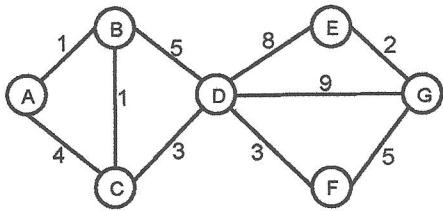


Matthew Olson

Q1. [9 pts] Search



Node	h_1	h_2
A	9.5	10
B	9	12
C	8	10
D	7	8
E	1.5	1
F	4	4.5
G	0	0

Consider the state space graph shown above. A is the start state and G is the goal state. The costs for each edge are shown on the graph. Each edge can be traversed in both directions. Note that the heuristic h_1 is consistent but the heuristic h_2 is not consistent.

(a) [4 pts] Possible paths returned

For each of the following graph search strategies (*do not answer for tree search*), mark which, if any, of the listed paths it could return. Note that for some search strategies the specific path returned might depend on tie-breaking behavior. In any such cases, make sure to mark *all* paths that could be returned under some tie-breaking scheme.

Search Algorithm	A-B-D-G	A-C-D-G	A-B-C-D-F-G
Depth first search	✗	✗	✗
Breadth first search	✗	✗	
Uniform cost search			✗
A* search with heuristic h_1			✗
A* search with heuristic h_2			✗

→ any path
 → shallowest path
 → least cost path
 → optimal path (consistent h_2)
 → optimal path (inconsistent h_2)

(b) Heuristic function properties

Suppose you are completing the new heuristic function h_3 shown below. All the values are fixed except $h_3(B)$.

Node	A	B	C	D	E	F	G
h_3	10	?	9	7	1.5	4.5	0

For each of the following conditions, write the set of values that are possible for $h_3(B)$. For example, to denote all non-negative numbers, write $[0, \infty]$, to denote the empty set, write \emptyset , and so on.

(i) [1 pt] What values of $h_3(B)$ make h_3 admissible?

$$\text{Actual cost} = 1+3+3+5 = 12$$

$$0 \leq h_3 \leq 12$$

(ii) [2 pts] What values of $h_3(B)$ make h_3 consistent?

$$h(A) - h_3(B) \leq \text{cost}(A, B) \Rightarrow h_3(B) \geq 9$$

AND

$$h(B) \leq h(A) \quad h_3(B) \leq 10$$

$$9 \leq h_3(B) \leq 10$$

(iii) [2 pts] What values of $h_3(B)$ will cause A* graph search to expand node A, then node C, then node B, then node D in order?

	c	h	f
A	0	10	10
AB	1	$h_3(B)$	$1+h_3(B)$
AC	4	9	13
ACB	5	$h_3(B)$	$5+h_3(B)$
ACD	7	7	14

$$1+h_3(B) > B$$

OR

$$\begin{cases} 5+h_3(B) \leq 14 \\ 1+h_3(B) \geq 14 \end{cases}$$

3

$$h_3(B) > 12$$

or

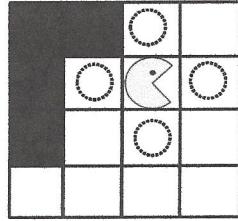
$$\begin{cases} h_3(B) \leq 9 \\ h_3(B) \leq B \end{cases}$$

$$12 \leq h_3(B) \leq 13$$

Q2. [9 pts] InvisiPac

Pacman finds himself to have an invisible “friend”, InvisiPac. Whenever InvisiPac visits a square with a food pellet, InvisiPac will eat that food pellet—giving away its location at that time. Suppose the maze’s size is $M \times N$ and there are F food pellets at the beginning.

Pacman and InvisiPac alternate moves. Pacman can move to any adjacent square (including the one where InvisiPac is) that are not walls, just as in the regular game. After Pacman moves, InvisiPac can teleport into any of the four squares that are adjacent to Pacman, as marked with the dashed circle in the graph. InvisiPac can occupy wall squares.



- (a) For this subquestion, whenever InvisiPac moves, it chooses randomly from the squares adjacent to Pacman. The dots eaten by InvisiPac don’t count as Pacman’s score. Pacman’s task is to eat as many food pellets as possible.

- (i) [1 pt] Which of the following is best suited to model this problem from Pacman’s perspective?

state space search CSP minimax game MDP RL

- (ii) [2 pts] What is the size of a minimal state space for this problem? Give your answer as a product of factors that reference problem quantities such as M , N , F , etc. as appropriate. Below each factor, state the information it encodes. For example, you might write $4 \times MN$ and write *number of directions* underneath the first term and *Pacman’s position* under the second.

$$2^F \cdot M \times N$$

↓ ↓
Food Position of Pacman

- (b) For this subquestion, whenever InvisiPac moves, it always moves into the same square relative to Pacman. For example, if InvisiPac starts one square North of Pacman, InvisiPac will always move into the square North of Pacman. Pacman knows that InvisiPac is stuck this way, but doesn’t know which of the four relative locations he is stuck in. As before, if InvisiPac ends up being in a square with a food pellet, it will eat it and Pacman will thereby find out InvisiPac’s location. Pacman’s task is to find a strategy that minimizes the worst-case number of moves it could take before Pacman knows InvisiPac’s location.

- (i) [1 pt] Which of the following is best suited to model this problem from Pacman’s perspective?

state space search CSP minimax game MDP RL

- (ii) [2 pts] What is the size of a minimal state space for this problem? Give your answer as a product of factors that reference problem quantities such as M , N , F , etc. as appropriate. Below each factor, state the information it encodes. For example, you might write $4 \times MN$ and write *number of directions* underneath the first term and *Pacman’s position* under the second.

$$2^F \cdot M \times N \cdot 2^4$$

↓ ↓ ↓
Food Position
Pacman InvisiPac

- (c) For this subquestion, whenever InvisiPac moves, it can choose freely between any of the four squares adjacent to Pacman. InvisiPac tries to eat as many food pellets as possible. Pacman’s task is to eat as many food pellets as possible.

- (i) [1 pt] Which of the following is best suited to model this problem from Pacman’s perspective?

state space search CSP minimax game MDP RL

- (ii) [2 pts] What is the size of a minimal state space for this problem? Give your answer as a product of factors that reference problem quantities such as M , N , F , etc. as appropriate. Below each factor, state the information it encodes. For example, you might write $4 \times MN$ and write *number of directions* underneath the first term and *Pacman’s position* under the second.

$$2^F \cdot M \times N$$

↓ ↓
Food Pacman Position

Q3. [27 pts] CSPs

(a) Pacman's new house

After years of struggling through mazes, Pacman has finally made peace with the ghosts, Blinky, Pinky, Inky, and Clyde, and invited them to live with him and Ms. Pacman. The move has forced Pacman to change the rooming assignments in his house, which has 6 rooms. He has decided to figure out the new assignments with a CSP in which the variables are Pacman (P), Ms. Pacman (M), Blinky (B), Pinky (K), Inky (I), and Clyde (C), the values are which room they will stay in, from 1-6, and the constraints are:

- i) No two agents can stay in the same room
- ii) $P > 3$
- iii) K is less than P
- iv) M is either 5 or 6
- v) $P > M$
- vi) B is even
- vii) I is not 1 or 6
- viii) $|I-C| = 1$
- ix) $|P-B| = 2$

- (i) [1 pt] **Unary constraints** On the grid below cross out the values from each domain that are eliminated by enforcing unary constraints.

P	X	X	X	4	5	6
B	X	2	X	4	X	6
C	1	2	3	4	5	6
K	1	2	3	4	5	6
I	X	2	3	4	5	X
M	X	X	X	X	5	6

- (ii) [1 pt] **MRV** According to the Minimum Remaining Value (MRV) heuristic, which variable should be assigned to first?

P

B

C

K

I

M

- (iii) [2 pts] **Forward Checking** For the purposes of decoupling this problem from your solution to the previous problem, assume we choose to assign P first, and assign it the value 6. What are the resulting domains after enforcing unary constraints (from part i) and running forward checking for this assignment?

P					6
B	X	2	X	4	X
C	1	2	3	4	X
K	1	2	3	4	X
I	X	2	3	4	X
M	X	2	X	4	X

- (iv) [3 pts] **Iterative Improvement** Instead of running backtracking search, you decide to start over and run iterative improvement with the min-conflicts heuristic for value selection. Starting with the following assignment:

P:6, B:4, C:3, K:2, I:1, M:5

First, for each variable write down how many constraints it violates in the table below.

Then, in the table on the right, for all variables that could be selected for assignment, put an x in any box that corresponds to a possible value that could be assigned to that variable according to min-conflicts. When marking next values a variable could take on, only mark values different from the current one.

Variable	# violated
P	0
B	0
C	1
K	0
I	2
M	0

	1	2	3	4	5	6
P						
B						
C						
K						
I						
M						

(b) Variable ordering

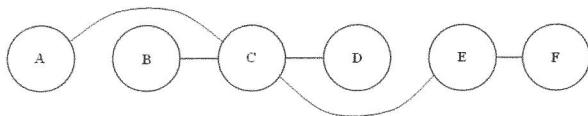
We say that a variable X is backtracked if, after a value has been assigned to X, the recursion returns at X without a solution, and a different value must be assigned to X.

For this problem, consider the following three algorithms:

1. Run backtracking search with no filtering
2. Initially enforce arc consistency, then run backtracking search with no filtering
3. Initially enforce arc consistency, then run backtracking search while enforcing arc consistency after each assignment

(i) [5 pts]

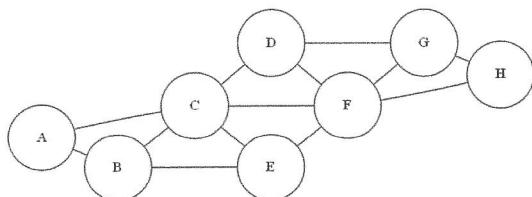
For each algorithm, circle all orderings of variable assignments that guarantee that no backtracking will be necessary when finding a solution to the CSP represented by the following constraint graph.



Algorithm 1	Algorithm 2	Algorithm 3
A-B-C-D-E-F	A-B-C-D-E-F	A-B-C-D-E-F
F-E-D-C-B-A	F-E-D-C-B-A	F-E-D-C-B-A
C-A-B-D-E-F	C-A-B-D-E-F	C-A-B-D-E-F
B-D-A-F-E-C	B-D-A-F-E-C	B-D-A-F-E-C
D-E-F-C-B-A	D-E-F-C-B-A	D-E-F-C-B-A
B-C-D-A-E-F	B-C-D-A-E-F	B-C-D-A-E-F

(ii) [5 pts]

For each algorithm, circle all orderings of variable assignments that guarantee that no more than two variables will be backtracked when finding a solution to the CSP represented by the following constraint graph.

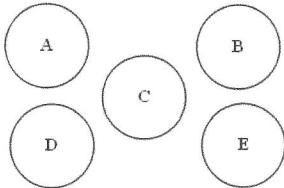


Algorithm 1	Algorithm 2	Algorithm 3
C-F-A-B-E-D-G-H	C-F-A-B-E-D-G-H	C-F-A-B-E-D-G-H
F-C-A-H-E-B-D-G	F-C-A-H-E-B-D-G	F-C-A-H-E-B-D-G
A-B-C-E-D-F-G-H	A-B-C-E-D-F-G-H	A-B-C-E-D-F-G-H
G-C-H-F-B-D-E-A	G-C-H-F-B-D-E-A	G-C-H-F-B-D-E-A
A-B-E-D-G-H-C-F	A-B-E-D-G-H-C-F	A-B-E-D-G-H-C-F
A-D-B-G-E-H-C-F	A-D-B-G-E-H-C-F	A-D-B-G-E-H-C-F

- (c) **All Satisfying Assignments** Now consider a modified CSP in which we wish to find every possible satisfying assignment, rather than just one such assignment as in normal CSPs. In order to solve this new problem, consider a new algorithm which is the same as the normal backtracking search algorithm, except that when it sees a solution, instead of returning it, the solution gets added to a list, and the algorithm backtracks. Once there are no variables remaining to backtrack on, the algorithm returns the list of solutions it has found.

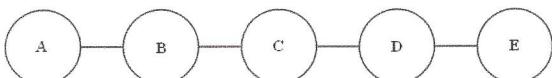
For each graph below, select whether or not using the MRV and/or LCV heuristics could affect the number of nodes expanded in the search tree in this new situation.

(i) [2 pts]



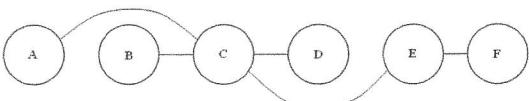
- Neither MRV nor LCV can have an effect.
- Only MRV can have an effect.
- Only LCV can have an effect .
- Both MRV and LCV can have an effect.

(ii) [2 pts]



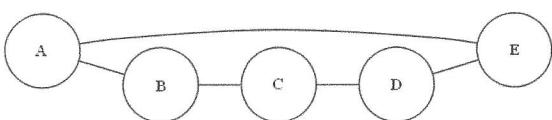
- Neither MRV nor LCV can have an effect.
- Only MRV can have an effect.
- Only LCV can have an effect .
- Both MRV and LCV can have an effect.

(iii) [2 pts]



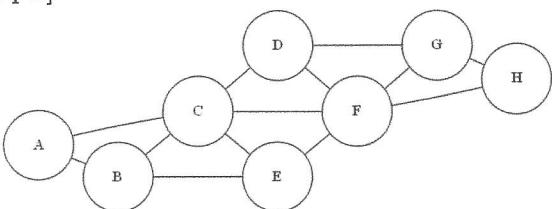
- Neither MRV nor LCV can have an effect.
- Only MRV can have an effect.
- Only LCV can have an effect .
- Both MRV and LCV can have an effect.

(iv) [2 pts]



- Neither MRV nor LCV can have an effect.
- Only MRV can have an effect.
- Only LCV can have an effect .
- Both MRV and LCV can have an effect.

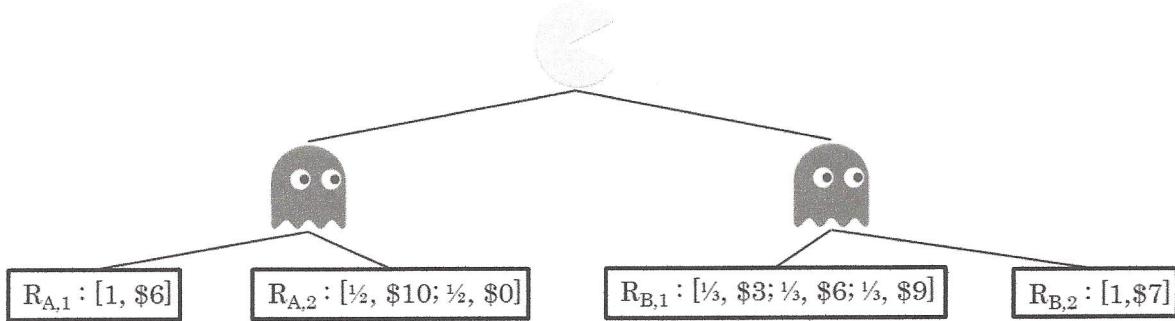
(v) [2 pts]



- Neither MRV nor LCV can have an effect.
- Only MRV can have an effect.
- Only LCV can have an effect .
- Both MRV and LCV can have an effect.

Q4. [10 pts] Utilities

Pacman is buying a raffle ticket from the ghost raffle ticket vendor. There are two ticket types: A and B , but there are multiple specific tickets of each type. Pacman picks a ticket type, but the ghost will then choose which specific ticket Pacman will receive. Pacman's utility for a given raffle ticket is equal to the utility of the lottery of outcomes for that raffle ticket. For example, ticket $R_{A,2}$ corresponds to a lottery with equal chances of yielding 10 and 0, and so $U(R_{A,2}) = U([\frac{1}{2}, 10; \frac{1}{2}, 0])$. Pacman, being a rational agent, wants to maximize his expected utility, but the ghost may have other goals! The outcomes are illustrated below.



- (a) Imagine that Pacman's utility for money is $U(m) = m$.

- (i) [2 pts] What are the utilities to Pacman of each raffle ticket?

$$U(R_{A,1}) = 6$$

$$U(R_{A,2}) = 5$$

$$U(R_{B,1}) = 6$$

$$U(R_{B,2}) = 7$$

- (ii) [1 pt] Which raffle ticket will Pacman receive under optimal play if the ghost is trying to minimize Pacman's utility (and Pacman knows the ghost is doing so)? (circle one)

$R_{A,1}$

$R_{A,2}$

$R_{B,1}$

$R_{B,2}$

- (iii) [1 pt] What is the equivalent monetary value of raffle ticket $R_{B,1}$?

\$6

- (b) Now imagine that Pacman's utility for money is given by $U(m) = m^2$.

- (i) [2 pts] What are the utilities to Pacman of each raffle ticket?

$$U(R_{A,1}) = 36$$

$$U(R_{A,2}) = 50$$

$$U(R_{B,1}) = \frac{9+36+81}{3} = 42$$

$$U(R_{B,2}) = 49$$

- (ii) [1 pt] The ghost is still trying to minimize Pacman's utility, but the ghost mistakenly thinks that Pacman's utility is given by $U(m) = m$, and Pacman is aware of this flaw in the ghost's model. Which raffle ticket will Pacman receive? (circle one)

$R_{A,1}$

$R_{A,2}$

$R_{B,1}$

$R_{B,2}$

- (iii) [1 pt] What is the equivalent monetary value of raffle ticket $R_{B,1}$? (you may leave your answer as an expression)

$\sqrt{42}$

- (c) [2 pts] Pacman has the raffle with distribution $[0.5, \$100; 0.5, \$0]$. A ghost insurance dealer offers Pacman an insurance policy where Pacman will get \$100 regardless of what the outcome of the ticket is. If Pacman's utility for money is $U(m) = \sqrt{m}$, what is the maximum amount of money Pacman would pay for this insurance?

$$U(100 - c) = \sqrt{100 - c} = 0.5\sqrt{100} = 5$$

$$c = 75$$

Q5. [9 pts] Games: Three-Player Cookie Pruning

Three of your TAs, Alvin, Sergey, and James rent a cookie shuffler, which takes in a set number of cookies and groups them into 3 batches, one for each player. The cookie shuffler has three levers (with positions either UP or DOWN), which act to control how the cookies are distributed among the three players. Assume that 30 cookies are initially put into the shuffler.

Each player controls one lever, and they act in turn. Alvin goes first, followed by Sergey, and finally James. Assume that all players are able to calculate the payoffs for every player at the terminal nodes. Assume the payoffs at the leaves correspond to the number of cookies for each player in their corresponding turn order. Hence, an utility of (7,10,13) corresponds to Alvin getting 7 cookies, Sergey getting 10 cookies, and James getting 13 cookies. No cookies are lost in the process, so the sum of cookies of all three players must equal the number of cookies put into the shuffler. Players want to maximize their own number of cookies.

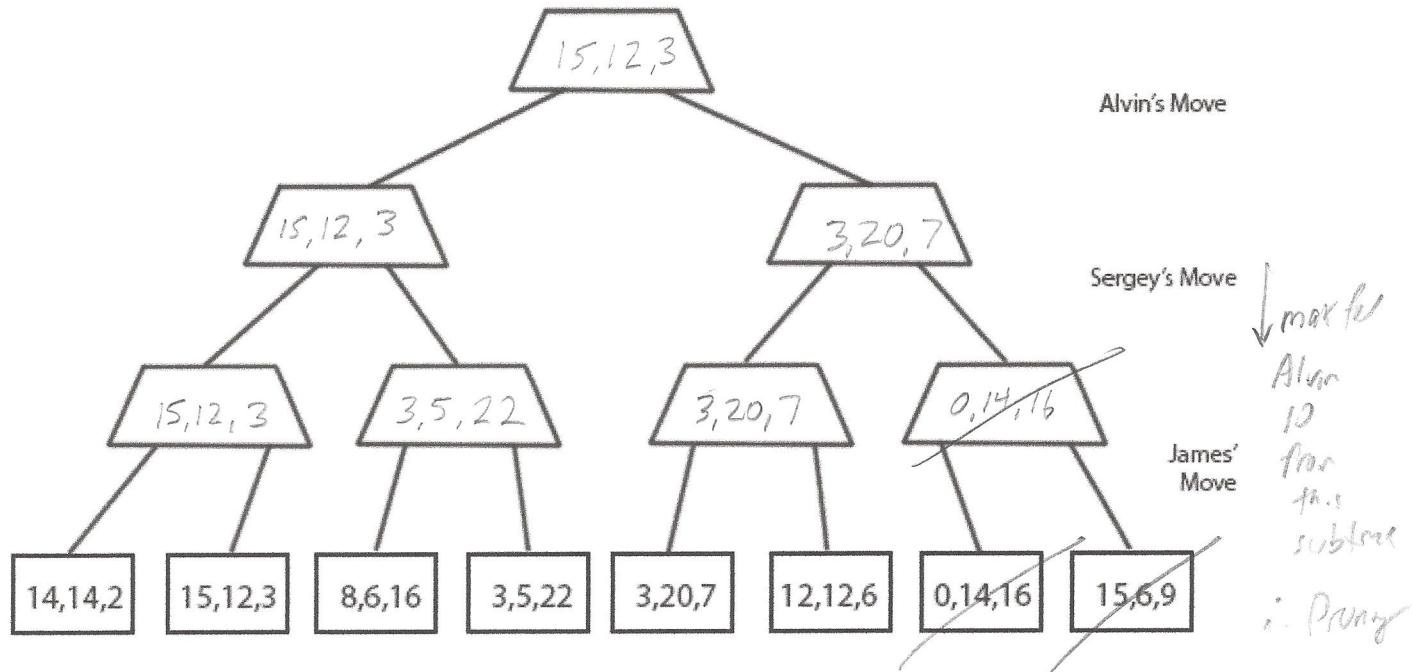
- (a) [3 pts] What is the utility triple propagated up to the root?

$15, 12, 3$

- (b) [6 pts] Is pruning possible in this game? Fill in "Yes" or "No". If yes, cross out all nodes (both leaves and intermediate nodes) that get pruned. If no, explain in one sentence why pruning is not possible. Assume the tree traversal goes from left to right.

Yes.

No, Reasoning:



Q6. [10 pts] The nature of discounting

Pacman is stuck in a friendlier maze where he gets a reward every time he visits state $(0,0)$. This setup is a bit different from the one you've seen before: Pacman can get the reward multiple times; these rewards do not get "used up" like food pellets and there are no "living rewards". As usual, Pacman can not move through walls and may take any of the following actions: go North (\uparrow), South (\downarrow), East (\rightarrow), West (\leftarrow), or stay in place (\circ). State $(0,0)$ gives a total reward of 1 every time Pacman takes an action in that state regardless of the outcome, and all other states give no reward.

You should not need to use any other complicated algorithm/calculations to answer the questions below. We remind you that geometric series converge as follows: $1 + \gamma + \gamma^2 + \dots = 1/(1 - \gamma)$.

- (a) [2 pts] Assume finite horizon of $h = 10$ (so Pacman takes exactly 10 steps) and no discounting ($\gamma = 1$).

Fill in an optimal policy:

0 1 2

0	\circ	\leftarrow	\leftarrow
1	\uparrow	\downarrow	\downarrow
2	\uparrow	\leftarrow	\leftarrow

(available actions: $\uparrow, \downarrow, \rightarrow, \leftarrow, \circ$)

Fill in the value function:

0 1 2

0	10	9	8
1	9	6	5
2	8	7	6

- (b) The following Q-values correspond to the value function you specified above.

$$Q = R + \sqrt{\gamma} (9 \text{ steps})$$

- (i) [1 pt] The Q value of state-action $(0,0)$, (East) is: 9
(ii) [1 pt] The Q value of state-action $(1,1)$, (East) is: 4

- (c) Assume finite horizon of $h = 10$, no discounting, but the action to stay in place is temporarily (for this sub-point only) unavailable. Actions that would make Pacman hit a wall are not available. Specifically, Pacman can not use actions North or West to remain in state $(0,0)$ once he is there.

- (i) [1 pt] [true or false] There is just one optimal action at state $(0,0)$ E and S
(ii) [1 pt] The value of state $(0,0)$ is: 5

- (d) [2 pts] Assume infinite horizon, discount factor $\gamma = 0.9$.

The value of state $(0,0)$ is: 10

$$\frac{1}{1-0.9} = 10$$

- (e) [2 pts] Assume infinite horizon and no discount ($\gamma = 1$). At every time step, after Pacman takes an action and collects his reward, a power outage could suddenly end the game with probability $\alpha = 0.1$.

The value of state $(0,0)$ is: 10

$$\frac{1}{\alpha} = 10$$

Q7. [10 pts] The Value of Games

Pacman is the model of rationality and seeks to maximize his expected utility, but that doesn't mean he never plays games.

- (a) [3 pts] **Q-Learning to Play under a Conspiracy.** Pacman does tabular Q-learning (where every state-action pair has its own Q-value) to figure out how to play a game against the adversarial ghosts. As he likes to explore, Pacman always plays a random action. After enough time has passed, every state-action pair is visited infinitely often. The learning rate decreases as needed. For any game state s , the value $\max_a Q(s, a)$ for the learned $Q(s, a)$ is equal to (for complete search trees)

- The minimax value where Pacman maximizes and ghosts minimize.
- The expectimax value where Pacman maximizes and ghosts act uniformly at random.
- The expectimax value where Pacman plays uniformly at random and ghosts minimize.
- The expectimax value where both Pacman and ghosts play uniformly at random.
- None of the above.

- (b) [3 pts] **Feature-based Q-Learning the Game under a Conspiracy.** Pacman now runs feature-based Q-learning. The Q-values are equal to the evaluation function $\sum_{i=1}^n w_i f_i(s, a)$ for weights w and features f . The number of features is much less than the number of states. As he likes to explore, Pacman always plays a random action. After enough time has passed, every state-action pair is visited infinitely often. The learning rate decreases as needed. The value $\max_a Q(s, a)$ for the learned $Q(s, a)$ is equal to (for complete search trees)

- The minimax value where Pacman maximizes and ghosts minimize and the same evaluation function is used at the leaves.
- The expectimax value where Pacman maximizes and ghosts act uniformly at random and the same evaluation function is used at the leaves.
- The expectimax value where Pacman plays uniformly at random and ghosts minimize and the same evaluation function is used at the leaves.
- The expectimax value where both Pacman and ghosts play uniformly at random and the same evaluation function is used at the leaves.
- None of the above.

- (c) [2 pts] **A Costly Game.** Pacman is now stuck playing a new game with only costs and no payoff. Instead of maximizing expected utility $V(s)$, he has to minimize expected costs $J(s)$. In place of a reward function, there is a cost function $C(s, a, s')$ for transitions from s to s' by action a . We denote the discount factor by $\gamma \in (0, 1)$. $J^*(s)$ is the expected cost incurred by the optimal policy. Which one of the following equations is satisfied by J^* ?

- $J^*(s) = \min_a \sum_{s'} [C(s, a, s') + \gamma \max_{a'} T(s, a', s') * J^*(s')]$
- $J^*(s) = \min_{s'} \sum_a T(s, a, s')[C(s, a, s') + \gamma * J^*(s')] \quad -$
- $J^*(s) = \min_a \sum_{s'} T(s, a, s')[C(s, a, s') + \gamma * \max_{s'} J^*(s')] \quad -$
- $J^*(s) = \min_{s'} \sum_a T(s, a, s')[C(s, a, s') + \gamma * \max_{s'} J^*(s')] \quad -$
- $J^*(s) = \min_a \sum_{s'} T(s, a, s')[C(s, a, s') + \gamma * J^*(s')] \quad \rightarrow Q \text{ learning min}$
- $J^*(s) = \min_{s'} \sum_a [C(s, a, s') + \gamma * J^*(s')]$

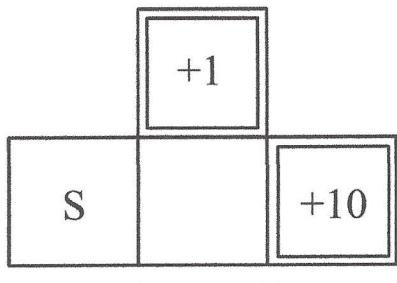
- (d) [2 pts] **It's a conspiracy again!** The ghosts have rigged the costly game so that once Pacman takes an action they can pick the outcome from all states $s' \in S'(s, a)$, the set of all s' with non-zero probability according to $T(s, a, s')$. Choose the correct Bellman-style equation for Pacman against the adversarial ghosts.

- $J^*(s) = \min_a \max_{s'} T(s, a, s')[C(s, a, s') + \gamma * J^*(s')]$
- $J^*(s) = \min_{s'} \sum_a T(s, a, s')[\max_{s'} C(s, a, s') + \gamma * J^*(s')]$
- $J^*(s) = \min_a \min_{s'} [C(s, a, s') + \gamma * \max_{s'} J^*(s')]$
- $J^*(s) = \min_a \max_{s'} [C(s, a, s') + \gamma * J^*(s')] \quad \rightarrow M_{non-min}$
- $J^*(s) = \min_{s'} \sum_a T(s, a, s')[\max_{s'} C(s, a, s') + \gamma * \max_{s'} J^*(s')]$
- $J^*(s) = \min_a \min_{s'} T(s, a, s')[C(s, a, s') + \gamma * J^*(s')]$

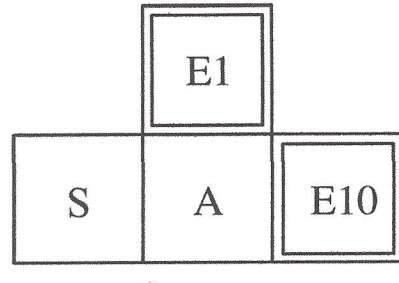
Q8. [16 pts] Infinite Time to Study

Pacman lives in a calm gridworld. S is the start state and double-squares are exit states. In exits, the only action available is exit, which earns the associated reward and transitions to a terminal state X (not shown). In normal states, the actions are to move to neighboring squares (for example, S has the single action \rightarrow) and they always succeed. There is no living reward, so all non-exit actions have reward 0.

Throughout the problem the discount $\gamma = 1$.



The calmworld



State names

- (a) [2 pts] What are the optimal values of S and A?

$$V^*(S) = 10$$

$$V^*(A) = 10$$

Pacman doesn't know the details of this gridworld so he does Q-learning with a learning rate of 0.5 and all Q-values initialized to 0 to figure it out.

Consider the following sequence of transitions in the calmworld:

s	a	s'	r
S	\rightarrow	A	0
A	\uparrow	E1	0
E1	exit	X	1
S	\rightarrow	A	0
A	\rightarrow	E10	0
E10	exit	X	10

- (b) [2 pts] Circle the Q-values that are non-zero after these episodes.

$$Q(S, \rightarrow) \quad Q(A, \uparrow) \quad Q(A, \rightarrow) \quad Q(E1, \text{exit}) \quad Q(E10, \text{exit})$$

- (c) [2 pts] What do the Q-values converge to if these episodes are repeated infinitely with a constant learning rate of 0.5? Write *none* if they do not converge.

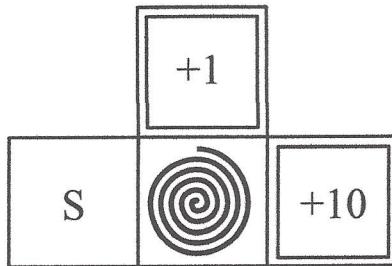
$$Q(S, \rightarrow) = 10$$

$$Q(A, \leftarrow) = 0$$

$$Q(A, \uparrow) = |$$

(Q-learning details reminder: assume $\alpha = 0.5$ and the Q-values are initialized to 0.)

It's vortex season in the gridworld. In the vortex state the only action is escape, which delivers Pacman to a neighboring state uniformly at random.



The vortexworld

(d) [2 pts] What are the optimal values of S and A in the vortex gridworld?

$$V^*(S) = 5.5$$

$$V^*(A) = 5.5$$

Consider the following sequences of transitions in the vortexworld:

S1				S2			
s	a	s'	r	s	a	s'	r
S	\rightarrow	A	0	S	\rightarrow	A	0
A	escape	E1	0	A	escape	E1	0
E1	exit	X	1	E1	exit	X	1
S	\rightarrow	A	0	S	\rightarrow	A	0
A	escape	E10	0	A	escape	E10	0
E10	exit	X	10	E10	exit	X	10

(e) [2 pts] What do the Q-values converge to if the sequence S1 is repeated infinitely with appropriately decreasing learning rate? Write *never* if they do not converge.

$$Q^{S1}(S, \rightarrow) = 5.5 \quad Q^{S1}(A, \text{escape}) = 5.5$$

(f) [2 pts] What if the sequence S2 is repeated instead?

$$Q^{S2}(S, \rightarrow) = 7 \quad Q^{S2}(A, \text{escape}) = 7$$

(g) [2 pts] Which is the true optimum $Q^*(S, \rightarrow)$ in the vortex gridworld? Circle the answer.

$Q^{S1}(S, \rightarrow)$ $Q^{S2}(S, \rightarrow)$ other

(h) [2 pts] Q-learning with constant $\alpha = 1$ and visiting state-actions infinitely often converges

in calmworld in vortexworld in neither world