tionary game theory (Smith, 1982; Weibull, 1995) looks at strategy drift over time: if your opponent's strategy is changing, how should you react? Textbooks on game theory from an economics point of view include those by Myerson (1991), Fudenberg and Tirole (1991), Osborne (2004), and Osborne and Rubinstein (1994); Mailath and Samuelson (2006) concentrate on repeated games. From an AI perspective we have Nisan *et al.* (2007), Leyton-Brown and Shoham (2008), and Shoham and Leyton-Brown (2009).

The 2007 Nobel Memorial Prize in Economics went to Hurwicz, Maskin, and Myerson "for having laid the foundations of mechanism design theory" (Hurwicz, 1973). The tragedy of the commons, a motivating problem for the field, was presented by Hardin (1968). The revelation principle is due to Myerson (1986), and the revenue equivalence theorem was developed independently by Myerson (1981) and Riley and Samuelson (1981). Two economists, Milgrom (1997) and Klemperer (2002), write about the multibillion-dollar spectrum auctions they were involved in.

Mechanism design is used in multiagent planning (Hunsberger and Grosz, 2000; Stone *et al.*, 2009) and scheduling (Rassenti *et al.*, 1982). Varian (1995) gives a brief overview with connections to the computer science literature, and Rosenschein and Zlotkin (1994) present a book-length treatment with applications to distributed AI. Related work on distributed AI also goes under other names, including collective intelligence (Tumer and Wolpert, 2000; Segaran, 2007) and market-based control (Clearwater, 1996). Since 2001 there has been an annual Trading Agents Competition (TAC), in which agents try to make the best profit on a series of auctions (Wellman *et al.*, 2001; Arunachalam and Sadeh, 2005). Papers on computational issues in auctions often appear in the ACM Conferences on Electronic Commerce.

## **EXERCISES**

**17.1** For the  $4 \times 3$  world shown in Figure 17.1, 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 15.2.1) for a hidden Markov model.

**17.2** Select a specific member of the set of policies that are optimal for R(s) > 0 as shown in Figure 17.2(b), and calculate the fraction of time the agent spends in each state, in the limit, if the policy is executed forever. (*Hint*: Construct the state-to-state transition probability matrix corresponding to the policy and see Exercise 15.2.)

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

**17.5** For the environment shown in Figure 17.1, find all the threshold values for R(s) such that the optimal policy changes when the threshold is crossed. You will need a way to calculate the optimal policy and its value for fixed R(s). (*Hint*: Prove that the value of any fixed policy varies linearly with R(s).)

- 17.6 Equation (17.7) on page 654 states that the Bellman operator is a contraction.
  - **a**. Show that, for any functions *f* and *g*,

 $\left|\max_{a} f(a) - \max_{a} g(a)\right| \le \max_{a} \left|f(a) - g(a)\right|.$ 

**b**. Write out an expression for  $|(BU_i - BU'_i)(s)|$  and then apply the result from (a) to complete the proof that the Bellman operator is a contraction.

**17.7** This exercise considers two-player MDPs that correspond to zero-sum, turn-taking games like those in Chapter 5. Let the players be A and B, and let R(s) be the reward for player A in state s. (The reward for B is always equal and opposite.)

- **a**. Let  $U_A(s)$  be the utility of state *s* when it is *A*'s turn to move in *s*, and let  $U_B(s)$  be the utility of state *s* when it is *B*'s turn to move in *s*. All rewards and utilities are calculated from *A*'s point of view (just as in a minimax game tree). Write down Bellman equations defining  $U_A(s)$  and  $U_B(s)$ .
- **b**. Explain how to do two-player value iteration with these equations, and define a suitable termination criterion.
- c. Consider the game described in Figure 5.17 on page 197. Draw the state space (rather than the game tree), showing the moves by A as solid lines and moves by B as dashed lines. Mark each state with R(s). You will find it helpful to arrange the states  $(s_A, s_B)$  on a two-dimensional grid, using  $s_A$  and  $s_B$  as "coordinates."
- d. Now apply two-player value iteration to solve this game, and derive the optimal policy.

**17.8** Consider the  $3 \times 3$  world shown in Figure 17.14(a). The transition model is the same as in the  $4 \times 3$  Figure 17.1: 80% of the time the agent goes in the direction it selects; the rest of the time it moves at right angles to the intended direction.

Implement value iteration for this world for each value of r below. Use discounted rewards with a discount factor of 0.99. Show the policy obtained in each case. Explain intuitively why the value of r leads to each policy.

- a. r = 100
  b. r = -3
  c. r = 0
- **d**. r = +3

r	-1	+10		+50	-1	-1	-1	•••	-1	-1	-1	-1	
-1	-1	-1		Start				•••					
-1	-1	-1		-50	+1	+1	+1	•••	+1	+1	+1	+1	
	(a)			(b)									
<b>Figure 17.14</b> (a) $3 \times 3$ world for Exercise 17.8. The reward for each state is indicated. The upper right square is a terminal state. (b) $101 \times 3$ world for Exercise 17.9 (omitting 93 identical columns in the middle). The start state has reward 0.													

**17.9** Consider the  $101 \times 3$  world shown in Figure 17.14(b). In the start state the agent has a choice of two deterministic actions, *Up* or *Down*, but in the other states the agent has one deterministic action, *Right*. Assuming a discounted reward function, for what values of the discount  $\gamma$  should the agent choose *Up* and for which *Down*? Compute the utility of each action as a function of  $\gamma$ . (Note that this simple example actually reflects many real-world situations in which one must weigh the value of an immediate action versus the potential continual long-term consequences, such as choosing to dump pollutants into a lake.)

**17.10** Consider an undiscounted MDP having three states, (1, 2, 3), with rewards -1, -2, 0, respectively. State 3 is a terminal state. In states 1 and 2 there are two possible actions: *a* and *b*. The transition model is as follows:

- In state 1, action *a* moves the agent to state 2 with probability 0.8 and makes the agent stay put with probability 0.2.
- In state 2, action *a* moves the agent to state 1 with probability 0.8 and makes the agent stay put with probability 0.2.
- In either state 1 or state 2, action *b* moves the agent to state 3 with probability 0.1 and makes the agent stay put with probability 0.9.

Answer the following questions:

- **a**. What can be determined *qualitatively* about the optimal policy in states 1 and 2?
- **b**. Apply policy iteration, showing each step in full, to determine the optimal policy and the values of states 1 and 2. Assume that the initial policy has action *b* in both states.
- **c**. What happens to policy iteration if the initial policy has action *a* in both states? Does discounting help? Does the optimal policy depend on the discount factor?
- **17.11** Consider the  $4 \times 3$  world shown in Figure 17.1.
  - **a**. Implement an environment simulator for this environment, such that the specific geography of the environment is easily altered. Some code for doing this is already in the online code repository.