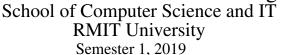


COSC1125/1127 Artificial Intelligence





Tutorial Sheet 8 Markov Decision Processes (MDPs)

1. Recall a Markov Decision Process (MDP) represents a scenario where actions are non-deterministic. Consider the following minor gridworld example:

	1	2	3
1	X		X
2			
3	X		

Cells are indexed by (row, column). Cells (1,1), (3,1) and (1,3) are terminal cells (marked with X), and have rewards of 20, -20 and 5, respectively. For all other cells, there is an reward of -1. The agent starts at cell (3,3). The agent can move in north, east, south and west directions or stay where it is. The actions are stochastic and have the following outcomes:

- If go north, 70% of the time the agent will go north, 30% will go west.
- If go east, 60% of the time the agent will go east, 40% of the time will go south.
- If go south, 80% of the time the agent will go south, 20% of the time will go east.
- If go west, 70% of the time the agent will go west, 30% of the time will go north.
- If stay in current cell, 100% of the time will stay.

If an action causes the agent to travel into a wall, the agent will stay in their current cell.

Construct a MDP for this instance of gridworld. Recall a MDP consists of S (set of states), s_0 (starting state), possible terminal states, A (set of actions), T (transition model/function) and R (reward function), so all of these have to be specified. To help, consider constructing the MDP in steps:

a) What are the states in this gridworld?

Answer

Each grid corresponds to a state.

b) What are the starting and terminal states?

Answer

Starting state is (3,3), terminal states are (1,1), (3,1) and (1,3).

c) What are the set of actions?

Answer

For each state, actions are {north, east, south, west, stay}

d) With cell (3,3) as the current state s, construct the transition model/function for all actions and subsequent states, i.e., T(s,a,s'), where a are all possible actions for (3,3) and s' are all possible successor states from (3,3).

Answer

- T((3,3), north, (2,3)) = 0.7 (70%, go north)
- T((3,3), north, (3,2)) = 0.3 (30%, go west)
- T((3,3), east, (3,3)) = 0.6 (60%, go east)
- T((3,3), east, (3,3)) = 0.4 (40%, go south)
- T((3,3), south, (3,3)) = 0.8 (80%, go south)
- T((3,3), south, (3,3)) = 0.20 (20%, go east)
- T((3,3), west, (3,2)) = 0.7 (70%, go west)
- T((3,3), west, (2,3)) = 0.3 (30%, go north)
- T((3,3), stay, (3,3)) = 1.0 (100%, stay)
- e) What is the reward function?

Answer

- R((1,1)) = +20
- R((3,1)) = -20
- R((1,3)) = +5
- R(other states) = -1
- 2. Using the MDP built for question 1, consider the following sequences of states that our agent progressed through. For each sequence, compute the *additive reward* and *discounted reward* (with a decay factor γ of 0.5). Take note of the differences between the two approaches.
 - a) (3,3), (2,3), (2,2), (2,1), (1,1)
 - b) $(3,3), (3,2), (2,2), (2,3), \dots$ (repeat)

Answer

- (a) Additive: (-1) + (-1) + (-1) + (-1) + (+20) = 16Discount: $(-1) + 0.5(-1) + 0.5^2(-1) + 0.5^3(-1) + 0.5^4(+20) = -0.625 (10/16)$
- (b) Additive: $(-1) + (-1) + (-1) + (-1) + \dots = -\infty$

Discount:
$$(-1) + 0.5(-1) + 0.5^2(-1) + 0.5^3(-1) + 0.5^4(-1) + \dots = \sum_{t=0}^{\infty} 0.5^t * (-1) = \frac{-1}{1-0.5} = -2$$

- 3. Consider the MDP from question 1, with a decay factor γ of 1.
 - a) Using value iteration, compute the utilities for each state after two iterations. For this question take the value iteration equation:

$$V_{k+1}(s) = \max_{a} \left\{ \sum_{s' \in S} T(s, a, s') \cdot [R(s, a, s') + \gamma V_k(s')] \right\}.$$

Initially, we initialize $V_0(s) = 0$ for all states $s \in S$. Remember terminal states do not get more rewards once reached.

b) Repeat and determine the values after three iterations.

Answer

Here is the detailed working for cell(1,2). Note that, in this case, it is clear that "moving west" is the best action (i.e., maximizes the second term in formula above), but in general one would need to calculate the term for every action and then keep the max one (that is the best action to do!).

Assume $V_0(s) = 0, \forall s \in S$.

Compute utility value V for state (1,2) - row 1, column 2 - using R(s,a,s'):

$$\begin{split} V_{k+1}(s) &= \max_{a} \left\{ \sum_{s' \in S} T(s, a, s') \cdot [R(s, a, s') + \gamma V_k(s')] \right\} \\ V_1((1, 2)) &= 0.7 \cdot [R((1, 2), \text{west}, (1, 1)) + \gamma V_0((1, 1))] + 0.3 \cdot [R((1, 2), \text{north}, (1, 2)) + \gamma V_0((1, 2))] \\ V_1((1, 2)) &= 0.7 \cdot [-1 + 1 \cdot 0] + 0.3 \cdot [-1 + 1 \cdot 0] = -1 \end{split}$$

$$V_1((1,1)) = 1.0[R((1,1),(stay),(1,1)) + \gamma V_0((1,1))]$$
 (terminal state so must "stay") $V_1((1,1)) = 1.0[20 + 1 \cdot 0] = 20$

$$V_2((1,2)) = 0.7 \cdot [R((1,2), \text{west}, (1,1)) + \gamma V_1((1,1))] + 0.3 \cdot [R((1,2), \text{north}, (1,2)) + \gamma V_1((1,2))]$$

$$V_2((1,2)) = 0.7 \cdot [-1 + 1 \cdot 20] + 0.3 \cdot [-1 + 1 \cdot (-1)] = 12.7$$

 $V_2((1,1)) = 20$ (terminal state so cannot take further action)

$$\begin{split} V_3((1,2)) &= 0.7 \cdot [R((1,2), \text{west}, (1,1)) + \gamma V_2((1,1))] + 0.3 \cdot [R((1,2), \text{north}, (1,2)) + \gamma V_2((1,2))] \\ V_2((1,2)) &= 0.7 \cdot [-1 + 1 \cdot 20] + 0.3 \cdot [-1 + 1 \cdot 12.7)] = 16.81 \end{split}$$

4. Consider the utilities after two iterations of value iteration from question 3 (we typically only perform this when values have converged, but for this question, do this after 2 iterations). Using the Maximum Expected Utility principle, find the best action (i.e., policy) for each state.

Answer

For the cell (1,2) we can derive that the best action to do there after 2 iterations is west, that is, $\pi((1,2)) = \text{west}$.

- 5. Consider the following policy for the MDP of question 1, with a decay factor γ of 1:
 - $\pi((2,1)) = \text{west.}$
 - $\pi((2,2)) = \text{west.}$
 - $\pi((2,3)) = \text{north.}$
 - $\pi((3,2)) = \text{north.}$
 - $\pi((3,3)) = \text{east.}$
 - $\pi((1,2)) =$ south.

Evaluate this policy, i.e., compute the utility for each state, for one and two iterations (we typically can solve it as a system of simultaneous equations but for this question we use an iterative approach).

Answer

Here we are given the policy (it is fixed), so the task is to calculate its value at every state, that is, $V_{\pi}(s)$ or $U_{\pi}(s)$. We can do that as follows:

$$V_{k+1}^{\pi}(s) = \sum_{s' \in S} T(s, \pi(s), s') \cdot [R(s, \pi(s), s') + \gamma V_k(s')].$$

Observe we do not take the \max of actions because we know, at every state s, what we should do, namely, $\pi(s)!$

BE CAREFUL, THIS HAS NOT BEEN RE-CHECKED CAREFULLY!

$$\begin{array}{c|ccccc}
 & 1 & 2 & 3 \\
1 & 20 & -0.8 & 5 \\
2 & 4.3 & -2 & -2.2 \\
3 & -20 & -2 & -2
\end{array}$$

6. Will the stay action be used for an optimal policy? Hint: consider the reward function for non-terminal cells.

Answer

If the reward function is positive valued and large in magnitude for non terminating states, and γ is small, then it is possible for an agent to choose to stay in a state.