

## Advanced Topics in AI

### Exercise 4

#### Question 1: MDPs: Micro-Blackjack

In micro-blackjack, you repeatedly draw a card (with replacement) that is equally likely to be a 2, 3, or 4. You can either Draw or Stop if the total score of the cards you have drawn is less than 6. If your total score is 6 or higher, the game ends, and you receive a utility of 0. When you Stop, your utility is equal to your total score (up to 5), and the game ends. When you Draw, you receive no utility. There is no discount ( $\gamma = 1$ ). Let's formulate this problem as an MDP with the following states: 0, 2, 3, 4, 5 and a Done state, for when the game ends.

- a. What is the transition function and the reward function for this MDP?

**Solution:**

The transition function is

$$T(s, \text{STOP}, \text{DONE}) = 1$$

$$T(0, \text{DRAW}, s') = 1/3 \text{ for } s' \in \{2, 3, 4\}$$

$$T(2, \text{DRAW}, s') = 1/3 \text{ for } s' \in \{4, 5, \text{DONE}\}$$

$$T(3, \text{DRAW}, s') = \begin{cases} 1/3 & \text{if } s' = 5 \\ 2/3 & \text{if } s' = \text{DONE} \end{cases}$$

$$T(4, \text{DRAW}, \text{DONE}) = 1$$

$$T(5, \text{DRAW}, \text{DONE}) = 1$$

$$T(s, a, s') = 0 \text{ otherwise}$$

The reward function is

$$R(s, \text{STOP}, \text{DONE}) = s, s \leq 5$$

$$R(s, a, s') = 0 \text{ otherwise}$$

- b. Fill in the following table of value iteration values for the first 4 iterations.

States	0	2	3	4	5
$V_0$	0	0	0	0	0
$V_1$	0	2	3	4	5
$V_2$	3	3	3	4	5
$V_3$	$10/3$	3	3	4	5
$V_4$	$10/3$	3	3	4	5

- c. You should have noticed that value iteration converged above. What is the optimal policy for the MDP?

States	0	2	3	4	5
$\pi^*$	DRAW	DRAW	STOP	STOP	STOP

- d. Perform one iteration of policy iteration for one step of this MDP, starting from the fixed policy below:

States	0	2	3	4	5
$\pi_i$	DRAW	STOP	DRAW	STOP	DRAW
$V^{\pi_i}$	2	2	0	4	0
$\pi_{i+1}$	DRAW	STOP	STOP	STOP	STOP

- e. Is the policy  $\pi_{i+1}$  optimal?

**Solution:**

No it isn't. Compare the policy with the policy in task c.

## Question 2: Golf as an MDP

In this exercise we will formulate golf as an MDP as follows:

- State Space : {Tee, Fairway, Sand, Green}
- Actions : {Conservative shot, Power shot}
- Initial State : Tee
- Terminal State : Green
- Transition model : (note that any successor state not on this list has a transition probability 0, and “Conservative” stands for “Conservative shot”)

$s$	$a$	$s'$	$T(s, a, s')$
Tee	Conservative	Fairway	0.9
Tee	Conservative	Sand	0.1
Tee	Power shot	Green	0.5
Tee	Power shot	Sand	0.5
Fairway	Conservative	Green	0.8
Fairway	Conservative	Sand	0.2
Sand	Conservative	Green	1.0

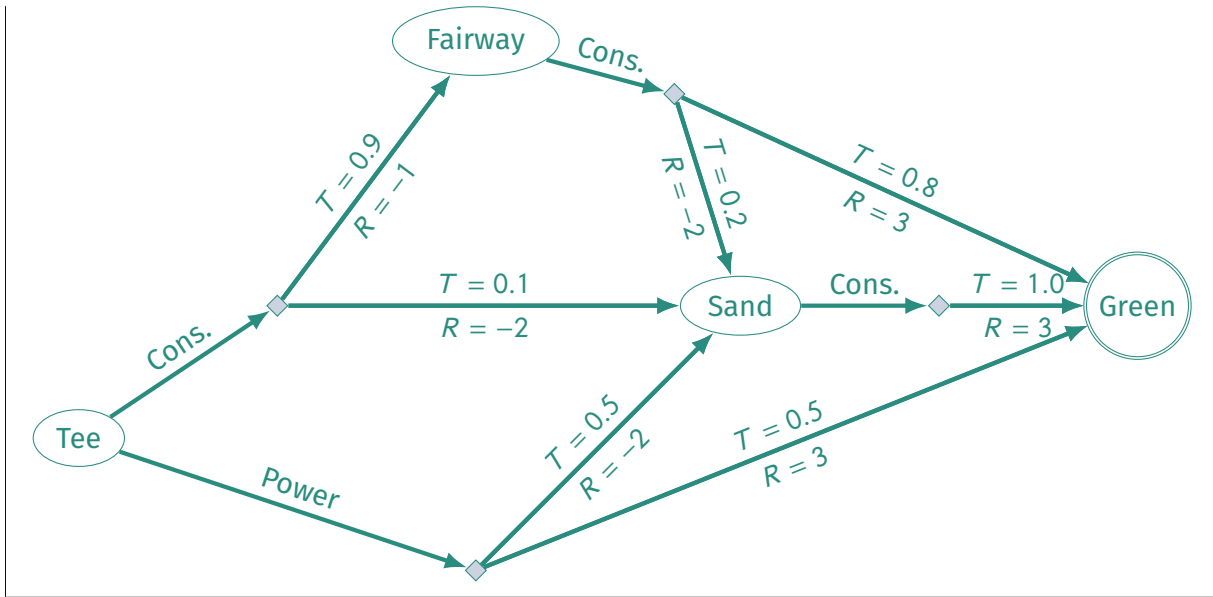
- Rewards: (note:  $R(\cdot, \cdot, s)$  means that the reward is received for transitioning to state  $s$ , regardless of the action taken or previous state)

$s'$	$R(\cdot, \cdot, s')$
Fairway	-1
Sand	-2
Green	3

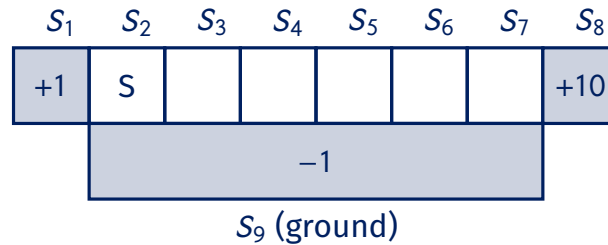
Draw a state graph defining this MDP problem. A state graph shows the states as nodes and has the actions as arcs labeled with T and R values. Remember, in a state graph no states repeat.

### Solution:

In the graph, “Cons.” stands for “Conservative Shot” and “Power” for “Power Shot”. Absent edges represent a transition of 0 probability.



### Question 3: Robot Balancing



Consider the above MDP, representing a robot on a balance beam. Each grid square is a state and the available actions are right and left. The agent starts in state  $s_2$ , and all states have reward 0 aside from the ends of the grid  $s_1$  and  $s_8$  and the ground state, which have the rewards shown. Moving left or right results in a move left or right (respectively) with probability  $p$ . With probability  $1 - p$ , the robot falls off the beam (transitions to ground, and receives a reward of -1). Falling off, or reaching either endpoint, result in the end of the episode (i.e., they are terminal states). Note that terminal states receive no future reward.

- a. For what values of  $p$  is the optimal action from  $s_2$  to move right if the discount  $\gamma$  is 1?

**Solution:**

$$U[\text{left}] = (+1)(p) + (-1)(1 - p) = 2p - 1$$

$$U[\text{always go right}] = 10p^6 - 1 \cdot (1 - p^6) > 2p - 1$$

$$\Leftrightarrow (11/2)p^5 > 1 \Leftrightarrow p > 0.71 \text{ (approximately)}$$

- b. For what values of  $\gamma$  is the optimal action from  $s_2$  to move right if  $p = 1$ ?

**Solution:**

$$10 * \gamma^6 > \gamma \Leftrightarrow \gamma > 0.63 \text{ (approximately)}$$

- c. Given initial value estimates of zero, show the results of one, then two rounds of value iteration (assume  $\gamma = 1$ ).

**Solution:**

States	$S_1$	$S_2$	$S_3$	$S_4$	$S_5$	$S_6$	$S_7$	$S_8$	$S_9$
$V_0$	0	0	0	0	0	0	0	0	0
$V_1$	0	$2p - 1$	$p - 1$	$p - 1$	$p - 1$	$p - 1$	$11p - 1$	0	0
$V_2$	0	$2p - 1$	$2p^2 - 1$	$p^2 - 1$	$p^2 - 1$	$11p^2 - 1$	$11p - 1$	0	0

### Question 4: Dice Bonanza

A casino is considering adding a new game to their collection, but need to analyze it before releasing it on their floor. They have hired you to execute the analysis. On each round of the game, the player has the option of rolling a fair 6-sided die. That is, the die lands on values 1 through 6 with equal probability. Each roll costs 1 dollar, and the player must roll the very first round. Each time the player rolls the die, the player has two possible actions:

1. *Stop*: Stop playing by collecting the dollar value that the die lands on, or
2. *Roll*: Roll again, paying another 1 dollar.

Having taken "Introduction to Artificial Intelligence", you decide to model this problem using an infinite horizon Markov Decision Process (MDP). The player initially starts in state *Start*, where the player only has one possible action: *Roll*. State  $s_i$  denotes the state where the die lands on  $i$ . Once a player decides to *Stop*, the game is over, transitioning the player to the *End* state.

- (a) In solving this problem, you consider using policy iteration. Your initial policy  $\pi$  is in the table below. Evaluate the policy at each state, with  $\gamma = 1$ . Build a system of linear equations and solve directly for  $v^\pi(s_i)$ ,  $i \in \{1, \dots, 6\}$ .

State	$s_1$	$s_2$	$s_3$	$s_4$	$s_5$	$s_6$
$\pi(s)$	<i>Roll</i>	<i>Roll</i>	<i>Stop</i>	<i>Stop</i>	<i>Stop</i>	<i>Stop</i>
$v^\pi(s)$	3	3	3	4	5	6

**Solution:**

We have that  $v^\pi(s_i) = i$  for  $i \in \{3, 4, 5, 6\}$ , since the player will be awarded no further rewards according to the policy. From the Bellman equations, we have that  $V(s_1) = -1 + 1/6(V(s_1) + V(s_2) + 3 + 4 + 5 + 6)$  and that  $V(s_2) = -1 + 1/6(V(s_1) + V(s_2) + 3 + 4 + 5 + 6)$ . Solving this linear system yields  $V(s_1) = V(s_2) = 3$ .

- (b) Having determined the values, perform a policy update to find the new policy  $\pi'$ . The table below shows the old policy  $\pi$  and has filled in parts of the updated policy  $\pi'$  for you. If both Roll and Stop are viable new actions for a state, write down both *Roll/Stop*. In this part as well, we have  $\gamma = 1$ .

State	$s_1$	$s_2$	$s_3$	$s_4$	$s_5$	$s_6$
$\pi(s)$	<i>Roll</i>	<i>Roll</i>	<i>Stop</i>	<i>Stop</i>	<i>Stop</i>	<i>Stop</i>
$\pi'(s)$	<i>Roll</i>	<b>Roll</b>	<b>Roll/Stop</b>	<b>Stop</b>	<b>Stop</b>	<i>Stop</i>

**Solution:**

For each  $s_i$  in part (a), we compare the values obtained via Rolling and Stopping. The value of Rolling for each state  $s_i$  is  $-1 + 1/6(3 + 3 + 3 + 4 + 5 + 6) = 3$ . The value of Stopping for each state  $s_i$  is  $i$ . At each state  $s_i$ , we take the action that yields the largest value; so, for  $s_1$  and  $s_2$ , we Roll, and for  $s_4$  and  $s_5$ , we stop. For  $s_3$ , we Roll/Stop, since the values from Rolling and Stopping

are equal.

(c) Is  $\pi(s)$  from part (a) optimal? Explain why or why not.

**Solution:**

Yes, the old policy is optimal. Looking at part (b), there is a tie between 2 equally good policies that policy iteration considers employing. One of these policies is the same as the old policy. This means that both new policies are as equally good as the old policy, and policy iteration has converged. Since policy iteration converges to the optimal policy, we can be sure that  $\pi(s)$  from part (a) is optimal.

(d) Suppose that we were now working with some  $\gamma \in [0, 1]$  and wanted to run **value iteration**. Select the **one** statement that would hold true at convergence, or write the correct answer next to Other if none of the options are correct.

- $V^*(s_i) = \max \left\{ -1 + \frac{i}{6}, \sum_j \gamma V^*(s_j) \right\}$
- $V^*(s_i) = \max \left\{ i, \frac{1}{6} \left[ -1 + \sum_j \gamma V^*(s_j) \right] \right\}$
- $V^*(s_i) = \max \left\{ -\frac{1}{6} + i, \sum_j \gamma V^*(s_j) \right\}$
- $V^*(s_i) = \max \left\{ i, -\frac{1}{6} + \sum_j \gamma V^*(s_j) \right\}$
- $V^*(s_i) = \frac{1}{6} \cdot \sum_j \max \left\{ i, -1 + \gamma V^*(s_j) \right\}$
- $V^*(s_i) = \frac{1}{6} \cdot \sum_j \max \left\{ -1 + i, \sum_k V^*(s_k) \right\}$
- $V^*(s_i) = \sum_j \max \left\{ -1 + i, \frac{1}{6} \gamma V^*(s_j) \right\}$
- $V^*(s_i) = \sum_j \max \left\{ \frac{i}{6}, -1 + \gamma V^*(s_j) \right\}$
- $V^*(s_i) = \max \left\{ i, -1 + \frac{\gamma}{6} \sum_j V^*(s_j) \right\}$
- $V^*(s_i) = \sum_j \max \left\{ i, -\frac{1}{6} + \gamma V^*(s_j) \right\}$
- $V^*(s_i) = \sum_j \max \left\{ -\frac{i}{6}, -1 + \gamma V^*(s_j) \right\}$

**Solution:**

At convergence,

$$\begin{aligned}
 V^*(s_i) &= \max Q^*(s_i, a) \\
 &= \max \{ Q^*(s_i, stop), Q^*(s_i, roll) \} \\
 &= \max \left\{ R(s_i, stop), R(s_i, roll) + \gamma \sum_j T(s_i, roll, s_j) V^*(s_j) \right\} \\
 &= \max \left\{ i, -1 + \frac{\gamma}{6} \sum_j V^*(s_j) \right\}
 \end{aligned}$$

## Question 5: Policy Evaluation with Q-Values

In this question, you will be working in an MDP with states  $S$ , actions  $A$ , discount factor  $\gamma$ , transition function  $T$ , and reward function  $R$ .

We have some fixed policy  $\pi : S \rightarrow A$ , which returns an action  $a = \pi(s)$  for each state  $s \in S$ . We want to learn the  $Q$  function  $Q^\pi(s, a)$  for this policy: the expected reward from taking action  $a$  in state  $s$  and then continuing to act according to  $\pi : Q^\pi(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma Q^\pi(s', \pi(s'))]$ . The policy  $\pi$  will not change.

Can we guarantee anything about how the values  $Q^\pi$  compared to the values  $Q^*$  for an optimal policy  $\pi^*$  ?

- $Q^\pi(s, a) \leq Q^*(s, a)$  for all  $s, a$
- $Q^\pi(s, a) = Q^*(s, a)$  for all  $s, a$
- $Q^\pi(s, a) \geq Q^*(s, a)$  for all  $s, a$
- None of the above guaranteed

### Solution:

We know that the optimal policy maximizes the sum of discounted rewards. If we take any action and then act according to a suboptimal policy the future Q-values will be smaller or equal to the optimal ones. So in total  $Q^\pi(s, a)$  will be smaller or equal to  $Q^*(s, a)$ .