

Assignment5 – Markov Decision Procedures

Given: May 25 Due: June 4

Problem 5.1 (Markov Decision Processes)

0 pt

1. Give an optimal policy π^* for the following MDP:
 - set of states: $S = \{0, 1, 2, 3, 4, 5\}$ with initial state 0
 - set of actions for $s \in S$: $A(s) = \{-1, 1\}$
 - transition model for $s, s' \in S$ and $a \in A(s)$: $P(s'|s, a)$ is such that
 - $s' = (s + a) \bmod 6$ with probability $2/3$,
 - $s' = (s + 3) \bmod 6$ with probability $1/3$.
 - reward function: $R(5) = 1$ and $R(s) = -0.1$ for $s \in S \setminus \{5\}$
2. State the Bellman equation.
3. Complete the following high-level description of the value iteration algorithm:
 - The algorithm keeps a table $U(s)$ for $s \in S$, that is initialized with

 - In each iteration, it uses the

 - in order to

 - $U(s)$ will converge to the

Solution:

1. $\pi^*(s) = 1$ if $s \in \{3, 4\}$ and $\pi^*(s) = -1$ if $s \in \{0, 1\}$ and arbitrary for $s \in \{2, 5\}$
2. $U(s) = R(s) + \gamma \max_{a \in A(s)} \sum_{s' \in S} U(s')P(s'|s, a)$
3.
 - The algorithm keeps a table $U(s)$ for $s \in S$, that is initialized with arbitrary values, e.g. all 0 or the rewards.
 - In each iteration, it uses the Bellman equation in order to update $U(s)$.
 - $U(s)$ will converge to the expected utility of s .

Problem 5.2 (Bellman Equation)

20 pt

State the Bellman Equation and explain every symbol in the equation and what the equation is used for and how.

Solution:

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

The meaning of the components is as follows:

- $U(s)$: the utility of the state s (long-term, global)
- $R(s)$: the reward at state s (short-term, local)
- $A(s)$: the set of actions available in state s
- $\max_{a \in A(s)}$: take the maximum over all available actions in state s
- $P(s'|s, a)$: the probability that taking action a in state s yields state s'
- $U(s')$: the utility in successor state s'
- $(\sum_{s'} P(s'|s, a) \cdot U(s'))$: the expected utility of action a by summing over all possible successor states

The equation is used to compute the utility of every state. The algorithm uses the equation as an iteration operator that computes new values for every $U(s)$ by evaluating the right hand side for the current values of U . If this leads to a *fixpoint*, a solution for the utilities has been found.

Problem 5.3 (MDP Example)

40 pt

Consider the following world:

+50	-1	-1	-1	...	-1	-1	-1	-1
<i>Start</i>				...				
-50	+1	+1	+1	...	+1	+1	+1	+1

The world is 101 fields wide (i.e., 203 fields in total). In the *Start* state an agent has two possible actions, *Up* and *Down*. It cannot return to *Start* though and the cannot pass gray fields, so after the first move the only possible action is *Right*.

1. Model this world as a Markov Decision Process, i.e., give the components S , s_0 , A , P , and R .
2. For what discount factors γ should the agent choose *Up* and for which *Down*? Compute the utility of each action (i.e., the utility of the successor state) as a function of γ .
3. What is the optimal policy if the upper path is better?

Solution:

1.
 - Set of states $S = \{-101, \dots, 0, \dots, 101\}$.
Note that the set of states can be swapped out arbitrarily against any other set of the same size. The choice made is practical because it allows using 0 as the start state and n as the state $|n|$ steps away from the start.
 - Initial state: $s_0 = 0$.
 - Reward function R : $R(0) = 0$, $R(1) = 50$, $R(-1) = -50$, $R(s) = -1$ for $s \in \{2, \dots, 101\}$, $R(s) = 1$ for $s \in \{-2, \dots, -101\}$
 - Possible actions in each state: $A(0) = \{Up, Down\}$, $A(n) = \{Right\}$ for all $n \in S \setminus \{0\}$
 - Transition model $P(s'|s, a)$: This world is deterministic — the successor state of each action is uniquely determined. Therefore, all probabilities are either 1 or 0.
 - current state $s = 0$: $P(1|0, Up) = 1$, $P(1|0, Down) = 0$, $P(-1|0, Up) = 0$, $P(-1|0, Down) = 1$
and $P(s'|0, a)$ for all $s' \in S \setminus \{-1, 0, 1\}$ and $a \in \{Up, Down\}$
 - current state $s \neq 0$:
 - * $P(s+1|s, Right) = 1$ for $s \in \{1, \dots, 100\}$, $P(101|101, Right) = 1$
 - * $P(s-1|s, Right) = 1$ for $s \in \{-1, \dots, -100\}$, $P(-101|-101, Right) = 1$
2. We have $U(s) = R(s) + \gamma \max_a (\sum_{s'} U(s'))$, since all transitions are deterministic. Then

$$U(1) = 50 + \gamma(-1 + \gamma(-1 + \dots)) = 50 - \sum_{i=1}^{100} \gamma^i$$

$$U(-1) = -50 + \gamma(1 + \gamma(1 + \dots)) = -50 + \sum_{i=1}^{100} \gamma^i$$

and for $i > 0$:

$$U(i) = \sum_{k=1}^i \gamma^k \quad U(-i) = -\sum_{k=1}^i \gamma^k$$

So we need to solve the following equation for γ :

$$\begin{aligned} 50 - \sum_{i=1}^{100} \gamma^i &= -50 + \sum_{i=1}^{100} \gamma^i \\ 50 &= \sum_{i=1}^{100} \gamma^i \end{aligned}$$

We get $\gamma \approx 0.984397669$, and we should go Up if γ is smaller.

3. The optimal policy π^* maps each $s \in S$ to an element of $A(s)$. Because most states have only one action, we immediately have $\pi^*(s) = \text{Right}$ for $s \neq 0$. For $s = 0$, we have $\pi^*(0) = Up$.

Problem 5.4 (Value Iteration for Navigation)

50 pt

Implement value iteration for an agent navigating worlds like the 4x3 world from the lecture notes. The agent has four possible actions: *right*, *up*, *left*, *down*. The probability of actually moving in the intended direction is p and the probability of moving in one of the orthogonal directions is $\frac{1-p}{2}$ respectively. For example, if $p = 0.8$ and the chosen action is *up*, the agent will actually move up with a probability of $p = 0.8$ and will move left and right with a probability of 0.1 each. If the agent ends up moving in a direction that has no free adjacent square, it will remain on its current square instead. For example, if the agent is on square (0, 0) with the action *up*, it will end up on square (0, 1) with a probability of p , on square (1, 0) with a probability of $\frac{1-p}{2}$ and on square (0, 0) with a probability of $\frac{1-p}{2}$.

(0, 2) -0.040 → 0.647	(1, 2) -0.040 → 0.753	(2, 2) -0.040 → 0.855	(3, 2) 1.000 T 1.000
(0, 1) -0.040 ↑ 0.557	(1, 1) W	(2, 1) -0.040 ↑ 0.569	(3, 1) -1.000 T -1.000
(0, 0) -0.040 ↑ 0.465	(1, 0) -0.040 ← 0.386	(2, 0) -0.040 ↑ 0.451	(3, 0) -0.040 ← 0.230

Results for 4x3 world with $p = 0.8$, $\gamma = 0.95$, $\epsilon = 0.001$.

A skeleton implementation with technical instructions can be found at <https://kwarc.info/teaching/AI/resources/AI2/mdp/>. It also allows the visualization of the computed utilities and policy (see figure above): Each square is annotated with the coordinates, the reward, the computed policy and the computed utility. Walls and terminal nodes don't have a policy and are marked with *W* and *T* respectively.

Hint: You will also have to compute a policy based on the utilities obtained from value iteration. For that, you should pick the actions that *maximize* the expected utility. A common mistake is the assumption that the best policy is always to go in the direction of the square with the maximal utility.

Solution: See <https://kwarc.info/teaching/AI/resources/AI2/mdp/>