

CSE 518 - Artificial Intelligence Homework

Instructor: Shashi Prabh

Chapter 17. Making Complex Decisions, MDP

17.1 For the 4×3 world shown in Figure 1, calculate which squares can be reached from (1,1) by the action sequence $[Up, Up, Right, Right, Right]$ and with what probabilities. Explain how this computation is related to the prediction task (see Section 14.2) for a hidden Markov model.

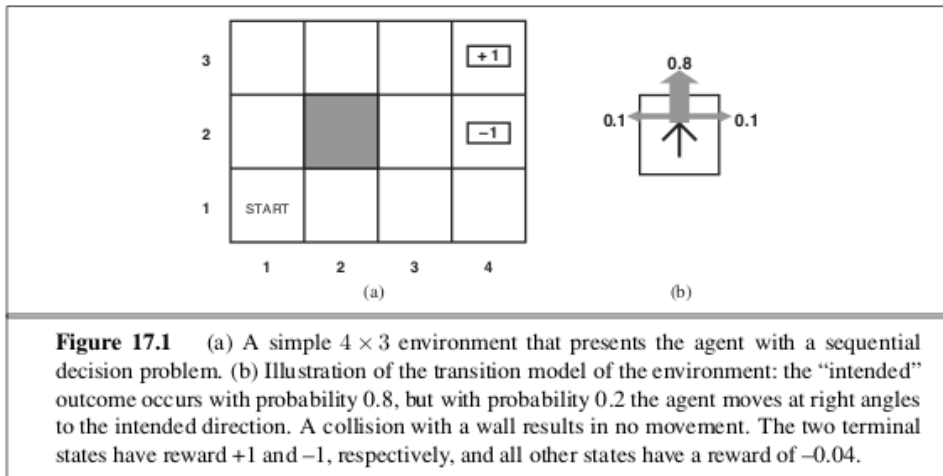


Figure 1: Exercise 17.1

17.2 Select a specific member of the set of policies that are optimal for $R(s) > 0$ as shown in Figure 2(b), and calculate the fraction of time the agent spends in each state, in the limit, if the policy is executed forever.

17.3 Suppose that we define the utility of a state sequence to be the *maximum* reward obtained in any state in the sequence. Show that this utility function does not result in stationary preferences between state sequences. Is it still possible to define a utility function on states such that MEU decision making gives optimal behavior?

17.4 Sometimes MDPs are formulated with a reward function $R(s, a)$ that depends on the action taken or with a reward function $R(s, a, s')$ that also depends on the outcome state.

- a. Write the Bellman equations for these formulations.
- b. Show how an MDP with reward function $R(s, a, s')$ can be transformed into a different MDP with reward function $R(s, a)$, such that optimal policies in the new MDP correspond exactly to optimal policies in the original MDP.

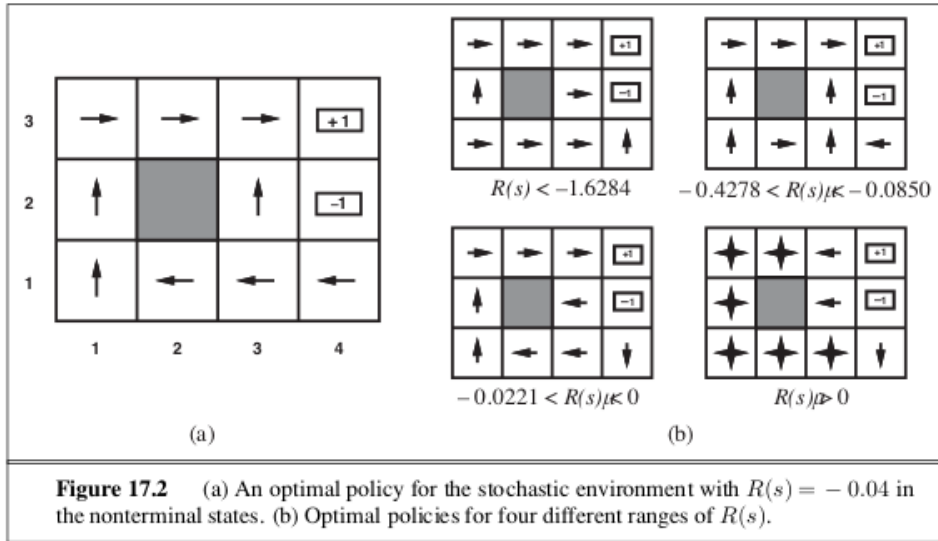


Figure 2: Exercise 17.2

c. Now do the same to convert MDPs with $R(s, a)$ into MDPs with $R(s)$.