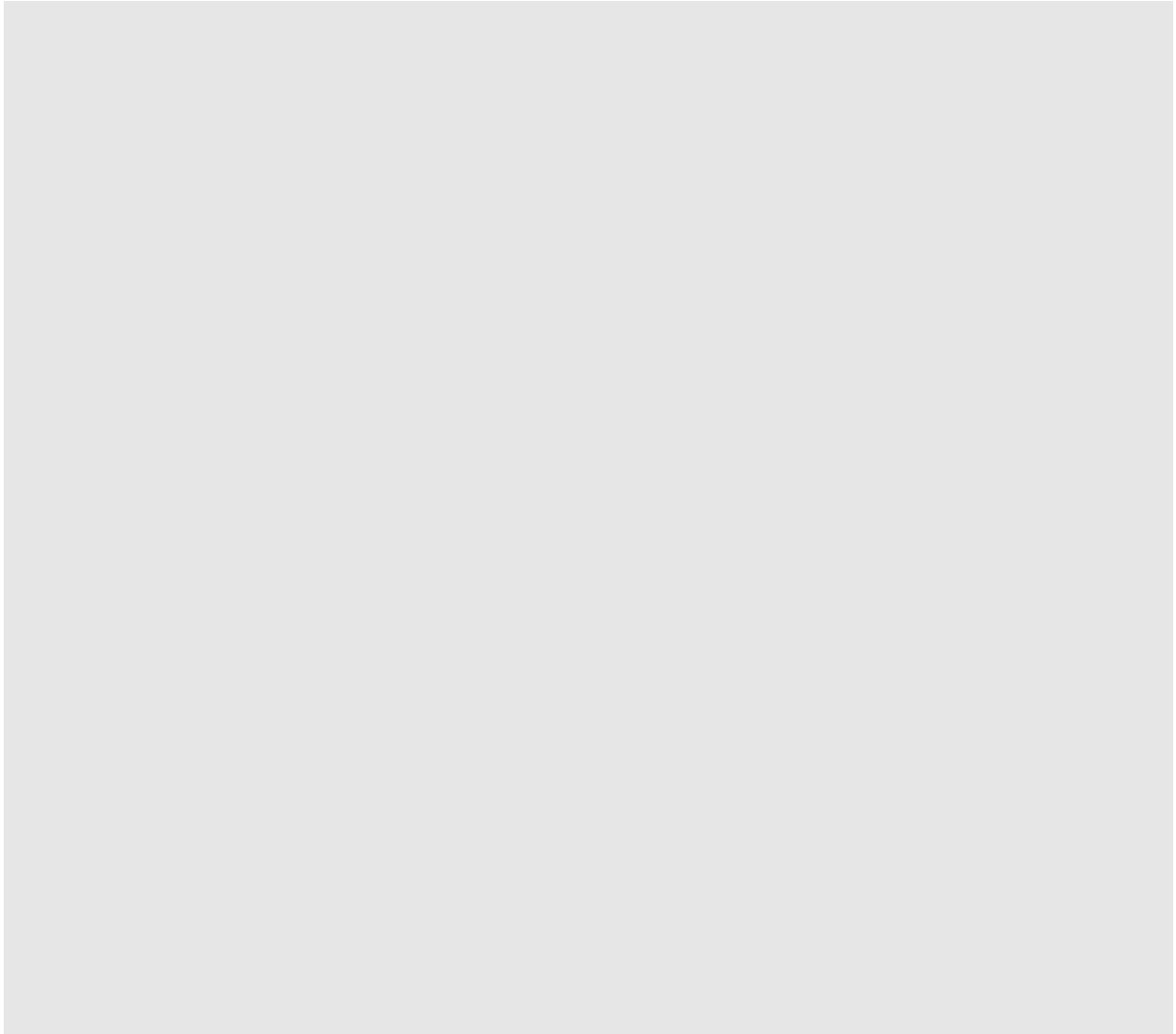
Algorithm 1 The Value Iteration Algorithm

```

1: procedure VALUE ITERATION( $\mathcal{P} :: \text{MDP}, k_{\max}$ )
2:    $U(s) \leftarrow 0$  for all  $s \in \mathcal{S}$ 
3:   for  $k \leftarrow 1, k_{\max}$  do
4:     for all  $s \in \mathcal{P}.\mathcal{S}$  do
5:        $U_{k+1}(s) = \max_a (R(s, a) + \gamma \sum_{s'} T(s' | s, a) U_k(s'))$ 
6:     end for
7:   end for
8: end procedure
  
```



a) Define the tuple for this MDP



b) Perform two iterations of value iteration:

c) What is our policy after two rounds of value iteration?

d) What is the time complexity of value iteration?