

ECE433/COS435 Introduction to RL

Assignment 1: MDP

Spring 2024

Fill me in

Your name here.

Due February 12, 2024

Collaborators

Fill me in

Please fill in the names and NetIDs of your collaborators in this section.

Instructions

You should work alone for this this assignment. Writeups should be typesetted in Latex and submitted as PDFs. You can work with whatever tool you like for the code, but **please submit the asked-for snippet and answer in the solutions box as part of your writeup. We will only be grading your write-up.** Make sure to still also attach your notebook/code with your submission.

Question 1. Markov Chain

Let the transition probability matrix of a two-state Markov chain be specified by We have a Markov chain with two states, s_1 and s_2 . The probability of transitioning from s_1 to s_2 is p , and vice versa. We can summarize the transition probabilities in the table shown below.

$$P = \begin{bmatrix} p & 1-p \\ 1-p & p \end{bmatrix},$$

the value of $P_{i,j}$ indicates the probability of transiting from state i to state j , for any $i, j \in [1, 2]$.

Question 1.a

Use the principle of induction¹ to show that

$$P^{(n)} = \begin{bmatrix} \frac{1}{2} + \frac{1}{2}(2p-1)^n & \frac{1}{2} - \frac{1}{2}(2p-1)^n \\ \frac{1}{2} - \frac{1}{2}(2p-1)^n & \frac{1}{2} + \frac{1}{2}(2p-1)^n \end{bmatrix}.$$

Solution

Your solution here...

Question 1.b

Expectation of State Occupancy

- Compute the expected number of times the process is in s_1 after n transitions, starting from s_1 .
- Compute the expected number of times the process is in s_2 after n transitions, starting from s_1 .
- Discuss how the expectations change as n approaches infinity and the implications for the state occupancy for the above two cases.

Solution

Your solution here...

Question 1.c

Probability of First Visit

- Compute the probability that the process visits s_2 for the first time on the k -th transition, given it starts in s_1 .
- How does this probability change as k increases?

Solution

Your solution here...

Question 1.d

Conditional Expectations

- Given that the chain is in s_2 at the n -th step, compute the conditional expectation of the number of visits to s_1 in the next m steps.
- Explore how this expectation varies with different values of p .

¹https://en.wikipedia.org/wiki/Mathematical_induction

Solution

Your solution here...

Question 1.e

Expected rewards. When transiting from one state to another, assume we receive a reward of 1 for reaching s_2 and -1 for reaching s_1 .

- Compute the expected total reward after n transitions (i.e., the summation of rewards), starting from s_1 .

Solution

Your solution here...

Question 2. Secretary problem

Suppose you are hiring one secretary and going to interview the candidates one by one sequentially. After an interview is over, you have to decide whether to hire the current candidate. If you hire the current candidate, the whole process stops. Otherwise, the interview continues, but the candidate will not return and cannot be hired. In other words, you can only hire the current candidate but cannot hire the past candidate. In total, there are n candidates, and you have a strict preference among them. It means you can tell who is better than whom, and no two candidates are equal. When meeting a new candidate, you compare with past candidates, but you do not know the ranking of current candidates in all people. For example, if you have interviewed c_1, c_2, c_3 , then you can rank these three, say $c_1 > c_3 > c_2$, but you do not know the ranking of c_1 among all n candidates. Candidates come in a uniform random order. The objective is to find a policy (when to stop) in order to maximize the probability of hiring the best candidate. The problem is known as the optimal stopping problem. It can be formulated by MDP.

Specifically, denote t as the time after we interview the t -th candidate, $t = 1, \dots, n$. We introduce a state variable $s_t \in \{-1, 0, 1\}$. $s_t = -1$ means the position is filled. $s_t = 0$ means the position is not filled, and the current t -th candidate is not the best candidate so far. $s_t = 1$ means the position is not filled, and the current t -th candidate is the best candidate so far. Initially, we set $s_t = 1$ because after interviewing the first candidate, he/she must be the best among past candidates. Suppose we do not hire anyone and keep interviewing.

Question 2.a

Compute the following probabilities: $P_t(s, 1) = \mathbb{P}(s_{t+1} = 1 | s_t = s)$ and $P_t(s, 0) = \mathbb{P}(s_{t+1} = 0 | s_t = s)$. (Hint: you may use t in the expressions.)

Solution

Your solution here...

Question 2.b

At time t , our action is either to hire ($a_t = 1$) or to continue interviewing ($a_t = 0$). Show the value of $P(s_{t+1} = s' | s_t = s, a_t = a)$ for $\forall s, s' \in \{-1, 0, 1\}, a_t \in \{0, 1\}$. (Hint: you may keep $P_t(s, s')$ in the expressions.)

Solution

Your solution here...

Question 3. Grid World Example

In this exercise, you will work with a simple reinforcement learning environment called "Gridworld." Gridworld is a 4x4 grid where an agent moves to reach a goal state. The agent can take four actions at each state (up, down, left, right) and receive a reward for each action. Moving into a wall (the edge of the grid) keeps the agent in its current state.

Grid Layout:

- The grid is a 4x4 matrix.
- Start state (S): Top left cell (0,0).
- Goal state (G): Bottom right cell (3,3).

The agent receives a reward of -1 for each action until it reaches the goal state.

Question 3.a

Formulate the problem as a Markov Decision Process (MDP). Define the states, actions, transition probabilities (assume deterministic transitions), rewards, and policy.

Solution

Your solution here...

Question 3.b

Define the policy.

Solution

Your solution here...

Question 3.c

How many unique (deterministic) policies are there in total?

Solution

Your solution here...