

Question 1

a) - Explore then commit: The total number of pulls, N , is the sum of wins and losses which is 20. Since this is less than $k=30$, we are still in the random exploration phase, so the probabilities are

$$\begin{array}{l} 1: \frac{1}{3} \\ 2: \frac{1}{3} \\ 3: \frac{1}{3} \end{array}$$

- ϵ -greedy

The best arm is 2, so it will be pulled w.p. $1-\epsilon+\frac{\epsilon}{3}$

$$\begin{array}{l} 1: 0.033 \\ 2: 0.933 \\ 3: 0.033 \end{array}$$

- UCB, $c=2$

$$N=20, \log N \approx 3$$

$$\mu_a + c \sqrt{\frac{\log N}{N(a)}} = \begin{bmatrix} 0.5 \\ 0.8 \\ 0.5 \end{bmatrix} + 2 \begin{bmatrix} \sqrt{\frac{3}{2}} \\ \sqrt{\frac{3}{10}} \\ \sqrt{\frac{3}{8}} \end{bmatrix} = \begin{bmatrix} 2.9 \\ 1.9 \\ 1.7 \end{bmatrix}$$

Antenna 1 has the highest UCB value, so it will be used w.p. 1.

$$1: 1.0, 2: 0, 3: 0.$$

b) There are a total of $N=10$ tries of antennas 1 and 3; $N=10, \log N = 2.3$

$$\mu_a + c \sqrt{\frac{\log N}{N(a)}} = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} + 2 \begin{bmatrix} \sqrt{\frac{2.3}{2}} \\ \sqrt{\frac{2.3}{8}} \end{bmatrix} = \begin{bmatrix} 2.64 \\ 1.57 \end{bmatrix}$$

Antenna 1 will be used.

Question 2

a) Q-learning update

$$Q(1,1) \leftarrow Q(1,1) + \alpha (r + \gamma \max_{a'} Q(s',a') - Q(1,1))$$
$$8 + 0.1(1 + 0.9 \cdot 8 - 8)$$

$$\boxed{\begin{aligned} Q(1,1) &= 8.02 \\ Q(1,2) &= 6 \quad (\text{unchanged}) \end{aligned}}$$

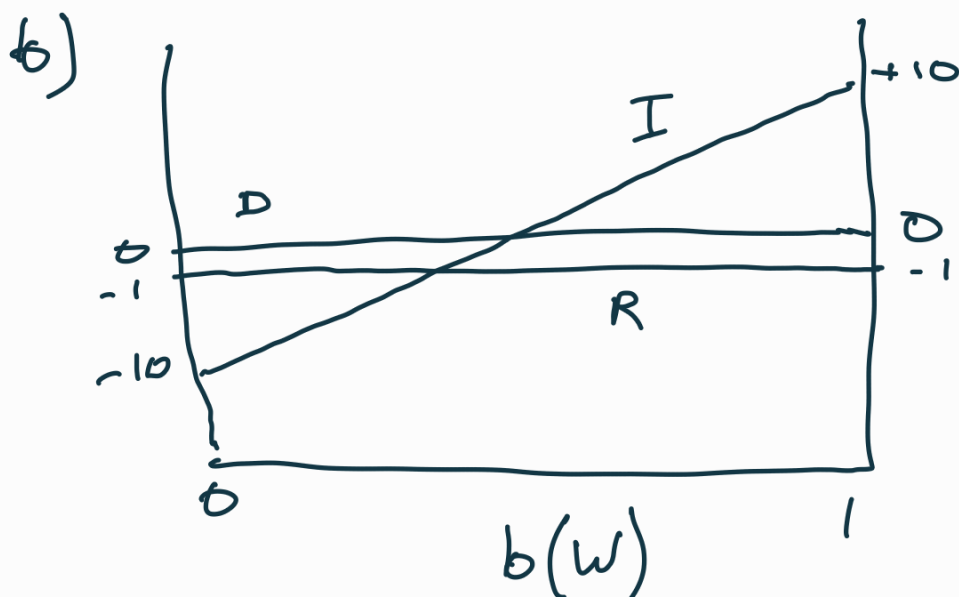
b) SARSA update

$$Q(1,1) \leftarrow Q(1,1) + \alpha (r + \gamma Q(1,2) - Q(1,1))$$
$$8 + 0.1(1 + 0.9 \cdot 6 - 8)$$

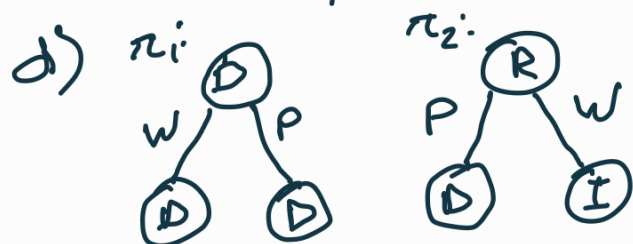
$$\boxed{\begin{aligned} Q(1,1) &= 7.84 \\ Q(1,2) &= 6 \quad (\text{unchanged}) \end{aligned}}$$

Question 3

a) $\alpha_D = [0, 0]$
 $\alpha_I = [10, -10]$
 $\alpha_R = [-1, -1]$



c) R will never be chosen because its α -vector is dominated over the entire belief space.



e) $U^\pi(s) = R(s, \pi(1)) + \gamma \left[\sum_{s'} T(s'|s, \pi(1)) \sum_0 Z(0|\pi(1), s') U^{\pi(0)}(s') \right]$

For π_1 = always decline, $s' = s$, and $\pi(0) = 0$

~~$U^{\pi_1}(s) = R(s, D) + \gamma \left[\sum_0 Z(0|D, s) U^D(0) \right]$~~

$= 0 \therefore \boxed{\alpha_{\pi_1} = [0, 0]}$

For π_2 :
~~can ignore this sum since $s = s'$~~

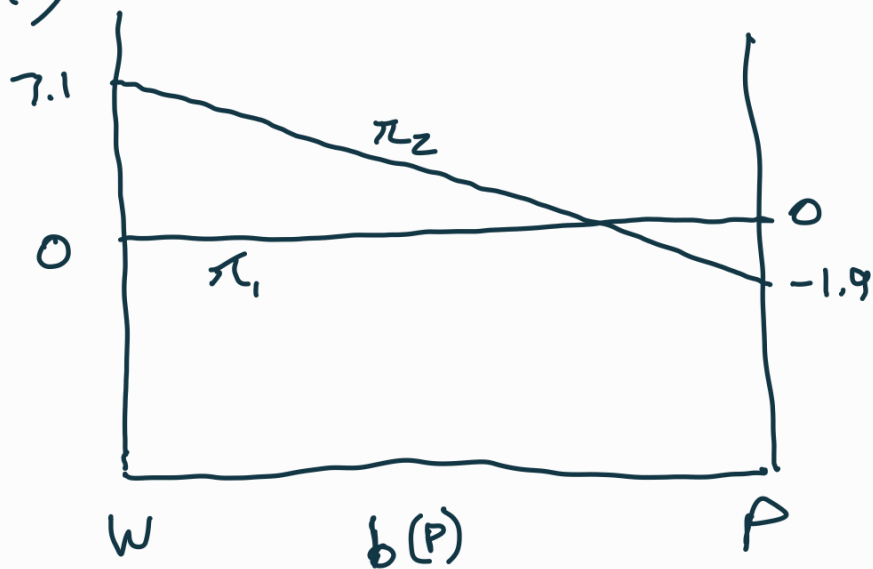
~~$U^{\pi_2}(s) = R(s, R) + \gamma \left[\sum_{s'} T(s'|s, R) (Z(P|R, s) U^D(s) + Z(W|R, s) U^I(s)) \right]$~~

$U^\pi(p) = -1 + \gamma [0.1 \cdot -10] = -1.9$

$U^\pi(w) = -1 + \gamma [0.9 \cdot 10] = 7.1$

$\boxed{\alpha_{\pi_2} = [7.1, -1.9]}$

f)



g) $b = [0.5, 0.5]$

$$\pi^* = \arg \max_{\pi \in \Gamma} \alpha_{\pi} \cdot b$$

$$\pi_1: [6, 0] \cdot [0.5, 0.5] = 0$$

$$\pi_2: [7.1, -1.9] \cdot [0.5, 0.5] = 2.6$$

$$\pi^* = \pi_2$$

$\pi^*(\cdot) = \boxed{R}$ would be selected.

Question 4

Q-learning is more appropriate because it is an off-policy algorithm so it can learn from data collected with any policy.

Policy Gradient is on-policy, so it can only learn from data consistent with the current policy.

Thus Q-learning is easier to use.

(Note: it may be possible to adapt policy gradient or the data to be used together, but this will not necessarily be easy)

Question 5

This function is less expressive than a neural network because it can only approximate linear functions of x . Neural networks have nonlinearities that the inputs are passed through so that they can approximate nonlinear functions.