

UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN  
CS440/ECE448 Artificial Intelligence  
**Exam 3**  
Spring 2022

May 13, 2022

---

**Your Name:** \_\_\_\_\_

**Your NetID:** \_\_\_\_\_

---

**Instructions**

- Please write your NetID on the top of every page.
- This is a **CLOSED BOOK** exam. You will be permitted to bring two 8.5x11 pages of handwritten notes (front & back).
- Calculators are not permitted. You need not simplify explicit numerical expressions.
- There is scratch paper available on the last page of the exam.

**Possibly Useful Formulas**

$$\text{Naïve Bayes: } P(X = x, Y = y) \approx P(Y = y) \prod_{i=1}^n P(W = w_i | Y = y)$$

$$\text{Laplace Smoothing: } P(w) = \frac{\text{Count}(w) + k}{\sum_w \text{Count}(w) + k(1 + \sum_w 1)}$$

$$\text{Perceptron: } \vec{w}_y = \vec{w}_y + \eta \vec{x}, \quad \vec{w}_{f(\vec{x})} = \vec{w}_{f(\vec{x})} - \eta \vec{x}$$

$$\text{Logistic Regression: } \nabla_{\vec{w}_c} \mathcal{L}_i = \nabla_{\vec{w}_c} \left( -\ln \frac{e^{\vec{w}_{c_i}^T \vec{x}_i}}{\sum_k e^{\vec{w}_k^T \vec{x}_i}} \right) = \left( \frac{e^{\vec{w}_{c_i}^T \vec{x}_i}}{\sum_k e^{\vec{w}_k^T \vec{x}_i}} - y_{i,c} \right) \vec{x}_i$$

$$\text{Neural Net: } \xi_j^{(l)} = b_j^{(l)} + \sum_k w_{j,k}^{(l)} h_k^{(l-1)}, \quad h_j^{(l)} = g^{(l)}(\xi_j^{(l)})$$

$$\text{Back-Propagation: } \frac{\partial \mathcal{L}}{\partial h_k^{(l-1)}} = \sum_j \frac{\partial \mathcal{L}}{\partial h_j^{(l)}} \frac{\partial h_j^{(l)}}{\partial h_k^{(l-1)}}$$

$$\text{Admissible Heuristic: } h(n) \leq d(n)$$

$$\text{Consistent Heuristic: } h(m) - h(n) \leq d(m) - d(n) \text{ if } d(m) - d(n) \geq 0$$

$$\text{Conditional Probability: } P(A, B, C) = P(A)P(B|A)P(C|A, B)$$

$$\text{Marginal Probability: } P(A, C) = \sum_b P(A, B = b, C)$$

$$\text{Viterbi Algorithm: } v_{j,t} = \max_i v_{i,t-1} a_{j,i} b_{j,k}$$

$$\text{Mixed Nash Equilibrium: } (1-p)w + px = (1-p)y + pz$$

$$\text{Expectiminimax: } U(s) = \max_a \sum_{s'} P(s'|s, a) U(s') \text{ or } U(s) = \min_a \sum_{s'} P(s'|s, a) U(s')$$

$$\text{Bellman's Equation: } U(s) = R(s) + \gamma \max_a \sum_{s'} P(s'|s, a) U(s')$$

$$\text{Value Iteration: } U_{t+1}(s) = R(s) + \gamma \max_a \sum_{s'} P(s'|s, a) U_t(s')$$

$$\text{Policy Evaluation: } U^\pi(s) = R(s) + \gamma \sum_{s'} P(s'|s, a) U^\pi(s')$$

$$\text{Policy Improvement: } \pi_{t+1}(s) = \arg \max_a \sum_{s'} P(s'|s, a) U^{\pi_t}(s')$$

$$\text{Q-Learning: } Q_{t+1}(s_t, a_t) = Q_t(s_t, a_t) + \alpha (Q_{local}(s_t, a_t) - Q_t(s_t, a_t))$$

$$\text{TD: } Q_{local}(s_t, a_t) = R_t(s_t) + \gamma \max_{a \in \mathcal{A}(s_{t+1})} Q_t(s_{t+1}, a')$$

$$\text{SARSA: } Q_{local}(s_t, a_t) = R_t(s_t) + \gamma Q_t(s_{t+1}, a_{t+1})$$

$$\text{Deep Q: } \mathcal{L} = \frac{1}{2} E [Q_t(\vec{s}_t, a_t) - Q_{local}(\vec{s}_t, a_t)]$$

$$\text{Imitation Learning: } \mathcal{L} = -\log \pi_{a_t}(\vec{s}_t)$$

$$\text{Actor-Critic: } \mathcal{L} = -\sum_a \pi_a(\vec{s}) Q_t(\vec{s}, a)$$

**Question 1** (7 points)

The planet Illus is home to two species of giant insects, called buzzers and biters. Only biters are dangerous; fortunately, biters are ten times less common than buzzers. A survey team has determined that buzzers are beige-colored with probability  $p$ , otherwise they are gray; biters are beige-colored with probability  $q$ , otherwise they are gray. The settlers of First Landing have discovered a nest of insects; they have observed ten, of which eight were beige and two were gray. Under what condition should they conclude that the insects are most probably biters? Your answer should be an inequality in terms of  $p$  and  $q$ .

**Solution:**

$$\left(\frac{1}{11}\right) q^8 (1-q)^2 > \left(\frac{10}{11}\right) p^8 (1-p)^2$$

**Question 2** (7 points)

Consider the following neural net:

$$\begin{bmatrix} h_1 \\ h_2 \end{bmatrix} = \text{ReLU} \left( \begin{bmatrix} w_{1,1} & w_{1,2} \\ w_{2,1} & w_{2,2} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right)$$

$$f = u_1 h_1 + u_2 h_2$$

$$\mathcal{L} = \frac{1}{2}(f - y)^2$$

Write  $\frac{\partial \mathcal{L}}{\partial w_{2,1}}$  in terms of any of the variables listed above. Be sure to consider separately, if necessary, the cases  $h_1 = 0$  vs.  $h_1 > 0$  and  $h_2 = 0$  vs.  $h_2 > 0$ .

**Solution:**

$$\frac{\partial \mathcal{L}}{\partial w_{2,1}} = \begin{cases} (f - y)u_2 x_1 & h_2 > 0 \\ 0 & h_2 = 0 \end{cases}$$

**Question 3** (7 points)

A rogue AI named Richard Feynman is loose on the internet, and is trying to find the safest path from Sandia to NCSA via intervening network routers. The total cost of the path from router  $i$  to NCSA is  $d_i = -\ln q_i$ , where  $q_i$  is the probability that Richard can successfully follow the entire path without getting caught at any of the routers between router  $i$  and NCSA. The probability of success is upper-bounded by  $q_i \leq e^{-m_i}$ , where  $m_i$  is the physical distance (in miles) between router  $i$  and NCSA. Consider an A\* search using the heuristic  $h(i) = m_i$ .

- (a) Is this information sufficient to prove that the heuristic is admissible? If so, prove it. If not, give a counter-example.

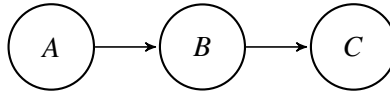
**Solution:** Yes.  $q_i \leq e^{-m_i}$ , so  $d_i = -\ln q_i \geq m_i$ .

- (b) Is this information sufficient to prove that the heuristic is consistent? If so, prove it. If not, give a counter-example.

**Solution:** No. We only know that  $d_m \geq h(m)$  and  $d_n \geq h(n)$ , but that doesn't mean that  $d_m - d_n \geq h(m) - h(n)$ . For example, we could have  $d_m = h(m)$ , but  $d_n > 0$  while  $h(n) = 0$ , then  $h(m) - h(n) = d_m > d_m - d_n$ .

**Question 4** (7 points)

Consider the Bayesian network shown below, where all three random variables ( $A$ ,  $B$  and  $C$ ) are binary.



Suppose the model has the following parameters:

$$P(A = T) = p$$

$$P(B = T|A = F) = q$$

$$P(B = T|A = T) = r$$

$$P(C = T|B = F) = s$$

$$P(C = T|B = T) = t$$

where  $p, q, r, s$  and  $t$  are some constants. Under what conditions is  $P(A = T|B = T, C = T) > P(A = T)$ ? Specify your answer as an inequality in terms of the model parameters  $p, q, r, s$ , and  $t$ .

**Solution:**

$$\frac{P(A = T, B = T, C = T)}{P(A = T, B = T, C = T) + P(A = F, B = T, C = T)} > P(A = T)$$

$$\frac{prt}{prt + (1 - p)qt} > p$$

**Question 5 (7 points)**

Alice and Bob play the game “Rock, Paper, Scissors” using the standard hand shapes (rock, paper, and scissors), but they use their own special rules to determine the winner of each game. Specifically, they use the following payoff matrix. In this matrix, the number above each diagonal gives the points won by Bob, the number below the diagonal gives the points won by Alice. For example, if Alice chooses “Rock” and Bob chooses “Paper” (denote this outcome as (A=Rock,B=Paper)) then Alice earns -6 points, and Bob earns 1.

Alice's Move	Bob's Move		
	Rock	Paper	Scissors
Rock	-8 / -2	-6 / 1	4 / -5
Paper	2 / 7	9 / -10	-9 / 3
Scissors	8 / 10	-7 / -4	3 / -3

- (a) Are there any fixed-strategy Nash equilibria for this game? If so, what are they?

**Solution:** (A=Scissors,B=Rock) is the only fixed-strategy equilibrium.

- (b) What is/are the Pareto optimal outcome(s) for this game?

**Solution:** (A=Scissors,B=Rock) and (A=Paper,B=Paper) are both Pareto-optimal.

- (c) Suppose that Bob chooses “Rock” with probability  $p$ , “Paper” with probability  $q$ , and “Scissors” with probability  $1 - p - q$ . Write two equations that could be solved to find values of  $p$  and  $q$  such that it is rational for Alice to choose her move at random. Do not simplify.

**Solution:** Any two of the following three equations is a correct answer:

$$-8p - 6q + 4(1 - p - q) = 2p + 9q - 9(1 - p - q)$$

$$-8p - 6q + 4(1 - p - q) = 8p - 7q + 3(1 - p - q)$$

$$2p + 9q - 9(1 - p - q) = 8p - 7q + 3(1 - p - q)$$

**Question 6** (7 points)

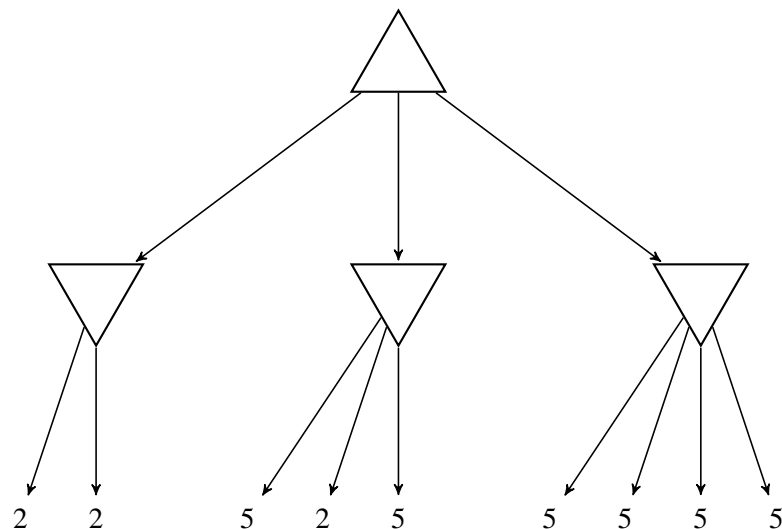
I have just invented a zero-sum, two-person sequential game called Tic-Tac-Go. The game is played on a three-by-three grid; denote each square as  $(r,c)$ , where the row is  $1 \leq r \leq 3$  and the column is  $1 \leq c \leq 3$ . In this game, Max places an X on one of the squares  $(1,1)$ ,  $(1,2)$ , or  $(2,2)$ . Min then places an O on one of the squares next to the X (right, left, up, or down). Max then earns one point for every empty square that is closer to the X than to the O (lower Manhattan distance). For example, the three tables below show cases in which Max wins 2, 5, and 5 points, respectively:

X	O	

O	X	

	X	O

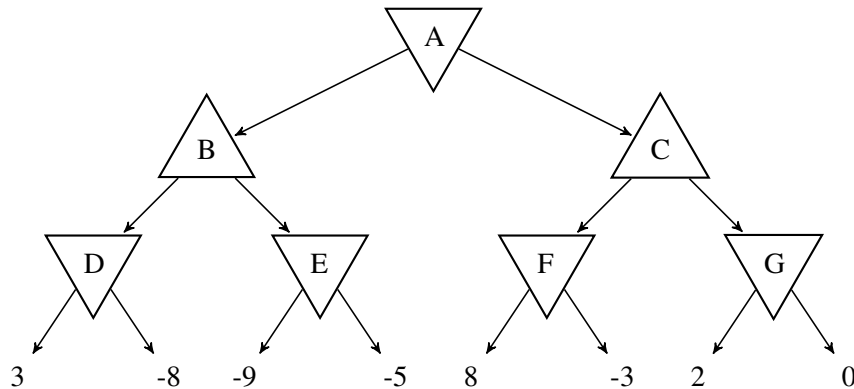
Draw a complete minimax tree for this game. Use upward-pointing triangles for Max nodes, and downward-pointing triangles for Min nodes. Write the numerical value, for Max, of each terminal node.

**Solution:**



**Question 7** (7 points)

Consider using the alpha-beta algorithm to search the tree below. When the alpha-beta algorithm starts, the  $A$  node has  $\alpha = -\infty$ ,  $\beta = \infty$ . When any node is expanded, it begins by inheriting  $\alpha$  and  $\beta$  from its parent, and may then update  $\alpha$  and  $\beta$  after evaluating its children.



- (a) When the  $E$  node is first expanded, what value of  $\alpha$  does it inherit?

**Solution:** The  $B$  node has previously set  $\alpha = -8$ , which is inherited by the  $E$  node.  
The wording of the problem makes the following alternative solution also correct: it is possible for the  $E$  node to be evaluated before the  $D$  node, in which case  $\alpha = -\infty$ .

- (b) When the  $C$  node is first expanded, what value of  $\beta$  does it inherit?

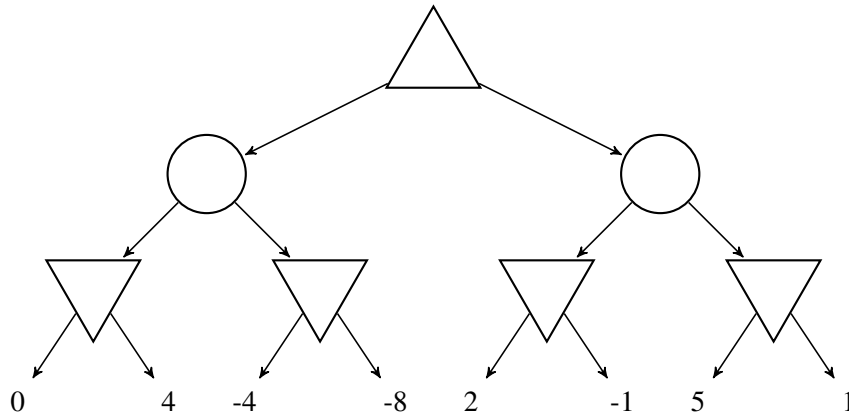
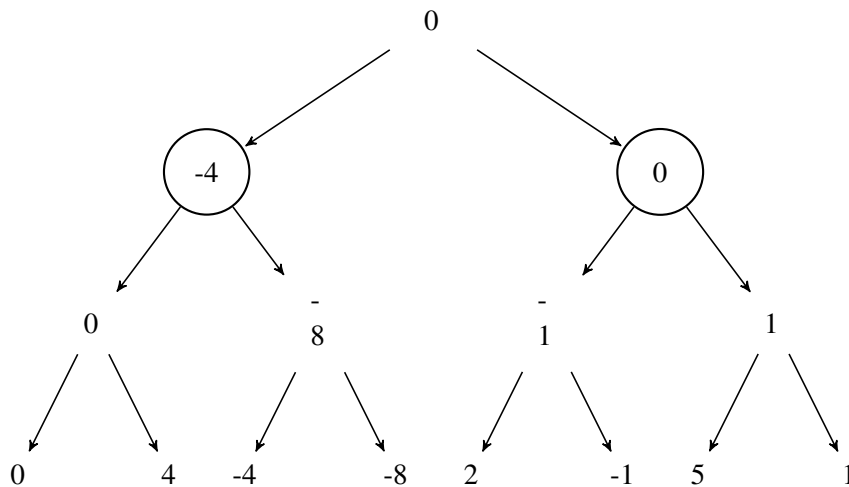
**Solution:** The  $A$  node has previously set  $\beta = -8$ , which is inherited by the  $C$  node.  
The wording of the problem makes the following alternative solution also correct: it is possible for the  $C$  node to be evaluated before the  $B$  node, in which case  $\beta = +\infty$ .

- (c) In the tree above, draw an X through any arrow leading to a node that will not be expanded during the alpha-beta search.

**Solution:** There should be an X through the arrow leading to the -5-valued terminal node, and an X through the arrow leading to node G.

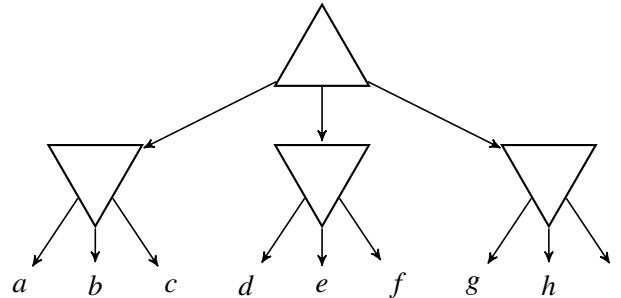
**Question 8** (7 points)

In the tree shown below, circles are chance nodes in which every child is equally likely. Upward-pointing triangles are max nodes, downward-pointing triangles are min nodes. Enter the numerical value of each node inside the corresponding node shape.

**Solution:**

**Question 9** (7 points)

Stochastic search (a.k.a. Monte Carlo tree search) is a useful heuristic, but sometimes gives a value that is quite different from the real value of a node. Consider the tree below, where  $a, b, c, d, e, f, g, h$ , and  $i$  are some real-valued constants.



Consider two different methods for estimating the value of the top node:

1. Minimax search
2. Monte Carlo tree search: choose one of the three Max moves uniformly at random, then, from the resulting child, choose a Min move uniformly at random

The value returned by Monte Carlo tree search is random, but there is one rather unique arrangement of the variables  $a, b, c, d, e, f, g, h$ , and  $i$  that will cause Monte Carlo tree search to *always* return a value that is greater than or equal to the value returned by a Minimax search. What is that arrangement? Your answer should be an equality or inequality, or a set of equalities or inequalities, specifying a relationship among the variables  $a, b, c, d, e, f, g, h$ , and  $i$  that will cause Monte Carlo tree search to *always* return a value greater than or equal to the value returned by Minimax.

**Solution:** Here is one correct answer:

$$\max(\min(a, b, c), \min(d, e, f), \min(g, h, i)) = \min(a, b, c, d, e, f, g, h, i)$$

Here is another way to say the same thing:

$$\min(a, b, c) = \min(d, e, f) = \min(g, h, i)$$

The wording of the problem only requires that the specified condition be sufficient, not necessary. Therefore the following solution is also a correct answer:

$$a = b = c = d = e = f = g = h = i$$

**Question 10** (7 points)

Once every ten years, Briar Rose awakens, and checks the value of her stock portfolio. With probability 60%, she finds that the Dow Jones Industrial Average (DJIA) has doubled while she slept; with probability 40%, she finds that it has halved. She then sells whatever assets she had during the past ten years, and uses all of the money for one of two things: either she purchases stocks (a portfolio whose value equals the DJIA), or she purchases gold (suppose that, in her world, the value of gold never changes). After taking one of these two actions, she goes back to sleep for ten more years.

Define the state of her finances in decade  $d$  to be a tuple specifying the value of her wealth ten years ago,  $w_{d-1}$ , and the value of her wealth now,  $w_d$ , thus  $s_d = (w_{d-1}, w_d)$ . Because the stock market either doubles or halves, while gold never changes, this tuple is always  $s \in \{(x, 0.5x), (x, x), (x, 2x)\}$  for some real number  $x$ . Define the reward to be the amount she has earned this decade:  $R(s_d) = w_d - w_{d-1}$ . Consider using value iteration to find  $U_t(s)$ , the  $t$ -step approximation of the utility of state  $s$ . Starting with  $U_0(s) = 0$  for all states, find  $U_1(s)$  and  $U_2(s)$  for each of the states  $s \in \{(x, 0.5x), (x, x), (x, 2x)\}$ . Your answer should be six equations, showing  $U_1((x, 0.5x)), \dots, U_2((x, 2x))$ , each as a function of  $x$  and  $\gamma$ , where  $\gamma$  is the discount factor.

**Solution:**

$$U_1((x, 0.5x)) = -0.5x$$

$$U_1((x, x)) = 0$$

$$U_1((x, 2x)) = x$$

For  $U_2(s)$ , the correct solution takes into account that  $s'$  is defined in terms of  $x'$ , which may be  $2x$  or  $0.5x$ . That yields the following solution:

$$U_2((x, 0.5x)) = -0.5x + \gamma \max(0, 0.6x' + 0.4(-0.5x'))$$

$$= -0.5x + \gamma \max(0, 0.3x + 0.4(-0.25x))$$

$$U_2((x, x)) = \gamma \max(0, 0.6x + 0.4(-0.5x))$$

$$U_2((x, 2x)) = x + \gamma \max(0, 0.6x' + 0.4(-0.5x'))$$

$$= x + \gamma \max(0, 1.2x + 0.4(-x))$$

**Question 11** (7 points)

Consider using model-based reinforcement learning to learn the maze in an online game. Suppose that this maze is constructed so that every room in the maze has three exits: left (L), right (R), and forward (F). The maze is magic: the ending point of a tunnel sometimes changes at random, so, for example, taking tunnel L from room  $s = 431$  will not always lead to the same place. Fortunately, each room in the maze has a sign specifying its room number, so that, even though you can't be sure where you're going, you can always be sure where you are. Every time you play the game, the software always starts you in room 431.

- (a) Define  $P(s'|s = 431, a = L)$  to be the probability of arriving in room  $s'$  if you start in room 431, and take tunnel L. Suppose you want to estimate  $P(s'|s = 431, a = L)$  with a precision of roughly  $\epsilon = 0.1$ . How many times do you need to perform action  $a = L$  in room  $s = 431$  in order to estimate  $P(s'|s = 431, a = L)$  with a precision of  $\epsilon = 0.1$ ?

**Solution:** Approximately  $\frac{1}{\epsilon} = 10$  times.

- (b) Consider the set of all rooms that can be reached by starting in room 431, and taking  $t$  tunnels in sequence (i.e., choose a tunnel, take it to a new room, then repeat this  $t - 1$  more times). Call the set of rooms that can be reached in this way  $\mathbb{S}_t$ . As a function of  $t$  and  $\epsilon$ , how many times must you play the game in order to estimate  $P(s'|s, a)$  for every room  $s \in \mathbb{S}_t$ , and for every action  $a \in \{L, R, F\}$ , with a precision of  $\epsilon$ ?

**Solution:** There are  $3^t$  rooms in the set, thus there are  $3^{t+1}$  possible combinations of  $(s, a)$ . For each one, we need to play the game  $1/\epsilon$  times to estimate the probability, thus in total we need to play the game  $3^{t+1}/\epsilon$  times.

**Question 12** (7 points)

One reason to use a replay buffer, in Q-learning, is that the Q-learning algorithm does not immediately learn about the long-term consequences of an action. For example, suppose that your replay buffer contains the following  $(s_t, a_t, R(s_t), s_{t+1}, a_{t+1})$  tuples:

$s_t$	$a_t$	$R(s_t)$	$s_{t+1}$	$a_{t+1}$
431	L	-0.04	1024	R
1024	R	0.5	516	R

Consider using SARSA Q-learning to estimate  $Q_t(s, a)$ , starting with  $Q_0(s, a) = 0$  for all  $s$  and  $a$ . When you estimate  $Q_t(s, a)$ , use every row of the replay buffer, not just the  $t^{\text{th}}$  row. Let  $\alpha$  be the learning rate, and let  $\gamma$  be the discount factor; as a function of  $\alpha$  and  $\gamma$ , find  $Q_1(431, L)$ ,  $Q_1(1024, R)$ , and  $Q_2(431, L)$ .

**Solution:**

$$\begin{aligned} Q_1(431, L) &= Q_0(431, L) + \alpha (R(431) + \gamma Q_0(1024, R) - Q_0(431, L)) \\ &= -0.04\alpha \end{aligned}$$

$$\begin{aligned} Q_1(1024, R) &= Q_0(1024, R) + \alpha (R(54) + \gamma Q_0(516, R) - Q_0(1024, R)) \\ &= 0.5\alpha \end{aligned}$$

$$\begin{aligned} Q_2(431, L) &= Q_1(431, L) + \alpha (R(431) + \gamma Q_1(1024, R) - Q_1(431, L)) \\ &= -0.04\alpha - 0.04\alpha + 0.5\alpha^2\gamma + 0.04\alpha^2 \end{aligned}$$

**Question 13** (7 points)

Consider a deep Q-learning algorithm, in which  $Q(\vec{s}, a)$  is approximated, after  $t$  iterations of training, by

$$Q_t(\vec{s}, a) = \vec{w}_t^T \vec{h}_t(\vec{s}, a),$$

where  $\vec{w}_t^T = [w_{t,1}, \dots, w_{t,5}]$  is a 5-dimensional weight vector, and  $\vec{h}_t^T(\vec{s}, a) = [h_{t,1}(\vec{s}, a), \dots, h_{t,5}(\vec{s}, a)]$  is a 5-dimensional hidden node activation vector. Consider two particular state vectors,  $\vec{s} = [0.3, 0.6]^T$  and  $\vec{s} = [-0.2, 0.5]^T$ , and two particular actions,  $a \in \{L, R\}$ . Suppose that, in response to these state vectors and these actions, the network currently computes the following hidden node activation vectors:

$\vec{s}^T$	$a$	$\vec{h}_t^T(\vec{s}, a)$
$[0.3, 0.6]$	L	$[a, b, c, d, e]$
$[0.3, 0.6]$	R	$[f, g, h, i, j]$
$[-0.2, 0.5]$	L	$[k, l, m, n, o]$
$[-0.2, 0.5]$	R	$[p, q, r, s, t]$

Suppose that, in the  $t^{\text{th}}$  iteration of deep Q-learning, using only the  $t^{\text{th}}$  entry from the replay buffer, the output weight vector is  $\vec{w}_t = [1, 0, 0, 0, 0]^T$ , the input state is  $\vec{s}_t = [0.3, 0.6]^T$ , the action is  $a_t = R$ , the reward is  $R(\vec{s}_t) = -0.04$ , the resulting state vector is  $\vec{s}_{t+1} = [-0.2, 0.5]^T$ , and the loss function is

$$\mathcal{L} = \frac{1}{2} (Q_t(\vec{s}_t, a_t) - Q_{\text{local}}(\vec{s}_t, a_t))^2$$

Assume a discount factor of  $\gamma$ , and assume TD-learning, not SARSA. In terms of any or all of the variables  $a, b, c, \dots, r, s, t$ , and/or  $\gamma$ , what is the value of  $\nabla_{\vec{w}_t} \mathcal{L}$ ?

**Solution:** First, find  $Q_{\text{local}}$ .

$$\begin{aligned} Q_{\text{local}}(\vec{s}_t, a_t) &= R(\vec{s}_t) + \gamma \max_{a'} Q_t(\vec{s}_{t+1}, a') \\ &= -0.04 + \gamma \max(k, p) \end{aligned}$$

Now in terms of  $Q_{\text{local}}$ , we can write:

$$\begin{aligned} \nabla_{\vec{w}_t} \mathcal{L} &= (Q_t(\vec{s}_t, a_t) - Q_{\text{local}}(\vec{s}_t, a_t)) \nabla_{\vec{w}_t} Q_t(\vec{s}_t, a_t) \\ &= (Q_t(\vec{s}_t, a_t) - Q_{\text{local}}(\vec{s}_t, a_t)) \vec{h}_t(\vec{s}_t, a_t) \\ &= (f + 0.04 - \gamma \max(k, p)) \begin{bmatrix} f \\ g \\ h \\ i \\ j \end{bmatrix} \end{aligned}$$

**Question 14** (7 points)

A robot vacuum cleaner is vacuuming a square room that is  $D$  meters on a side. The robot is a disc of radius  $R$  meters. The room contains just one obstacle: a rectangular box, of length  $L$  and width  $W$ , located more than  $2R$  meters away from any wall, with its sides parallel to the walls of the room. Because of the shape of the robot, it can never be closer than  $R$  to any wall, or to any edge of the box. Define the configuration space of the robot to be  $(x,y)$ , the two-dimensional position of the center of the robot. In order to plan its trajectory, the robot discretizes the room into a set of  $R \times R$  squares. How many such squares are reachable in the robot's configuration space?

**Solution:** The reachable area of the room is  $(D - 2R)^2 - (L + 2R)(W + 2R)$ . Dividing that area into  $R \times R$  squares gives approximately

$$\frac{(D - 2R)^2 - (L + 2R)(W + 2R)}{R^2}$$

squares.

If one considers that only the *center* of each square needs to be within the reachable area, then the total area covered by the entire *periphery* of each square could be as large as  $(D - R)^2 - (L + R)(W + R)$ . If one makes this assumption, then the number of squares is

$$\frac{(D - R)^2 - (L + R)(W + R)}{R^2}$$



SCRATCH PAPER

SCRATCH PAPER