



Northeastern University  
CS 7180 – Special Topics in AI (Reinforcement Learning)  
Fall 2018, Robert Platt

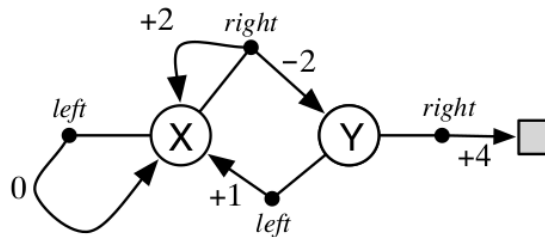
## MDPs Assignment

Name: \_\_\_\_\_

Problem	Points
TRAJECTORIES, RETURNS, AND VALUES.	/15
MAZE RUNNING.	/5
BROKEN VISION SYSTEM.	/5
UNDERSTANDING VALUE FUNCTIONS	/5
CALCULATING TOTAL RETURNS.	/5
BELLMAN, PART 1	/15
BELLMAN, PART 2	/5
ADDING A CONSTANT TO REWARDS.	/5
BELLMAN, PART 3	/10
BELLMAN, PART 4	/20
<b>Total</b>	<b>/90</b>

### Instructions

- Don't cheat!

**(15 pts.)** TRAJECTORIES, RETURNS, AND VALUES.

Consider the MDP above, in which there are two states,  $X$  and  $Y$ , two actions, *right* and *left*, and the deterministic rewards on each transition are as indicated by the numbers. Note that if action *right* is taken in state  $X$ , then the transition may be either to  $X$  with a reward of  $+2$  or to  $Y$  with a reward of  $-2$ . These two possibilities occur with probabilities  $2/3$  (for the transition to  $X$ ) and  $1/3$  (for the transition to state  $Y$ ).

Consider two deterministic policies,  $\pi_1$  and  $\pi_2$ :

$$\begin{aligned}\pi_1(X) &= \textit{left} \\ \pi_1(Y) &= \textit{right}\end{aligned}$$

$$\begin{aligned}\pi_2(X) &= \textit{right} \\ \pi_2(Y) &= \textit{right}\end{aligned}$$

- (2 pts) Show a typical trajectory (sequence of states, actions and rewards) from  $X$  for policy  $\pi_1$ :
- (2 pts) Show a typical trajectory (sequence of states, actions and rewards) from  $X$  for policy  $\pi_2$ :
- (2 pts) Assuming the discount-rate parameter is  $\gamma = 0.5$ , what is the return from the initial state for the second trajectory?
- (2 pts) Assuming  $\gamma = 0.5$ , what is the value of state  $Y$  under policy  $\pi_1$ ?
- (2 pts) Assuming  $\gamma = 0.5$ , what is the action-value of  $(Y, \textit{left})$  under policy  $\pi_1$ ?
- (5 pts) Assuming  $\gamma = 0.5$ , what is the value of state  $(X)$  under policy  $\pi_2$ ?

**(5 pts.) MAZE RUNNING.**

Exercise 3.7 in the SB textbook.

**(5 pts.)** BROKEN VISION SYSTEM.

Imagine that you are a vision system. When you are first turned on for the day, an image floods into your camera. You can see lots of things, but not all things. You can't see objects that are occluded, and of course you can't see objects that are behind you. After seeing that first scene, do you have access to the Markov state of the environment? Suppose your camera was broken that day and you received no images at all, all day. Would you have access to the Markov state then?

**(5 pts.)** UNDERSTANDING VALUE FUNCTIONS

Suppose we have two problems with the same state and action spaces. Let the optimal action-value functions of the two problems be denoted  $Q^*$  and  $\bar{Q}^*$ , and suppose it happens to be the case that  $\bar{Q}^*(s, a) = Q^*(s, a) + f(s), \forall s, a$  for some function  $f : S \rightarrow \mathbb{R}$ . What is the relationship between the optimal policies  $\pi^*$  and  $\bar{\pi}^*$  for the two problems?

**(5 pts.)** CALCULATING TOTAL RETURNS.

Exercise 3.9 in the SB textbook.

**(15 pts.)** BELLMAN, PART 1

Exercise 3.11, 3.13 in the SB textbook.

**(5 pts.)** BELLMAN, PART 2

Exercise 3.14 in the SB textbook.



**(5 pts.)** ADDING A CONSTANT TO REWARDS.

Exercise 3.15 in the SB textbook.

**(10 pts.)** BELLMAN, PART 3

Exercise 3.19 in the SB textbook.

**(20 pts.)** BELLMAN, PART 4

Exercise 3.26, 3.27, 3.28 in the SB textbook.