CS 188 Introduction to Summer 2021 Artificial Intelligence

Midterm

- You have approximately 110 minutes.
- The exam is open book, open calculator, and open notes. Outside internet use is not allowed.
- For multiple choice questions,
 - means mark **all options** that apply
 - \bigcirc means mark a single choice

First name	
Last name	
SID	

Q1.	Potpourri	/12
Q2.	Maximizing Different Reward Metrics	/20
Q3.	Q-Learning: Weights Amnesia	/15
Q4.	Search: Algorithms	/10
Q5.	Search: Bowser's Kingdom	/8
Q6.	Caramel Crush	/15
Q7.	Game Trees	/20
	Total	/100

For staff use only:

THIS PAGE IS INTENTIONALLY LEFT BLANK

Q1. [12 pts] Potpourri

- (a) [2 pts] Your friend Julia proposes a new algorithm, a modified version of Inference by Enumeration called **mod-IBE**. She first marginalizes out hidden variables, then select entries in the joint table that are consistent with evidence. And then, she normalizes. Is mod-IBE a valid approach for probabilistic inference?
 - Yes, since marginalizing hidden variables doesn't ever affect the rows consistent with evidence.
 - \bigcirc Yes, but not for the reason above.
 - No, since marginalizing combines entries consistent with evidence with those inconsistent with evidence.
 - \bigcirc No, but not for the reason above.
- (b) Consider a MDP with deterministic transitions, representing a Golden Bear walking around. S is the set of possible states and A is the set of possible actions. We have all the optimal Q-values Q(s, a) and values V(s), and so can reconstruct an optimal policy.

But unfortunately, Bears behave slightly irrationally. We'll model our Bear as choosing an action a to take from state s with probability $P_O(a|s)$.

(i) [3 pts] Which of the following options for $P_Q(a|s)$ would result in a valid probability distribution that models Bear's slightly irrational behavior?



- (ii) [3 pts] True or False: The modified MDP posed by this problem can be modeled as a vanilla expectimax tree. "Vanilla" is defined as an expectimax tree as seen in lecture: alternating layers of max and expectation nodes, with multiple children for every node.
 - True, since there are alternating turns between deterministic and nondeterministic behavior.
 - \bigcirc True, but not for the reason above.

 \bigcirc False, since introducing $P_Q(a|s)$ causes the expectimax tree to have two layers of expectation nodes for every max node.

- False, but not for the reason above.
- (iii) [2 pts] We're interested in our Bear's trajectories. A **trajectory** is a sequence of states and actions $\{s_1, a_1, s_2, a_2..., s_T, a_T, s_{T+1}\}$ which detail an agent's movements within the MDP.

What is the probability of observing a full trajectory $\{s_1, a_1, s_2, a_2, \dots, s_T, a_T, s_{T+1}\}$ given that we know the Bear has already gone through a sub-trajectory $\{s_1, a_1, s_2, a_2, \dots, s_t, a_t, s_{t+1}\}, 1 \le t < T$?

$$P(\{s_1, a_1, s_2, a_2 \dots s_{T+1}\} | \{s_1, a_1, s_2, a_2 \dots s_{t+1}\}) = \dots$$





(c) [2 pts] Select all of the following algorithms that would require a reinforcement learning agent to construct estimates of transition probabilities and rewards in the course of updating its policy:



Q-Learning

Q2. [20 pts] Maximizing Different Reward Metrics



Consider a robot navigating in a 3×4 2D grid with walls around the edges as shown above. The robot wants to move around and get high reward, but it needs your help to plan ahead to achieve this!

The robot's state is represented by (row, column), with the starting position at (1, 1). He only has two actions $a \in \{\text{right}, \text{up}\}$. The robot cannot leave the grid. The robot can choose to move into the wall, but its state will stay the same after the action.

Note that in this grid, the transitions described above are *deterministic*. We characterize this MDP's dynamics with a function f(s, a) = s'. When the robot transitions into state s, it receives a reward R(s) which depends only on the state it transitions to. R(s) for every cell is shown above. R(s) is defined to be 0 at all grid cells without a numerical label.

For example, the robot can incur a -10 reward by using the up action from (3, 3) without changing its position.

(i) [2 pts] What is the equivalent equation for computing optimal values in the MDP defined above in terms of transition (a) function f(s, a) and reward function R(s)? Assume no discounting $(\gamma = 1)$. Also, define x@y = max(x, y). For each letter (A), (B), (C), (D), (E), fill in a single entry for the term corresponding to the correct equation to implement value iteration for this MDP. Select one bubble per row to form the whole equation.

$$V^*(s) = (\underline{A})(\underline{B}) \left[R((\underline{C}))(\underline{D})V^*((\underline{E})) \right]$$
(1)

- (A): \bigcirc 1 \bigcirc max_a \bigcirc \mathbb{E}_a
- (B): $\bigcirc \sum_{s'} \bigcirc \sum_{a} \bigcirc 1$ (C): $\bigcirc f(s,a) \bigcirc f(s',a) \bigcirc V^*(s)$
- (D): $\bigcirc \times \bigcirc + \bigcirc @$
- (E): $\bigcirc f(s,a) \bigcirc R(s') \bigcirc R(s)$
- (ii) [2 pts] What is the optimal value of the starting state when you can at most take N steps in the environment?

N = 3?	
N = 5?	
N = 10	2

- (b) Normally, the optimal value function V^* represents the maximum possible sum of future rewards since we typically want to find a policy which maximizes the sum of future rewards in the MDP. However, suppose that we instead want to find a policy which maximizes the future **product** of rewards. Thus, we want to compute V^* such that it represents the maximum possible product of future rewards.
 - (i) [2 pts] Give an expression for a corresponding equation for V^* which will compute this quantity. Define x@y =max(x, y). For each letter (A), (B), (C), (D), (E), fill in a single entry for the term corresponding to the correct equation to implement value iteration to maximize the product of future rewards. Select one bubble per row to form the whole equation.

$$V^*(s) = (\mathbf{A})(\mathbf{B}) \left[R((\mathbf{C}))(\mathbf{D})V^*((\mathbf{E})) \right]$$
(2)

- (A): \bigcirc 1 \bigcirc max_a \bigcirc \mathbb{E}_a
- (B): $\bigcirc \sum_{s'} \bigcirc \sum_{a} \bigcirc 1$ (C): $\bigcirc f(s,a) \bigcirc f(s',a) \bigcirc V^*(s)$

- (D): $\bigcirc \times \bigcirc + \bigcirc @$ • (E): $\bigcirc f(s,a) \bigcirc R(s') \bigcirc R(s)$
- (ii) [2 pts] Given this new objective, what is the optimal value of the starting state when you can at most take N steps in the environment?



- (c) Now, suppose that we instead want to find a policy which maximizes the **highest** reward achieved at any timestep. Thus, we want to compute V^* such that it represents the maximum possible value of the **highest** reward the robot ever achieves.
 - (i) [2 pts] Give an expression for a corresponding equation for V^* which will compute this quantity. Define x@y =max(x, y). For each letter (A), (B), (C), (D), (E), fill in a single entry for the term corresponding to the correct equation to implement value iteration to maximize the highest reward achieved by the robot. Select one bubble per row to form the whole equation.

$$V^*(s) = (\mathbf{A})(\mathbf{B}) \left[R((\mathbf{C}))(\mathbf{D})V^*((\mathbf{E})) \right]$$
(3)

- (A): \bigcirc 1 \bigcirc max_a \bigcirc \mathbb{E}_a
- (B): $\bigcirc \Sigma_{s'} \bigcirc \Sigma_a \bigcirc 1$
- (C): $\bigcirc f(s,a) \bigcirc f(s',a) \bigcirc V^*(s)$
- (D): $\bigcirc \times \bigcirc + \bigcirc$ @
- (E): $\bigcirc f(s,a) \bigcirc R(s') \bigcirc R(s)$
- (ii) [2 pts] Given this new objective, what is the optimal value of the starting state when you can at most take N steps in the environment?



- (d) The standard Bellman equation provides a condition the optimal policy must satisfy by performing a one-step look-ahead by defining a recursive relationship between $V^*(s)$ and $V^*(s')$, where s' is a possible future next state. However, in general, one could define an analagous update with a two-step lookahead.
 - (i) [3 pts] Which of the following would be the correct way to formulate the Bellman equation with a two-step lookahead?

$$V^*(s) = (_\mathbf{A}_) \left[R((_\mathbf{B}_)) + (_\mathbf{C}_) \left[R((_\mathbf{D}_)) + V^*((_\mathbf{E}_)) \right] \right)$$
(4)

- (A): $\bigcirc \sum_a \bigcirc \max_a \bigcirc \mathbb{E}_a$ (B): $\bigcirc s \bigcirc f(s,a) \bigcirc s'$
- (C): $\bigcirc \Sigma_{a'} \bigcirc \max_{a'} \bigcirc \mathbb{E}_{a'}$
- (**D**): $\bigcirc f(s,a') \bigcirc f(f(s,a),a') \bigcirc f(f(s,a'),a)$
- (E): $\bigcirc f(s,a') \bigcirc f(f(s,a),a') \bigcirc f(f(s,a'),a)$
- (ii) [2 pts] If we perform value iteration using the above two-step Bellman equation, what is the first iteration at which state (1, 1) achieves a positive value? Assume as usual that $V_0(s) = 0$ for all s.
- (iii) [1 pt] How many iterations of two-step Bellman equation will it take for the value function to converge if N=6? Assume as usual that $V_0(s) = 0$ for all s.
- (iv) [2 pts] If we were to define the two-step look-ahead Bellman equation for highest reward and product of rewards, which of the three formulations will converge the fastest for the above MDP and N = 3?
 - Sum of rewards
 - Highest reward

- Product of rewards
 Two of them converge the fastest, in the same number of iterations
 All of them converge in the same number of iterations

Q3. [15 pts] Q-Learning: Weights Amnesia

Pacman travels to a simple gridworld with a dummy ghost who doesn't know how to move! Unfortunately, he caught a fever on its way which messed up the weights $w \in \mathbb{R}^2$ it learned for approximate Q-Learning. The only thing he remembered is the feature extraction function $f(s, a) = \begin{bmatrix} d_{ghost} + 1 \\ d_{food} + 1 \end{bmatrix}$, where d_{ghost} and d_{food} are the Manhattan distances to the ghost and the food, respectively, after taking action *a* from state *s*.

Pacman tells you that its learning rate $\alpha = 0.2$, its discount factor $\gamma = 0.9$, and its current weights $w = \begin{bmatrix} -3 \\ 2 \end{bmatrix}$. Work out the following s to help it learn a better weight!

You also notice some properties of this gridworld:

- All actions will succeeed with probability 1. If an action makes Pacman move into a wall, then Pacman's position will remain unchanged for this action.
- If Pacman runs into a ghost, the game ends, and Pacman will receive a -10 penalty. (E is a terminal state)
- If Pacman runs into a food, the game ends, and Pacman will receive a +10 reward. (A is a terminal state)
- Otherwise, Pacman will receive a -1 penalty.



Figure 1

(a) (i) [2 pts] Compute the following features and Q-values.



- (i) [1 pt] Now, in order to learn a good policy, you tell Pacman to use ε-greedy policy. Assume that ε = 0.6, what's the probability assigned to the action with the highest Q value?
 P(a*|S) =
 - (ii) [2 pts] Pacman then randomly sampled a RIGHT(→) action and gathers a transition experience, what will be Pacman's updated weights?



(iii) [4 pts] After updating the weights, rank the 4 new Q-values at state S in a decreasing order.



- (c) Assume that you train Pacman with an infinite number of episodes, a proper learning rate α , and everything else kept unchanged.
 - (i) [2 pts] Is it able to learn the true optimal Q-values?
 - Yes, since we visited every state infinitely many times, Q-learning will converge to optimal.
 - \bigcirc Yes, but not for the reason above.
 - \bigcirc No, since the *e*-greedy policy we use for learning is not optimal.
 - \bigcirc No, but not for the reason above.
 - (ii) [2 pts] Is it able to learn the true optimal state values?
 - Yes, because the learned Q-values could be optimal.
 - \bigcirc Yes, but not for the reason above.
 - \bigcirc No, because the learned Q-values are never optimal.
 - \bigcirc No, but not for the reason above.

(iii) [2 pts] Is it able to learn the true optimal policy?

- Yes, because it could learn the optimal Q-values.
- \bigcirc Yes, but not for the reason above.
- No, because it never learns the optimal Q-values.
- \bigcirc No, but not for the reason above.

Q4. [10 pts] Search: Algorithms



Figure 2: Search graph

	Α	В	\mathbf{C}	\mathbf{S}
H-1	0	0	0	0
H-2	6	7	1	7
H-3	7	7	1	7
H-4	4	7	1	7

Figure 3: Heuristics table

(a) Consider the search graph and heuristics shown above. Select all of the goals that could be returned by each of the search algorithms below.

For this question, if there is a tie on the fringe, assume the tie is broken randomly.

- (i) $\begin{bmatrix} 1 \text{ pt} \end{bmatrix} \text{DFS}$ $\Box G_1$ $\Box G_2$ $\Box G_3$
- (ii) $\begin{bmatrix} 1 & \text{pt} \end{bmatrix}$ BFS $\begin{bmatrix} G_1 \\ G_2 \\ G_3 \end{bmatrix}$

(iii) $\begin{bmatrix} 1 & \text{pt} \end{bmatrix}$ UCS $\begin{bmatrix} G_1 \\ G_2 \\ G_3 \end{bmatrix}$



SID:

 $\Box G_3$

(v) $\begin{bmatrix} 1 \text{ pt} \end{bmatrix}$ Greedy with H-3 $\begin{bmatrix} G_1 \\ G_2 \\ G_3 \end{bmatrix}$ (vi) $\begin{bmatrix} 1 \text{ pt} \end{bmatrix} A^*$ with H-1 $\begin{bmatrix} G_1 \\ G_2 \\ G_3 \end{bmatrix}$ (vii) $\begin{bmatrix} 1 \text{ pt} \end{bmatrix} A^*$ with H-2 $\begin{bmatrix} G_1 \\ G_2 \\ G_3 \end{bmatrix}$ (vii) $\begin{bmatrix} 1 \text{ pt} \end{bmatrix} A^*$ with H-2

(b) For each heuristic, indicate whether it is consistent, admissible, or neither (select more than one option if appropriate):



(ii)	[1 pt] H-2	
	Consistent	
	Admissible	
	Neither	

(iii)	[1]	pt] H-3
		Consistent
		Admissible

Neither

Q5. [8 pts] Search: Bowser's Kingdom

(a) Mario is looking for Princess Peach in Bowser's N x M kingdom. Along the way, he has to collect all K > 0 stationary stars by landing in the same spot as the star. Every time he collects a star, he gets a power up that allows him to move one extra space per timestep. For example, if Mario collects 1 star, he can now move up to 2 spaces per time step, when he collects a 2nd star he can move 3 spaces per time step and so on. The power up comes into effect in the timestep after Mario collects the star and lasts for the rest of the search.

Princess Peach is a capable independent woman who has already escaped and is now **also moving around one space per timestep** in the kingdom. Princess Peach cannot collect stars. Both Mario and Princess Peach can only move North, South, East, and West. The search ends when Mario and Peach land in the same space, and Mario has collected every star in the kingdom.

(i) [3 pts] What is the size of the minimum state space for this problem?

$\bigcirc (MN)^2$	$\bigcirc (MN)^2(2^K)$	$\bigcirc 4 * (MNK)^2$
$\bigcirc 2^{MN}$	$\bigcirc (MNK)^2$	$\bigcirc 4 * 2^{K(MN)^2}$
$\bigcirc MN(2^K)$	$\bigcirc 4 * MN(2^K)$	$\bigcirc 4 * (MN)^2 (2^K)$

Now Bowser has been alerted and has begun scanning the kingdom by quadrants. At each time step Bowser will scan a quadrant in clockwise order starting from the top right quadrant. If Mario is caught, i.e. he is in the quadrant that is currently being scanned, then it will be game over.

(ii) [1 pt] If state space in the previous part was X, state the size of the minimum state space with this updated game rule?

\bigcirc	4X	\bigcirc	2^4X
0	X	\bigcirc	4 ^{<i>X</i>}
\bigcirc	$\frac{X}{4}$	\bigcirc	None of the above

- (iii) [4 pts] Assume that the scanning is still in place. Check all of the following that are admissible heuristics (Note that <u>any</u> distances from Mario to stars or from stars to Princess Peach are 0 if there are no remaining stars):
 - 1 for every state.

Manhattan distance from Mario to Princess Peach divided by 2.

Manhattan distance from Mario to Princess Peach divided by (K + 2).

Sum of Manhattan distances from Mario to all stars and to Princess Peach.

(Sum of Manhattan distances from Mario to each star + Manhattan distance from Mario to Princess Peach)/(3K).

(Manhattan distance from Mario to furthest star + Manhattan distance from that furthest star to Princess Peach)/(K + 2).

(Manhattan distance from Mario to closest star + Manhattan distance from that closest star to Princess Peach)/K.

Q6. [15 pts] Caramel Crush

Mesut is on his phone during a staff meeting and decides to play his favorite mobile game - Caramel Crush. He decides to use informed search to solve the game optimally.

Caramel Crush has an N by N grid of squares that be either be empty, or one of C - 1 differently colored gems. At each turn, Mesut can destroy a colored gem from the grid, also recursively destroying any touching gems of the same color. After each destruction, the gem will fall down (such that there is never an empty space under any gem). The goal of the game is to destroy all gems with as few moves as possible.



- (a) First, let's help Mesut understand the game.
 - (i) [2 pts] Does every state have a unique sequence of moves to the goal state?

○ Yes ○ 1	No
-----------	----

(ii) [2 pts] What is the size of the state space?

\bigcirc	N^2	\bigcirc	C^{N^2}
\bigcirc	CN^2	\bigcirc	N^{2^C}

- (b) Now, Mesut asks us to come up with a heuristic and search algorithm for him. Given a heuristic h(n), state whether it is optimal for A* tree search, A* graph search, both, or neither. Note: each action has a cost of 1
 - (i) [2 pts] h_1 = number of gems left

\bigcirc	optimality for A* tree search only	\bigcirc	optimality for neither
\bigcirc	optimality for A* graph search only	\bigcirc	optimality for both

(ii) [2 pts] h_2 = unique number of gem colors left

\bigcirc	optimality for A* tree search only	\bigcirc	optimality for neither
\bigcirc	optimality for A* graph search only	\bigcirc	optimality for both

(iii) [3 pts] h_3 = number of contiguous chunks of the same color gem (as outlined in the figure (iii)) h_3 = number of contiguous chunks of the same color gem (as outlined in the figure (iii)) h_3 = number of contiguous chunks of the same color gem (as outlined in the figure (iii)) h_3 = number of contiguous chunks of the same color gem (as outlined in the figure (iii)) h_3 = number of contiguous chunks of the same color gem (as outlined in the figure (iii)) h_3 = number of contiguous chunks of the same color gem (iiii)) h_3 = nu	ure)
---	------

\bigcirc	optimality for A* tree search only	\bigcirc	optimality for neither
\bigcirc	optimality for A* graph search only	\bigcirc	optimality for both

(c) Mesut gave up coding the Caramel Crush solver.

Instead, let's think about an abstract general search problem with the given properties.

(i) [4 pts] Suppose we construct a heuristic h_1 , such that for every edge n_i to n_j with cost c_{ij} , our heuristic follows the property $\frac{h_1(n_i)}{h_1(n_j)} = e^{c_{ij}}$. What properties would a new heuristic h_2 have if $h_2 = ln(h_1)$?

(Only admissible	\bigcirc	Both
(201	\bigcirc	NT - 141

Only consistent

O Neither

Q7. [20 pts] Game Trees

Pacman and the ghost are playing a game, where Pacman and the ghost have their own score. Both are trying to maximize their personal score. When breaking ties, they will also prefer a larger difference between their score and the other's.

A player can score anywhere between 1 and 10 points.

The figure below shows the game tree of pacman's max node (quadrilateral) followed by the ghost's nodes (hexagons). The scores shown at the leaf nodes follow $\begin{bmatrix} pacmanscore \\ ghostscore \end{bmatrix}$

For example, pacman would choose $\begin{bmatrix} 4\\2 \end{bmatrix}$ over $\begin{bmatrix} 3\\2 \end{bmatrix}$ and the ghost would choose $\begin{bmatrix} 5\\9 \end{bmatrix}$ over $\begin{bmatrix} 6\\9 \end{bmatrix}$.

(a) Fill in the blanks with the pair of scores preferred by each node of the game tree.



(i) [1 pt] What branch is chosen in node i)?

- () c)
- () d)
- () e)
- (ii) [1 pt] What branch is chosen in node ii)?
 - \bigcirc f)
 - () g)
 - () h)
- (iii) [1 pt] What branch is chosen in node iii)?
 - () j)
 - () k)
 - 0 1)
- (iv) [1 pt] What branch is chosen in node iv)?
 - () m)
 - () n)
 - () o)

- (v) [1 pt] What branch is chosen in node v)?
 - () a)
 - () p)
 - () q)
 - () b)
- (b) [3 pts] You decide to save time by pruning branches in your game tree search. Mark the branches that do not need to be explored. Assume that branches are explored from left to right.



(c) [2 pts] Suppose we have square function $f(x) = x^2$ that is applied to the leaves of the same game tree as shown below. Now determine which branches will be pruned by crossing out the line of branches that do not need to be explored. Again, assume that branches are explored from left to right.





SID:

(d) We want to identify the relationships between the values at the root (value of A and B) and variables in the game tree. Both game trees have the same variables (w, x, y, z, t) and $\mathbf{x} > \mathbf{t}$.

Upwards triangles represent max nodes while downwards triangles represent min nodes. Change nodes are circles and prefer the left action with probability 0.75.

For each of the following, indicate whether the relationship is always, sometimes, or never true.

- (i) [1 pt] *A* > *t*
 - Always true
 - \bigcirc Sometimes true
 - O Never true
- (ii) [1 pt] *A* < *t*
 - \bigcirc Always true
 - \bigcirc Sometimes true
 - \bigcirc Never true
- (iii) [1 pt] $A \ge x OR A \ge y$
 - \bigcirc Always true
 - \bigcirc Sometimes true
 - \bigcirc Never true
- (iv) [1 pt] $A \le max(w, x, y, z)$
 - \bigcirc Always true
 - Sometimes true
 - \bigcirc Never true
- (v) [1 pt] B < w
 - \bigcirc Always true
 - Sometimes true
 - O Never true
- (vi) [1 pt] $B \ge x$
 - \bigcirc Always true
 - Sometimes true

- \bigcirc Never true
- (**vii**) [2 pts] *B* < *t*
 - \bigcirc Always true
 - \bigcirc Sometimes true
 - \bigcirc Never true
- (viii) [1 pt] $B \ge max(w, x)$
 - Always true
 - \bigcirc Sometimes true
 - \bigcirc Never true
- (ix) [1 pt] A < B
 - \bigcirc Always true
 - \bigcirc Sometimes true
 - O Never true