

- You have approximately 80 minutes.
- The exam is closed book, closed calculator, and closed notes except your one-page crib sheet.
- Mark your answers ON THE EXAM ITSELF. If you are not sure of your answer you may wish to provide a *brief* explanation. All short answer sections can be successfully answered in a few sentences AT MOST.
- For multiple choice questions,
 - means mark **all options** that apply
 - means mark a **single choice**
 - When selecting an answer, please fill in the bubble or square **completely** (● and ■)

First name	
Last name	
SID	
Student to your right	
Student to your left	

Your Discussion TA (fill all that apply):

- | | |
|--|--|
| <input type="checkbox"/> Caryn (TTh 10:00am Wheeler) | <input type="checkbox"/> Arin (TTh 10:00am Dwinelle) |
| <input type="checkbox"/> Bobby (MW 3:00 pm) | <input type="checkbox"/> Mike(TTh 11:00am) |
| <input type="checkbox"/> Benson (MW 4:00 pm) | <input type="checkbox"/> Do not attend any |

For staff use only:

Q1. They See Me Rolling (Search Problem)	/21
Q2. Search Algorithms Potpourri	/19
Q3. MDP Solvers	/28
Q4. Model-Free RL	/14
Q5. Game Trees	/18
Total	/100

THIS PAGE IS INTENTIONALLY LEFT BLANK

Q1. [21 pts] They See Me Rolling (Search Problem)

Pacman buys a car to start Rolling in the Pac-City! But driving a car requires a new set of controls because he can now travel faster than 1 grid per turn (second). Instead of solely moving [North, South, East, West, Stop], Pacman's car has two distinct integers to control: throttle, and steering.

Throttle: $t_i \in \{1, 0, -1\}$, corresponding to {Gas, Coast, Brake}. This controls the **speed** of the car by determining its acceleration. The integer chosen here will be added to his velocity for the next state. For example, if Pacman is currently driving at 5 grid/s and chooses Gas he will be traveling at 6 grid/s in the next turn.

Steering: $s_i \in \{1, 0, -1\}$, corresponding to {Turn Left, Neutral, Turn Right}. This controls the **direction** of the car. For example, if he is facing North and chooses Turn Left he will be facing West in the next turn.

- (a) [11 pts] Suppose Pac-city has dimension m by n , but only $k < mn$ squares are legal roads. The speed limit of Pac-city is 3 grid/s. For this sub-part only, suppose Pacman is a law-abiding citizen, so $0 \leq v \leq 3$ at all time, and he only drives on legal roads.

- (i) [3 pts] Without any additional information, what is the **tightest upper bound** on the size of state space, if he wants to search a route (not necessarily the shortest) from his current location to anywhere in the city. Please note that your state space representation must be able to represent **all** states in the search space.

$4mn$
 $4k$
 $12mn$
 $12k$
 $16mn$
 $16k$
 $48mn$
 $48k$

Only legal grids count, so there are k legal position. At each legal position, there are 4 possible speed (0, 1, 2, 3), so a factor of 4 is multiplied. In addition, since change of direction depends on orientation of the car, another factor of 4 is multiplied. The size of state space is bounded by $k * 4 * 4 = 16k$

- (ii) [3 pts] What is the maximum branching factor? The answer should be an integer.

9

3 possible throttle inputs, and 3 possible steering inputs.

- (iii) [3 pts] Which algorithm(s) is/are guaranteed to return a path between two points, if one exists?

Depth First Tree Search
 Breadth First Tree Search
 Depth First Graph Search
 Breadth First Graph Search

- (iv) [2 pts] Is Breadth First Graph Search guaranteed to return the path with the shortest grid distance?

Yes
 No

The Breadth First Graph Search is guranteed to return the path with the shortest amount of time, because each edge here represent moving for 1 unit of time.

- (b) [5 pts] Now let's remove the constraint that Pacman follows the speed limit. Now Pacman's speed is limited by the mechanical constraints of the car, which is 6 grid/s, double the speed limit.

Pacman is now able to drive twice as fast on the route to his destination. How do the following properties of the search problem change as a result of being able to drive twice as fast?

- (i) [1 pt] Size of State Space:

Increases
 Stays the same
 Decreases

At each legal position, there are 7 possible speed (0,1, 2, 3, 4, 5, 6), so a factor of 7 is multiplied. size of state space is now $k * 7 * 4 = 28k$

- (ii) [1 pt] Maximum Branching Factor:

Increases
 Stays the same
 Decreases

Branching factor is independent of top speed

For the following part, **mark all choices that could happen on any graph**

- (iii) [3 pts] The number of nodes expanded with **Depth** First Graph Search:

Increases
 Stays the same
 Decreases

Anything can happen with Depth First Graph Search. For example, too fast can jump past an intersection that would lead to a shorter path to the goal

- (c) [5 pts] Now we need to consider that there are $p > 0$ police cars waiting at $p > 0$ distinct locations trying to catch Pacman riding dirty!! All police cars are stationary, but once Pacman takes an action which lands him in the same grid as one police car, Pacman will be arrested and the game ends.

Pacman wants to find a route to his destination, without being arrested. How do the following properties of the search problem change as a result of avoiding the police?

(i) [1 pt] Size of State Space:

- Increases Stays the same Decreases

The size of statespace is still the same because all the factors in state space calculation is the same

(ii) [1 pt] Maximum Branching Factor:

- Increases Stays the same Decreases

Branching factor is independent of police presense

For the following part, **mark all choices that could happen on any graph**

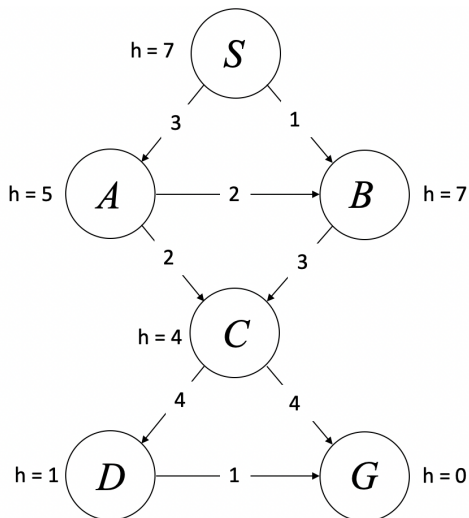
(iii) [3 pts] Number of nodes expanded with **Breadth** First Graph Search:

- Increases Stays the same Decreases

Again anything could happen with BFS. Policemen could block out misleading paths to decrease the number of expansion. They could also not affect anything. They could also block the closest goal causing you to explore more nodes.

Q2. [19 pts] Search Algorithms Potpourri

- (a) [8 pts] We will investigate various search algorithms for the following graph. Edges are labeled with their costs, and heuristic values h for states are labeled next to the states. S is the start state, and G is the goal state. In all search algorithms, assume ties are broken in alphabetical order.



- (i) [2 pts] Select all boxes that describe the given heuristic values.

admissible consistent Neither

- (ii) [2 pts] Given the above heuristics, what is the order that the states are going to be expanded in, assuming we run A* graph search with the heuristic values provided.

Index	1	2	3	4	5	Not Expanded
S	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
A	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
B	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
C	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>
D	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>
G	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>

- (iii) [1 pt] Assuming we run A* graph search with the heuristic values provided, what path is returned?

$S \rightarrow A \rightarrow B \rightarrow C \rightarrow D \rightarrow G$
 $S \rightarrow A \rightarrow C \rightarrow G$
 $S \rightarrow A \rightarrow C \rightarrow D \rightarrow G$
 $S \rightarrow B \rightarrow C \rightarrow G$
 $S \rightarrow A \rightarrow C \rightarrow D \rightarrow G$
 $S \rightarrow A \rightarrow C \rightarrow D \rightarrow G$
 $S \rightarrow A \rightarrow B \rightarrow C \rightarrow G$
 None of the above

- (iv) [2 pts] Given the above heuristics, what is the order that the states are going to be expanded in, assuming we run greedy graph search with the heuristic values provided.

Index	1	2	3	4	5	Not Expanded
S	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
A	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
B	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>
C	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
D	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>
G	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>

- (v) [1 pt] What path is returned by greedy graph search?

$S \rightarrow A \rightarrow B \rightarrow C \rightarrow D \rightarrow G$
 $S \rightarrow A \rightarrow C \rightarrow G$
 $S \rightarrow A \rightarrow C \rightarrow D \rightarrow G$
 $S \rightarrow A \rightarrow C \rightarrow D \rightarrow G$
 $S \rightarrow A \rightarrow C \rightarrow D \rightarrow G$
 None of the above

- (b) [4 pts] Consider a complete graph, K_n , the undirected graph with n vertices where all n vertices are connected (there is an edge between every pair of vertices), resulting in $\binom{n}{2}$ edges. Please select the maximum possible depth of the resulting tree when the following **graph** search algorithms are run (assume any possible start and goal vertices).

	1	$\lceil \frac{n}{2} \rceil$	$n - 1$	$\binom{n}{2}$	None of the above
BFS	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
DFS	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>

- (c) [3 pts] Given two admissible heuristics h_A and h_B .

- (i) [2 pts] Which of the following are guaranteed to also be admissible heuristics?

$h_A + h_B$
 $\frac{1}{2}(h_A)$
 $\frac{1}{2}(h_B)$
 $\frac{1}{2}(h_A + h_B)$
 $h_A * h_B$
 $max(h_A, h_B)$
 $min(h_A, h_B)$

- (ii) [1 pt] Consider performing A* **tree** search. Which is generally best to use if we want to expand the fewest number of nodes? **Note this was changed from graph to tree search during the exam**

$h_A + h_B$
 $\frac{1}{2}(h_A)$
 $\frac{1}{2}(h_B)$
 $\frac{1}{2}(h_A + h_B)$
 $h_A * h_B$
 $max(h_A, h_B)$
 $min(h_A, h_B)$

- (d) [4 pts] Consider performing tree search for some search graph. Let $depth(n)$ be the depth of search node n and $cost(n)$ be the total cost from the start state to node n . Let G_d be a goal node with minimum depth, and G_c be a goal node with minimum total cost.

- (i) [2 pts] For iterative deepening (where we repeatedly run DFS and increase the maximum depth allowed by 1), mark all conditions that are guaranteed to be true for every node n that could be expanded during the search, or mark "None of the above" if none of the conditions are guaranteed.

$cost(n) \leq cost(G_c)$
 $cost(n) \leq cost(G_d)$
 $depth(n) \leq depth(G_c)$
 $depth(n) \leq depth(G_d)$
 None of the above

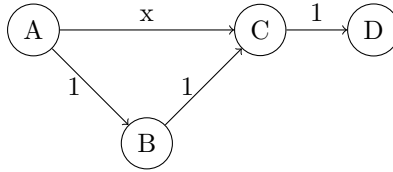
When running iterative deepening we will explore all nodes of depth k , before we explore any nodes of depth $k + 1$. As a result will never explore any nodes that have depth greater than G_d because we would stop exploring once we reached G_d . This also means we would never explore any nodes with depth greater than G_c because G_c has to have depth greater than or equal to the minimum depth goal, G_d

- (ii) [2 pts] What is necessarily true regarding iterative deepening on any search tree?

Complete as opposed to DFS tree search
 Strictly faster than DFS tree search
 Strictly faster than BFS tree search
 More memory efficient than BFS tree search
 A type of stochastic local search
 None of the above

Q3. [28 pts] MDP Solvers

- (a) [5 pts] Consider the following finite-state MDP, with states A , B , C , and D . Each edge indicates an action between states, and all transitions are deterministic. The edge weights denote the transition rewards. Assume $\gamma = 1$ for all parts.

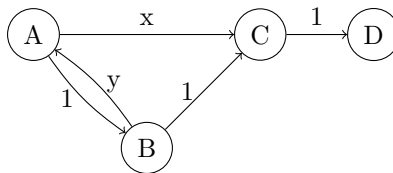


For each of the states, let k be the **first** iteration in Value Iteration where the **values** could have converged for some value of x (in other words, the smallest k such that $V_k(s) = V^*(s)$). List **all** possible k for each state. Mark $k = \infty$ if the value for the state never converges.

State	k
A	2, 3
B	2
C	1
D	0

C will find its optimal value in one iteration. B will find its optimal value one iteration after that (2 iterations total). A will find its optimal value one iteration after C (2 iterations total) or one iteration after B (3 iterations total), depending on the value of x .

- (b) [11 pts] Consider the following finite-state MDP, with states A , B , C , and D . Each edge indicates an action between states, and all transitions are deterministic. Let $x, y \geq 0$. The edge weights denote the transition rewards. Assume $\gamma = 1$ for all parts.



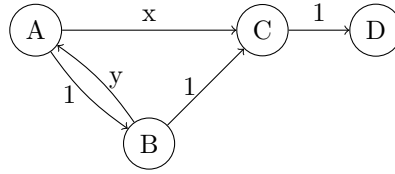
- (i) [5 pts]

For each of the states, let k be the **first** iteration in Value Iteration where the values could have converged for some nonnegative **value** of x or y (in other words, the smallest k such that $V_k(s) = V^*(s)$). Find all possible k for each state. Mark $k = \infty$ if the value for the state never converges.

State	k
A	∞
B	∞
C	1
D	0

A and B will never find their optimal value because they can get infinite value. C and D are the same as above.

- (ii) [6 pts] Additionally, assume $x + y < 1$. For each of the states, let k be the **first** iteration in Policy Iteration where the **policy** is optimal for a particular state (in other words, the smallest k such that $\pi_k(s) = \pi^*(s)$). Find k for each state. Write $k = \infty$ if the policy for the state never converges. Assume tie-breaking during policy improvement is alphabetical. (The graph is redrawn below for convenience)



Let π_0 be the following:

State (s)	Action ($\pi_0(s)$)	State	k
A	C	A	1
B	C	B	2
C	D	C	0
D	D	D	0

First, evaluate π_0 :

$$\begin{aligned}
 V^{\pi_0}(D) &= 0 \\
 \implies V^{\pi_0}(C) &= 1 + V^{\pi_0}(D) = 1 \\
 \implies V^{\pi_0}(B) &= 1 + V^{\pi_0}(C) = 2 \\
 \implies V^{\pi_0}(A) &= x + V^{\pi_0}(C) = x + 1
 \end{aligned}$$

Now do policy improvement to obtain π_1 :

$$\begin{aligned}
 \pi_1(D) &= D \\
 \pi_1(C) &= D \\
 \pi_1(B) &= \arg \max_{\{A, C\}} \{A : x + y + 1, C : 2\} = C \\
 \pi_1(A) &= \arg \max_{\{B, C\}} \{B : 2, C : x + 1\} = B
 \end{aligned}$$

Now, evaluate π_1 :

$$\begin{aligned}
 V^{\pi_1}(D) &= 0 \\
 \implies V^{\pi_1}(C) &= 1 + V^{\pi_1}(D) = 1 \\
 \implies V^{\pi_1}(B) &= 1 + V^{\pi_1}(C) = 2 \\
 \implies V^{\pi_1}(A) &= x + V^{\pi_1}(B) = x + 2
 \end{aligned}$$

Now run policy improvement to obtain π_2 :

$$\begin{aligned}
 \pi_2(D) &= D \\
 \pi_2(C) &= D \\
 \pi_2(B) &= \arg \max_{\{A, C\}} \{A : y + 3, C : 2\} = A \\
 \pi_2(A) &= \arg \max_{\{B, C\}} \{B : 3, C : x + 1\} = B
 \end{aligned}$$

Observe that this policy is optimal, because the value $V^{\pi_2}(A) = V^{\pi_2}(B) = +\infty$. The other values are trivially optimal because the agent has only one choice of action.

- (c) (i) [2 pts] Mark **all** of the statements that must be true for any MDP.

For no state s and for all policies π , $V^*(s) \geq V^\pi(s)$

- For some state s and some policy π , $V^*(s) \geq V^\pi(s)$
- For all states s and all policies π , $V^*(s) \geq V^\pi(s)$
- None of the above

(ii) [2 pts] Mark **all** of the statements that are true for value iteration.

- Each iteration of value iteration produces a value function that has higher value than the prior value functions **for all states**.
- Each iteration of value iteration produces a value function that has value at least as good as the prior value functions **for all states**.
- Value iteration can produce a value function that has lower value than the earlier value functions **for some state**.
- At convergence, the value iteration does not change the value function **for any state**.
- None of the above

(d) [8 pts] Consider an MDP alternating game, where Pacman takes turns with Ghost to stochastically take actions. Pacman takes actions to maximize the expected score, whereas Ghost takes actions to minimize the expected score. Both know the others strategy.

Let $V^*(s)$ be the value function for Pacman. Pacman and Ghost have identical transition, reward functions, and action space. (Discounting is applied for Pacmans actions)

Derive the new Bellman equation for $V^*(s)$ in terms of γ , R , T , and V^* :

$$V^*(s) = \boxed{\begin{array}{|c|c|c|c|c|c|c|c|} \hline \text{(i)} & \text{(ii)} & \text{(iii)} & \text{(iv)} & \text{(v)} & \text{(vi)} & \text{(vii)} & \text{(viii)} \\ \hline \end{array}}$$

For each blank (i) through (iv), mark the appropriate subexpression. If it is possible to write the expression for $V^*(s)$ without a particular sub-expression, mark "None".

(i) [1 pt] $\frac{1}{2}$ $\frac{1}{3}$ 2 -1 3 γ None

(ii) [1 pt] \max_a \min_a None

(iii) [1 pt] $\sum_{s'}$ $\sum_{s'} T(s, a, s')$ None

(iv) [1 pt] [$[R(s, a, s')+$ $[R(s, a, s')-$ None

(v) [1 pt] \max_a \min_a $\max_{a'}$ $\min_{a'}$ None

(vi) [1 pt] $\sum_{s'}$ $\sum_{s'} T(s, a, s')$ $\sum_{s''}$ $\sum_{s''} T(s', a', s'')$ None

(vii) [1 pt] [$[R(s, a, s')+$ $[R(s, a, s')-$
 $[R(s', a', s'')+$ $[R(s', a', s'')-$ None

(viii) [1 pt] $\frac{1}{2}V^*(s')]$ $\gamma V^*(s')]$ $V^*(s')]$ $\frac{1}{2}V^*(s'')]$
 $\gamma V^*(s'')]$ $V^*(s'')]$ None

$$\max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \min_{a'} \sum_{s''} T(s', a', s'') [R(s', a', s'') + \gamma V^*(s'')]]$$

Q4. [14 pts] Model-Free RL

(a) [8 pts]

Consider the following MDP with state space $\mathcal{S} = \{A, B, C, D, E, F\}$ and action space $\mathcal{A} = \{left, right, up, down, stay\}$. However, we do not know the transition dynamics or reward function (we do not know what the resulting next state and reward are after applying an action in a state).

A	B	C
D	E	F

- (i) [4 pts] We will run Q-learning to try to learn a policy that maximizes the expected reward in the MDP, with discount factor $\gamma = 1$ and learning rate $\alpha = 0.5$. Suppose we observe the following transitions in the environment. After observing each transition, we update the Q function, which is initially 0. Fill in the blanks with the corresponding values of the Q function after these updates.

Episode Number	State	Action	Reward	Next State
1	A	right	2	B
2	B	right	10	C
3	B	right	-3	E
4	C	stay	-5	A

State	Action	Q(state, action)
A	stay	0
A	right	1
B	right	1
C	stay	-2

- (ii) [4 pts] We now observe the following samples:

Episode Number	State	Action	Reward	Next State
5	B	down	-5	E
6	B	down	-5	E
7	B	up	-2	B
8	B	left	-8	B
9	B	left	-8	A
10	B	stay	-6	B
11	A	right	1	B

Continuing with the Q function from the last subproblem, we train the Q function on these samples. What is $Q(A, right)$?

1.5

The only transition that matters is the last one, because the $Q(B, right)$ is untouched and largest after these negative experiences in B. So we only need to update with the last sample. So, $Q(A, right) = 0.5 * 1 + 0.5 * (1 + 1) = 1.5$.

- (b) [6 pts] We are now given a policy π and would like to determine how good it is using Temporal Difference Learning with $\alpha = 0.5$ and $\gamma = 1$. We run it in the environment and observe the following transitions. After observing each transition, we update the value function, which is initially 0. Fill in the blanks with the corresponding values of the Value function after these updates.

Episode Number	State	Action	Reward	Next State
1	A	right	2	B
2	B	right	12	C
3	B	right	-8	E
4	C	down	-6	F
5	F	stay	12	F
6	C	down	-6	F

State	$V^\pi(state)$
A	1
B	-1
C	-1.5
D	0
E	0
F	6

Q5. [18 pts] Game Trees

(a) [4 pts] Please mark if the following statements are true or false, and provide a brief justification or counterexample. Note a function, $f(x)$, is monotonically increasing if $f(a) > f(b)$ for all $a > b$.

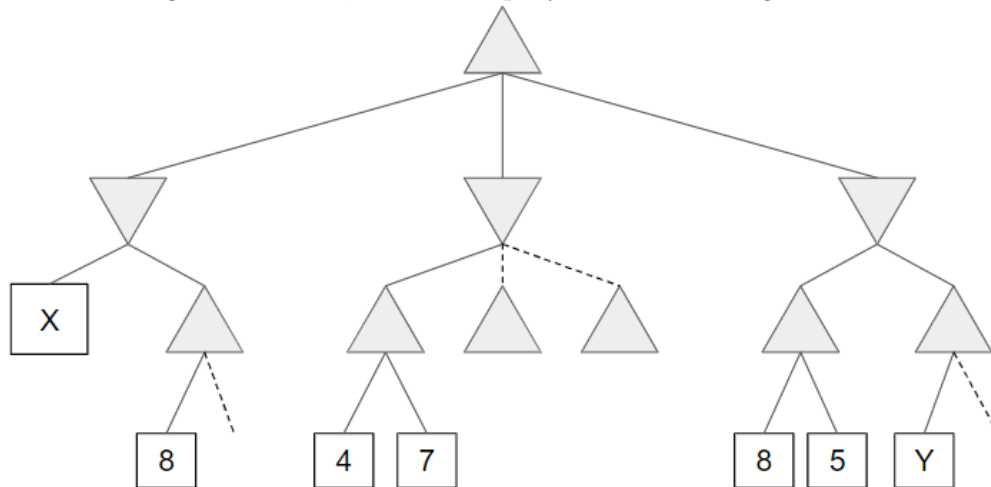
(i) [2 pts] True False Alpha-beta pruning will prune the same branches after a monotonically increasing function is applied to the leaves.

Alpha-beta pruning (and minimax in general) relies only on comparing the relative order of values which is preserved by monotonic transformation.

(ii) [2 pts] True False Expectimax ordering will choose the same actions after a monotonically increasing function is applied to the leaves.

$E(3,5) = E(4,4)$ however if squared $E(9, 25) > E(16,16)$. Only invariant under affine transformations.

(b) [6 pts] Consider the following minimax tree, where the top layer is a maximizing node



Please determine the range of X and Y such that the dashed lines are pruned in Left-to-Right alpha beta pruning. The solid lines should not be pruned for the values in your range. You might find $-\infty$ and ∞ helpful. Note, pruning occurs from left to right and latter nodes are pruned in the case of ties.

(i) [3 pts] The range of X will be

7

≤

X

<

8

(ii) [3 pts] The range of Y will be

8

≤

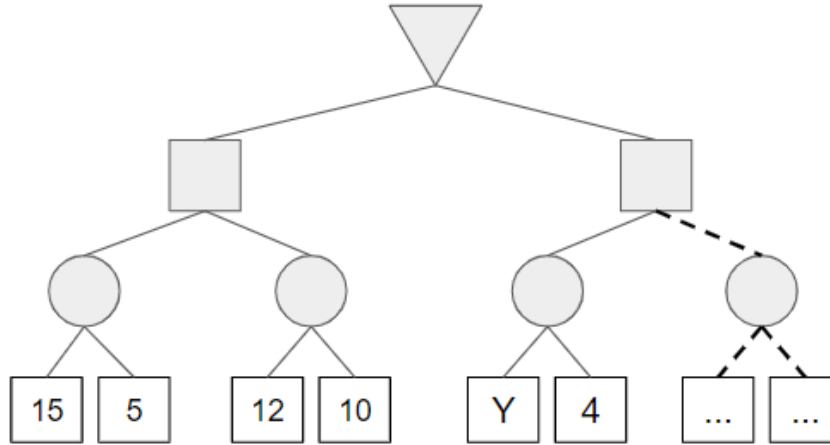
Y

<

∞

First pruned branch necessitates that $X \leq 8$. Second and third pruned branches necessitate that $7 \leq X$. Last pruned branch necessitates that $8 \leq Y$ and $X < 8$

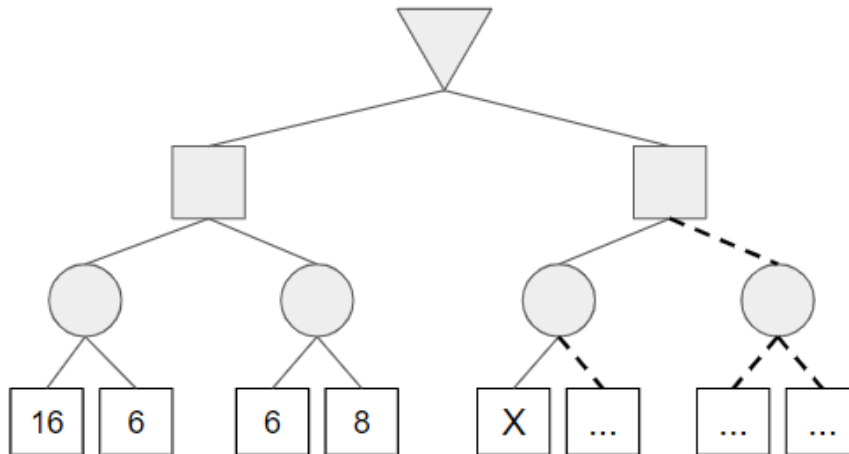
- (c) [8 pts] You are a casino owner (triangle) who is frequented by risky gamblers (squares). Your goal is to minimize the expected reward of the gamblers. Given the choice between lotteries, risky gamblers prefer the lottery (circle) that has the highest possible payout. Each reward under a lottery is equally likely. In the case of a tie between lotteries in maximum reward, the risky gambler will pick the lottery with the greater expected value. All rewards are **nonnegative integers**. Note, pruning occurs from left to right and latter nodes are pruned in the case of ties.



- (i) [4 pts] For the scenario depicted in the above game tree, write the smallest possible Y such that the dotted branches are pruned (and no solid branches are pruned). If you believe no such value of Y exists, write “N/A”.

$Y =$

After examining the left half we can realize that the risky gambler is going to choose the left lottery (15,5) because it has a higher maximum reward. This node has an expected value of 10. So once we realize the gambler will choose a lottery that has an expected value of 10 or greater for the right half we can start pruning. We reach this condition after checking $Y = 19$ and 4. Either the gambler chooses (19,4) and the expected value is 11.5 or he chooses the right lottery. The right lottery can only be chosen if its maximum reward is 19 or greater (if it is 19, its expected value would also have to be greater), in which case the expected value would be greater than 10 anyway. Note, we had to check the 4 because if that node was a zero the expected value could be less than 10. Also, note if Y was 18 or less we couldn't prune because then the right branch could be chosen with an expected value lower than 10 (if it was (19,0)).



- (ii) [4 pts] For the scenario depicted in the above game tree, write the smallest possible X such that the all dotted branches are pruned (and no solid branches are pruned). If you believe no such value of X exists,

write "N/A".

X =

22

After examining the left half we can realize that the risky gambler is going to choose the left lottery (16,6) which has an expected value of 11. So once we realize the gambler will choose a lottery that has an expected value of 11 or greater for the right half we can start pruning. We reach this condition after checking $X = 22$. No matter what the other node in that lottery is, the expected value of that lottery will be 11 or higher since all nodes are nonnegative. We also don't need to check the rightmost lottery since the only way it will be chosen is if its maximum reward is 22 or greater which would cause its expected value to be 11 or greater.