# CS221 Fall 2018 Homework 4

SUNet ID:   05794739

Name:   Luis Perez

Collaborators:

By turning in this assignment, I agree by the Stanford honor code and declare that all of this is my own work.

## Problem 1

(a) We give the value for each iteration. We note that $V_{\text{opt}}^0(s) = 0$ to start out. We also note that since for $s_t \in \{-2, 2\}$ we are at a terminal state, we'll have $V_{\text{opt}}(s_t) = 0$ for all iterations.

  (a) After iteration 0, we'll have:

$$V_{\text{opt}}^0(-1) = 0$$
$$V_{\text{opt}}^0(0) = 0$$
$$V_{\text{opt}}^0(1) = 0$$

  (b) After the first iteration, we'll have the following values:

$$V_{\text{opt}}^1(-1) = \max_{a \in \{-1,1\}} \{0.8[20 + V_{\text{opt}}^0(-2)] + 0.2[-5 + V_{\text{opt}}^0(0)], 0.7[20 + V_{\text{opt}}^0(-2)] + 0.3[-5 V_{\text{opt}}^0(0)]\}$$
$$= 15$$
$$V_{\text{opt}}^1(0) = \max_{a \in \{-1,1\}} \{0.8[-5 + V_{\text{opt}}^0(-1)] + 0.2[-5 + V_{\text{opt}}^0(1)], 0.7[-5 + V_{\text{opt}}^0(-1)] + 0.3[-5 + V_{\text{opt}}^0(1)]\}$$
$$= -5$$
$$V_{\text{opt}}^1(1) = \max_{a \in \{-1,1\}} \{0.8[-5 + V_{\text{opt}}^0(0)] + 0.2[100 + V_{\text{opt}}^0(2)], 0.7[-5 + V_{\text{opt}}^0(0)] + 0.3[100 + V_{\text{opt}}^0(2)]\}$$
$$= 26.5$$

  (c) Finally, after the second iteration, we'll have:

$$V_{\text{opt}}^2(-1) = \max_{a \in \{-1,1\}} \{0.8[20 + V_{\text{opt}}^1(-2)] + 0.2[-5 + V_{\text{opt}}^1(0)], 0.7[20 + V_{\text{opt}}^1(-2)] + 0.3[-5 + V_{\text{opt}}^1(0)]\}$$
$$= 14$$
$$V_{\text{opt}}^2(0) = \max_{a \in \{-1,1\}} \{0.8[-5 + V_{\text{opt}}^1(-1)] + 0.2[-5 + V_{\text{opt}}^1(1)], 0.7[-5 + V_{\text{opt}}^1(-1)] + 0.3[-5 + V_{\text{opt}}^1(1)]\}$$
$$= 13.45$$
$$V_{\text{opt}}^2(1) = \max_{a \in \{-1,1\}} \{0.8[-5 + V_{\text{opt}}^1(0)] + 0.2[100 + V_{\text{opt}}^1(2)], 0.7[-5 + V_{\text{opt}}^1(0)] + 0.3[100 + V_{\text{opt}}^1(2)]\}$$
$$= 23$$

(b) We interpret this question as asking for the resulting optimal policy for non-terminal states after two iterations. In that case, we have:

$$\pi_{\text{opt}}^2(-1) = \arg\max_{a \in \{-1,1\}} \{0.8[20 + V_{\text{opt}}^1(-2)] + 0.2[-5 + V_{\text{opt}}^1(0)], 0.7[20 + V_{\text{opt}}^1(-2)] + 0.3[-5 + V_{\text{opt}}^1(0)]\}$$

$$= -1$$

$$\pi_{\text{opt}}^2(0) = \arg\max_{a \in \{-1,1\}} \{0.8[-5 + V_{\text{opt}}^1(-1)] + 0.2[-5 + V_{\text{opt}}^1(1)], 0.7[-5 + V_{\text{opt}}^1(-1)] + 0.3[-5 + V_{\text{opt}}^1(1)]\}$$

$$= 1$$

$$\pi_{\text{opt}}^2(1) = \arg\max_{a \in \{-1,1\}} \{0.8[-5 + V_{\text{opt}}^1(0)] + 0.2[100 + V_{\text{opt}}^1(2)], 0.7[-5 + V_{\text{opt}}^1(0)] + 0.3[100 + V_{\text{opt}}^1(2)]\}$$

$$= 1$$