# CSE 518 - Artificial Intelligence
# Homework

### Instructor: Shashi Prabh

## Chapter 17. Making Complex Decisions, MDP

**17.1**  For the $4 \times 3$ 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 17.1**    (a) A simple $4 \times 3$ environment that presents the agent with a sequential decision problem. (b) Illustration of the transition model of the environment: the "intended" outcome occurs with probability 0.8, but with probability 0.2 the agent moves at right angles to the intended direction. A collision with a wall results in no movement. The two terminal states have reward +1 and −1, respectively, and all other states have a reward of −0.04.
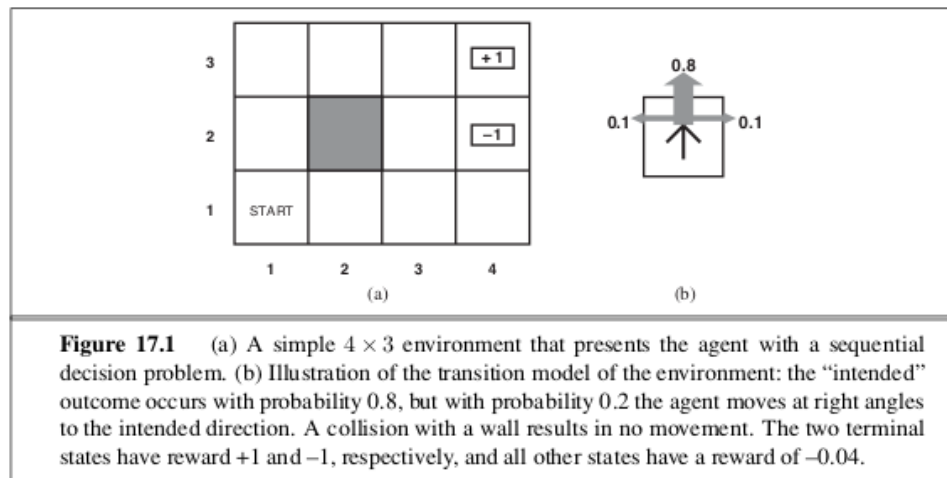
Figure 1: Exercise 1

**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.
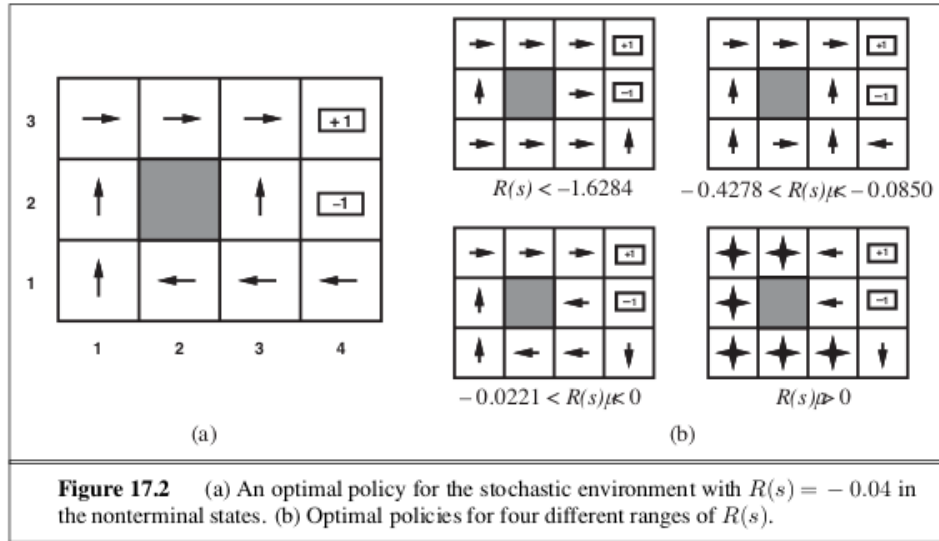
**Figure 17.2** (a) An optimal policy for the stochastic environment with $R(s) = -0.04$ in the nonterminal states. (b) Optimal policies for four different ranges of $R(s)$.

Figure 2: Exercise 1

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