

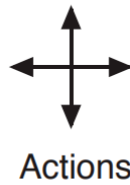
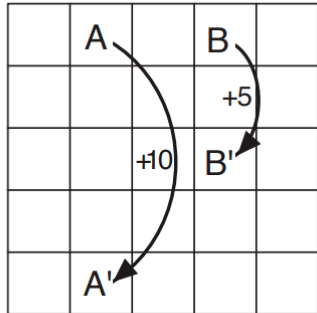
1. (*Exercise 3.12 in 2nd ed.*) Recall that the value $v_\pi(s)$ for state s when following policy π is the expected total reward (or discounted reward) the agent would receive when starting from state s and executing policy π . How can we write $v_\pi(s)$ in terms of the action-values $q_\pi(s, a)$?

2. In this question, you will take a word specification of an MDP, and write the formal terms and determine the optimal policy. Suppose you have a problem with two actions. The agent always starts in the same state, s_0 . From this state, if it takes action 1 it transitions to a new state s_1 and receives reward 10; if it takes action 2 it transitions to a new state s_2 and receives reward 5. From s_1 if it takes action 1 it receives a reward of 5 and terminates; if it takes action 2 it receives a reward of 10 and terminates. From s_2 if it takes action 1 it receives a reward of 10 and terminates; if it takes action 2 it receives a reward of 5 and terminates. Assume the agent cares equally about long term reward as about immediate reward.
 - (a) Draw the MDP for this problem. Is it an episodic or continuing problem? What is γ ?
 - (b) Assume the policy is $\pi(a = 1|s_i) = 0.3$ for all $s_i \in \{s_0, s_1, s_2\}$. What is $\pi(a = 2|s_i)$? And what is the value function for this policy? In other words, find $v_\pi(s)$ for all three states.
 - (c) What is the optimal policy in this environment?

Worksheet 4

CMPUT 397
September 27, 2019

3. Consider the gridworld and value function in the figure below. Using your knowledge of the transition dynamics and the values (numbers in each grid cell), write down the policy corresponding to taking the greedy action with respect to the values in each state. Create a grid with the same dimension as the figure and draw an arrow in each square denoting the greedy action.



3.3	8.8	4.4	5.3	1.5
1.5	3.0	2.3	1.9	0.5
0.1	0.7	0.7	0.4	-0.4
-1.0	-0.4	-0.4	-0.6	-1.2
-1.9	-1.3	-1.2	-1.4	-2.0

Worksheet 4

4. Consider the continuing MDP shown on the bottom. The only decision to be made is that in the top state, where two actions are available, left and right. The numbers show the rewards that are received deterministically after each action.
- (a) List and describe all the possible policies in this MDP.
 - (b) Is the following policy valid for this MDP (i.e. does it fit our definition of a policy): Choose *left* for five steps, then *right* for five steps, then *left* for five steps, and so on? Explain your answer.
 - (c) What policy is optimal if $\gamma = 0$? If $\gamma = 0.9$? If $\gamma = 0.5$?
 - (d) For each possible policy, what is the value of state s ? Write down the numeric value to two decimal places. *Hint*: write down the return under each policy starting in state s (don't forget γ). Simplify the infinite sum, using the fact that many rewards are zero. Then plug in the rewards and γ and compute the number.

